

How do monomorphic bacteria evolve? The *Mycobacterium tuberculosis* complex and the awkward population genetics of extreme clonality

CHRISTOPH STRITT, SEBASTIEN GAGNEUX

Swiss Tropical and Public Health Institute, Allschwil, Switzerland; University of Basel, Basel, Switzerland

December 14, 2022

Abstract

Exchange of genetic material through sexual reproduction or horizontal gene transfer is ubiquitous in nature. Among the few outliers that rarely recombine and mainly evolve by de novo mutation are a group of deadly bacterial pathogens, including the causative agents of leprosy, plague, typhoid, and tuberculosis. The interplay of evolutionary processes is poorly understood in these organisms. Population genetic methods allowing to infer mutation, recombination, genetic drift, and natural selection make strong assumptions that are difficult to reconcile with clonal reproduction and fully linked genomes consisting mainly of coding regions. In this review, we highlight the challenges of extreme clonality by discussing population genetic inference with the Mycobacterium tuberculosis complex, a group of closely related obligate bacterial pathogens of mammals. We show how uncertainties underlying quantitative models and verbal arguments affect previous conclusions about the way these organisms evolve. A question mark remains behind various quantities of applied and theoretical interest, including mutation rates, the interpretation of nonsynonymous polymorphisms, or the role of genetic bottlenecks. Looking ahead, we discuss how new tools for evolutionary simulations, going beyond the traditional Wright-Fisher framework, promise a more rigorous treatment of basic evolutionary processes in clonal bacteria.

INTRODUCTION

Mutation, recombination, genetic drift, and natural selection are the basic evolutionary processes that drive the evolution of life. It is the aim and "great obsession" of population genetics to infer these processes from patterns of genetic variation observed in nature (Gillespie, 2004). Since the Modern Synthesis of evolutionary biology in the 1930s, a variety of mathematical models have been developed for this purpose, which today are in wide use in the analysis of genome sequencing data (Templeton, 2021).

A problem in the application of population genetic models to empirical data is that modeling assumptions can be a far cry from the biology and life history of real organisms. Archaea and bacteria reproduce clonally through binary fission, frequently undergo horizontal gene transfer (HGT), and have genomes consisting mainly of coding regions. These characteristics are difficult to reconcile with models that are tailored to animals and plants (Woese and Goldenfeld, 2009) and commonly assume random mating, linkage equilibrium, and neutrality (Maynard Smith, 1995; Rocha, 2018). As a consequence, outside the laboratory, studies of bacterial population genetics have either remained descriptive, with much effort going into understanding the extent and effects of HGT (e.g. Denamur et al., 2021); or have resorted to models whose applicability remains an open question (discussed by Johri et al., 2022).

39 While the opportunistic, hardly predictable process of HGT has been highlighted as the most
40 problematic breach of assumptions (Maynard Smith, 1995), a different, less frequently discussed
41 challenge arises from the opposite extreme of the recombination spectrum: strictly clonal evolution,
42 or the absence of any gene flow. HGT is not a general characteristic of bacteria (Hanage, 2016).
43 Some bacteria are "monomorphic", that is, characterized by low levels of sequence diversity and an
44 apparent absence of genetic exchange (Achtman, 2008). The causative agents of several devastating
45 bacterial diseases of humans and animals belong to this group, including *Bacillus anthracis* (anthrax),
46 *Salmonella enterica* serotype typhi (typhoid), *Yersinia pestis* (plague), *Mycobacterium leprae* (leprosy),
47 and the members of the *Mycobacterium tuberculosis* complex (tuberculosis). Our understanding
48 of the evolution of these bacteria is hampered not only by the low information content in their
49 genomes, but also because there is little theoretical and conceptual work on population genetic
50 inference under extreme clonality.

51 It has been suggested that phylogenies are all that is needed to study non-recombining bacteria,
52 bacterial population genetics thus becoming "a branch of cladistics" (Maynard Smith, 1995). In
53 the absence of recombination, genetic linkage is complete and a genome is behaving as a single
54 non-recombining locus. This makes for neat phylogenies, since every part of the genome has the
55 same genealogical history. But it also complicates the inference of the processes underlying the
56 observed tree. Under strict clonality, the fate of a mutation arising in any of the few thousand
57 genes present in a typical bacterial genome is tied to all other sites in the genome. Selection
58 acting on this mutation affects the fixation probability of linked variation and interferes with
59 selection at other sites (Charlesworth, 2012; Neher, 2013). The dynamics and outcome of such
60 linked selection depend on a parameter that is usually unknown: the distribution of fitness effects
61 of new mutations (Eyre-Walker and Keightley, 2007). Periodic selection, for instance, results
62 when beneficial mutations are rare enough such that selective sweeps are well separated in time.
63 More complex dynamics emerge when beneficial mutations are frequent and co-occur in the same
64 population or on the same chromosome (Sniegowski and Gerrish, 2010).

65 Linked selection is rarely mentioned or investigated in the context of extremely clonal bacteria,
66 as already observed ten years ago (Charlesworth, 2012). This is an important omission, as it is
67 not all that clear how one would go about inferring evolutionary processes from fully linked
68 genomes. What biases are introduced when linkage equilibrium and neutrality are assumed when
69 analyzing clonal genomes? Can we meaningfully talk of "populations" when each bacterial cell is
70 a genetically isolated island?

71 *Mycobacterium tuberculosis* as a model for clonal evolution

72 In this review, we highlight the obligate pathogens of the *Mycobacterium tuberculosis* complex
73 (MTBC) as a model to study evolution under extreme clonality. The MTBC comprises a group of
74 closely related obligate pathogens that cause tuberculosis (TB) in humans and a range of wild and
75 domestic animals (Figure 1). Human TB mainly affects the global poor and has killed more than
76 1.6 million people in 2021 (World Health Organization, 2022), while it also affected large parts of
77 society in Europe and Northern America up to the early 20th century (Dubos and Dubos, 1952).
78 Today, the evolution of antibiotic resistance is a main challenge and focus of research in TB. The
79 genomes of thousands of MTBC strains from around the world have been sequenced, mainly to
80 study epidemiological dynamics and drug resistance evolution, but also to infer the origin and
81 biogeographic history of the species (Gagneux, 2018).

82 Members of the MTBC are among the more diverse of the predominantly clonal bacteria

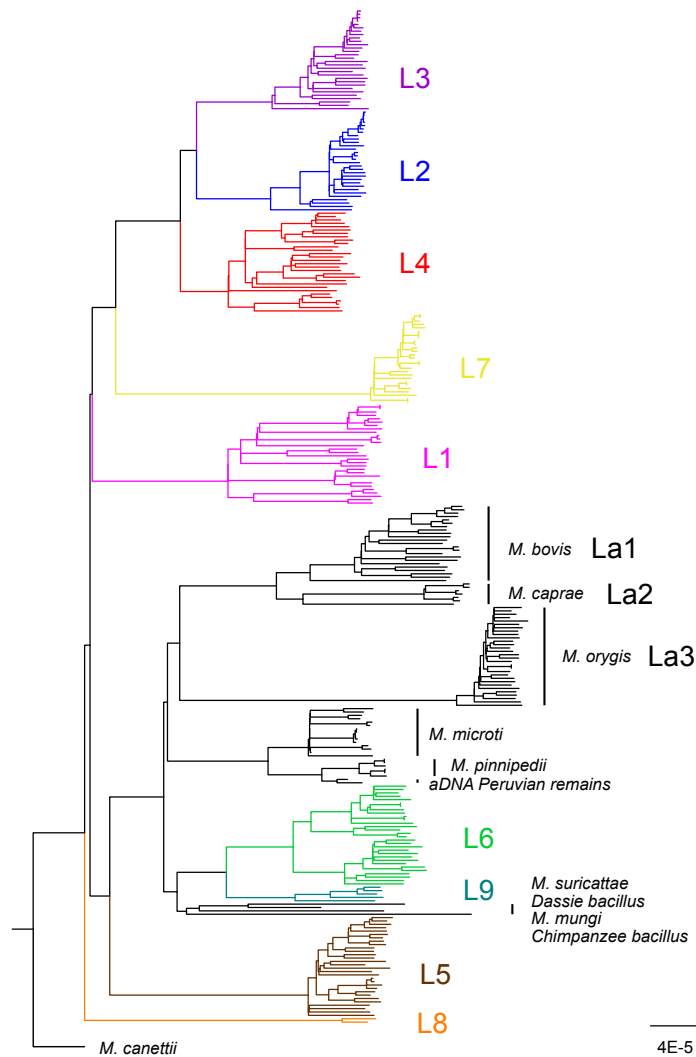


Figure 1: Rooted maximum likelihood phylogeny of the MTBC estimated from genome-wide SNPs (tree adapted from Zwyer et al. 2021). *M. canettii* is the outgroup, human-adapted lineages (L1 to L9) are shown in colors, animal-adapted lineages in black. Species names represent the historically grown nomenclature, lineage names are a more recent classification based on genomic data. Lineages 1 to 4 and 7 are also referred to as *M. tuberculosis sensu stricto*, lineages 5 and 6 as *M. africanum*. Bootstrap supports for the lineages are above 0.95 and are not displayed in the figure.

83 (Achtman, 2012), even though individual strains differ only by a maximum of ca. 2,400 SNPs
 84 across the 4.4 Mb genome (Figure 2a). At the molecular level, the MTBC is further characterized
 85 by a high GC content, a high proportion of nonsynonymous polymorphisms, and a low proportion
 86 of homoplastic mutations (Figure 2b-d). It seems that the low diversity of the MTBC has deterred
 87 evolutionary biologists from engaging with this bacterium. Many studies content with speculative
 88 invocations of genetic drift and natural selection, typically referring to the triad Hershberg et al.
 89 (2008), Namouchi et al. (2012) and Pepperell et al. (2013), who are among the few studies in the

90 large MTBC literature that have put basic evolutionary processes into focus.

91 In this review, we present the main hypotheses about what drives the evolution of the MTBC,
92 and how they have been arrived at. Particular attention is paid to models, their assumptions, and
93 the traits of the MTBC that might conflict with the latter. Evolutionary simulations are discussed as
94 a way to achieve a more quantitative treatment of frequently invoked processes such as purifying
95 selection or periodic bottlenecks.

96 RECOMBINATION

97 How "strict" is clonality in the MTBC? In the past, bacteria were classified as "clonal" or "monomor-
98 phic" based on a handful of housekeeping genes (Maynard Smith et al., 1993; Selander et al., 1987).
99 With the full resolution of whole genome sequences, this classification needs to be reassessed.
100 As discussed in the following, experimental and observational evidence agree that the MTBC is
101 predominantly clonal, and that few to no new genes have found their way into the MTBC since
102 the most recent common ancestor of the currently existing lineages. In contrast to interstrain
103 recombination, intrachromosomal recombination is common and increasingly recognized as an
104 important source of genetic variation.

105 Experimental evidence: genetic factors versus lack of opportunity

106 Most of the knowledge about the molecular mechanisms of HGT in mycobacteria stems from
107 research with *Mycobacterium smegmatis*, a fast-growing, non-pathogenic mycobacterium more easily
108 amenable to cultivation and genetic engineering than the bacteria of the MTBC. Mycobacteria lack
109 the traditional components of HGT, possibly because transfer through the complex cell envelopes
110 of these diderm bacteria requires other mechanisms (Madacki et al., 2021). Investigations of gene
111 transfer in *M. smegmatis* have led to the description of a previously unknown form of bacterial
112 conjugation: distributive conjugal transfer (DCT, reviewed by Gray and Derbyshire, 2018). DCT
113 involves the transfer of chromosomal DNA and gives rise to mosaic genomes, with hundreds
114 of pieces of DNA of variable sizes dispersed in the receiver genome. DCT thus challenges the
115 paradigm of "localized sex" in bacteria (Smith et al., 1991) and might explain the recombinogenic
116 population structure of many mycobacteria (Panda et al., 2018).

117 Of particular interest regarding the evolution of the MTBC is the observation of DCT in the
118 closely related *Mycobacterium canettii*. *M. canettii* shares an average nucleotide identity of 97.5%
119 with the MTBC yet is strikingly more diverse: a handful of *M. canettii* strains from eastern Africa
120 harbor more genetic diversity than the whole MTBC (Supply et al., 2013). Mating assays have
121 shown that DCT occurs in *M. canettii*, while no DCT was observed between three MTBC strains
122 (Boritsch et al., 2016). The same assays combining *M. canettii* and MTBC strains revealed that the
123 latter can act as donors but not as receivers of DNA during DCT, as pieces of MTBC DNA were
124 integrated into *M. canettii* genomes but not vice versa (Madacki et al., 2021). In *M. smegmatis*,
125 polymorphisms in the *esxI* secretion locus underlay self identity and conjugal compatibility
126 (Clark et al., 2022). In *M. canettii* and the MTBC, the molecular mechanisms underlying conjugal
127 compatibility do not depend on *esxI* and remain to be elucidated (Madacki et al., 2021).

128 Lack of opportunity has been proposed to explain why intracellular pathogens such as the
129 MTBC do not seem to recombine (Casadevall, 2008; Chiner-Oms et al., 2019b). Alternatively,
130 avoidance of HGT could be an evolutionary strategy with a genetic basis, an adaptation to
131 parasitism (Tibayrenc and Ayala, 2017). Against the first scenario, it can be argued that there

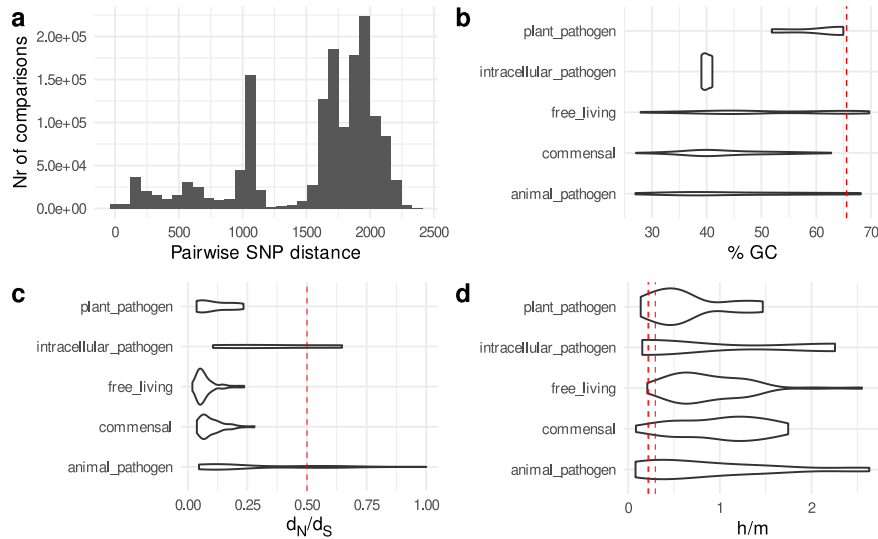


Figure 2: Genetic diversity and molecular characteristics of the MTBC. *a)* Pairwise genetic differences between the strains shown in Figure 1, based on single nucleotide polymorphisms. *b) to d)* show molecular characteristics of the MTBC compared to 150 other bacterial species with diverse lifestyles (data from Bobay and Ochman, 2018). Red lines show the values for the bacteria of the MTBC (*M. tuberculosis sensu stricto*, *M. bovis*, and *M. africanum*) along the distributions. *b)* GC content, *c)* d_N/d_S , the genome-wide ratio of nonsynonymous to synonymous polymorphisms, *d)* the ratio of homoplasic to non-homoplasic mutations, a proxy for recombination.

132 is indeed more opportunity to recombine than the label "intracellular pathogen" might suggest.
 133 The bacteria of the MTBC are not confined to intracellular environments, but are also present
 134 in large extracellular populations after the induction of necrosis (Orme, 2014). Furthermore,
 135 mixed infections do occur (Moreno-Molina et al., 2021; Tarashi et al., 2017), such that diverged
 136 strains might have the opportunity to recombine. Further investigation into the genetic and
 137 environmental determinants of extreme clonality would be worthwhile, and the *M. canettii*-MTBC
 138 system provides a great opportunity to elucidate the poorly understood evolutionary transition to
 139 extreme clonality characteristic of many obligate pathogens.

140 Recombination between closely related strains: how strict is clonality?

141 Genome sequences from diverse MTBC strains are an important complement to experimental data,
 142 which leave open the question how far the observed outcome depends on the specific conditions
 143 and strains used in the laboratory. Various studies have investigated the extent of HGT in natural
 144 strains of the MTBC, motivated by the observation how HGT accelerates resistance evolution
 145 in other bacterial pathogens (Davies and Davis, 2010). Some have suggested that interstrain
 146 recombination does occur. Liu et al. (2006), using datasets of 36 synonymous SNPs in 3,320 strains
 147 and 407 SNPs in 37 strains, found that mutation alone cannot explain the observed haplotype
 148 diversity, and identified a mosaic region in front of a *PPE* gene suggesting a recombination hotspot.
 149 They also point out the possibility that the pattern may have arisen through recombination
 150 between homologous sequences in the same genome. Namouchi et al. (2012) investigated 24

151 sequenced MTBC genomes and reported that "four different approaches showed evident signs of
152 recombination in *M. tuberculosis*", with recombination typically involving small tracts of around
153 50 bp. On the other hand, the most extensive investigation to date, using different methods on
154 genome-wide SNPs in 1,591 diverse strains, found "no measurable ongoing recombination among
155 the MTBC strains" (Chiner-Oms et al., 2019b).

156 Generalizing from these studies is difficult due to the diversity of datasets and methods used.
157 It has been suggested that the signs of recombination described by Namouchi et al. are mainly
158 artefacts as they are overrepresented in regions difficult to align or assemble, in particular repetitive
159 and low-complexity regions in insertion sequences and the expanded *PE/PPE* gene families
160 (Godfroid et al., 2018). Alternatively, signs of recombination can arise from gene conversion during
161 intrachromosomal recombination, to which these repetitive sequences are prone (Liu et al., 2006).
162 Gene conversion is the non-reciprocal transfer of DNA from one homologous sequence to another,
163 which in the MTBC might account for recombination signatures in *ESX*, *PE*, *PPE*, *PE/PGRS*
164 gene families (Karboul et al., 2008; Phelan et al., 2016; Uplekar et al., 2011). Intrachromosomal
165 recombination can also have more dramatic outcomes. More and more structural variants are
166 described in MTBC genomes, ranging from insertion sequence (McEvoy et al., 2007) and gene copy
167 number polymorphisms (Fishbein et al., 2015) to massive inversions (Merrikh and Merrikh, 2018)
168 and tandem duplications (Wang et al., 2022). This is a vast topic deserving a dedicated review.
169 It is brought up here to emphasize that recombination is an umbrella term for diverse processes
170 of inter- and intrachromosomal exchange; and that clonality does therefore not imply absence
171 of recombination, strictly speaking, but only of HGT. In the near future, long-read sequencing
172 should allow more extensive studies of the repetitive "dark matter" in the MTBC genome and how
173 it generates genetic variation intrachromosomally.

174 A basic limitation of methods to infer recombination is that they cannot distinguish *de novo*
175 mutations from allelic recombination between closely related individuals, which might involve the
176 exchange of a single nucleotide (Martin et al., 2011). Allelic recombination does not introduce new
177 genes, but it can affect the nucleotide landscape through recombination-associated processes like
178 biased gene conversion (Duret and Galtier, 2009) or increased mutation rates around strand breaks
179 (Fitzgerald and Rosenberg, 2019). While HGT between close relatives would be less restricted by
180 opportunity, genetic incompatibilities might prevent gene transfer between close relatives, as in *M.*
181 *smegmatis* (Clark et al., 2022).

182 MUTATION

183 While in some bacteria new variants are more likely to be generated by HGT than by mutation
184 (Vos and Didelot, 2009), under extreme clonality *de novo* mutations are the main source of genetic
185 diversity and adaptation. The speed and direction in which a clonal prokaryote evolves is thus
186 determined by the rate and spectrum of new mutations and by their effect on fitness. Numerous
187 studies have investigated mutagenesis in the MTBC (reviewed by Mcgrath et al., 2014). As
188 discussed below, in addition to methodological issues in estimating mutation rates, the life history
189 of the bacteria, which can include extended periods of dormancy, poses a main challenge in
190 understanding the rate at which variation originates *in vivo*.

191 In the MTBC literature, as elsewhere, the mutation rate is sometimes confounded with the
192 molecular clock rate. While the former refers to the rate at which mutations appear in the genome,
193 the latter stands for an allegedly constant rate at which mutations accumulate through time (Ho et
194 al., 2011). Both rates are subsumed in the more general concept of evolutionary rates. As discussed

195 below, the power law that describes the slowing of evolutionary rates as one considers longer
196 timescales is not as clear in the MTBC as in other bacteria: *in vitro* mutation rate estimates can
197 be similar to clock rate estimates from datasets including ancient DNA. How far methodological
198 biases or evolutionary processes underly this surprising finding remains to be understood.

199 Plasticity of mutation rates and generation times

200 In the model mycobacterium *M. smegmatis*, a mutation rate of 5.27×10^{-10} mutations per site per
201 generation was inferred in a mutation accumulation experiment (Kucukyildirim et al., 2016). For
202 the MTBC itself, no such experiment has been conducted yet. Fluctuation assays suggest that point
203 mutations in the MTBC appear at a rate of about 2.1×10^{-10} mutations per site per generation
204 and at a similar rate during active disease in macaques if a generation time of 20 h is assumed
205 (Ford et al., 2011, Figure 3). A later study, using the same fluctuation assay, found *in vitro* rates
206 of 6.01×10^{-10} in a lineage 4 and 2.16×10^{-9} in a lineage 2 strain, suggesting somewhat faster
207 and variable mutation rates within the MTBC (Ford et al., 2013). Comparatively fast rates were
208 also proposed in two additional experimental evolution studies. After serial passaging of a MTBC
209 strain through macrophage-like THP1 cells for 80 generations, Guerrini et al. (2016) inferred
210 a rate of 5.7×10^{-9} per bp per generation. Copin et al. (2016), passaging bacteria in mice and
211 assuming a generation time of 20 h, estimated a mutation rate of 3.8×10^{-9} in wild type mice
212 and of 7.7×10^{-10} in T cell-deficient mice, suggesting that the presence of T cells leads to elevated
213 mutation rates.

214 Overall, per-generation mutation rates estimated for the MTBC are well within the range of
215 those in other bacteria, which typically are in the order 10^{-10} (reviewed by Katju and Bergthorsson,
216 2019). When trying to scale mutation rates to calendar time, however, complications due to the
217 complex life history of these bacteria become apparent. The bacteria of the MTBC have long
218 generation times ranging from 18 h in nutrient rich medium to potentially much longer time-spans
219 *in vivo* (Colangeli et al., 2020). Assuming a generation time of 24 hours, a mutation rate of
220 2.1×10^{-10} translates to 7.7×10^{-8} mutations per site per year, or about 0.34 per genome per year,
221 which is indeed low compared to other bacteria (Duchêne et al., 2016; Lynch, 2010).

222 In contrast to pathogens employing a "hit and run" strategy, bacteria of the MTBC can enter a
223 state of reduced activity and persist for years in latent infections (Dutta and Karakousis, 2014). It is
224 unclear whether latency and longer generation times imply a reduced mutation rate, as expected
225 if mutation is driven by replication, or not, as expected if environmental stress drives mutation
226 (Weller and Wu, 2015). Ford et al. (2011), in their experimental infection of macaques, found
227 similar rates in latent and active disease (Figure 3), supporting stress-induced mutagenesis. A
228 more complex, two-phased scenario was suggested by Colangeli et al. (2020), who investigated 24
229 paired TB cases with latently infected household contacts: mutation rates remained high up to
230 two years, but then decrease with longer latency as the bacteria enter a quiescent state with longer
231 generation times (Figure 3).

232 In summary, mutation rates estimated for the MTBC should be interpreted with some caution.
233 Generation times are only known with confidence *in vitro*. At the same time, fluctuation assays
234 reflect the mutation rate of a single gene (*rpoB*, the main drug resistance target of rifampicin)
235 that might not be representative for the whole genome (Katju and Bergthorsson, 2019); and in
236 the absence of stress, which *in vivo* might alter both the rate and the spectrum of new mutations
237 (Fitzgerald and Rosenberg, 2019).

238 Why are MTBC genomes so GC-rich?

239 In bacteria, newly arising mutations are biased towards adenines and thymines (Hershberg and
240 Petrov, 2010; Hildebrand et al., 2010), which in the MTBC might reflect stress-induced mutagenesis
241 in an intracellular environment rich in reactive oxygen and nitrogen species (Chiner-Oms et al.,
242 2019a; Liu et al., 2020). If mutation bias and genetic drift alone would determine the nucleotide
243 landscape (mutation-drift equilibrium), the expected GC content in the MTBC would be 41.5%
244 (Hershberg and Petrov, 2010). MTBC genomes, however, consist to 65.6% of guanines and cytosines
245 (Figure 2b; Cole et al., 1998), with values of 80% at synonymous and 60% at nonsynonymous
246 sites. Such a discrepancy between observed and expected GC content is observed in many
247 prokaryotes, whose genomes vary hugely in GC content (Figure 2b). It implies that an unknown
248 process, unaccounted for in standard models of molecular evolution, affects the segregation of
249 polymorphisms through time (Rocha and Feil, 2010).

250 Several large-scale comparative studies have attempted to find a general explanation for
251 the discordance between expected and observed GC content in prokaryotes. One prominent
252 hypothesis is that nucleotide composition reflects adaptation to environmental conditions, for
253 example through selection for thermal stability of DNA (e.g. Reichenberger et al., 2015). An
254 intriguing twist to this idea was recently added by Weissman et al. (2019), who described a
255 correlation between GC content, environmental variables, and the presence of *Ku*, the key gene in
256 the non-homologous end-joining (NHEJ) pathway for DNA break repair. The authors propose
257 that high GC content could be beneficial in bacteria suffering stress-induced double strand breaks
258 in periods of slow or no growth, when NHEJ is required for repair because only a single copy of
259 the genome is present. This is an interesting scenario for the MTBC, where long periods of latency
260 can occur (see above) and the *Ku* gene is highly conserved.

261 An alternative explanation for GC bias that does not imply a selective advantage is GC-biased
262 gene conversion (gBGC). This process occurs during homologous recombination when mismatches
263 in heteroduplex DNA are preferentially resolved into guanines and cytosines (reviewed by Duret
264 and Galtier, 2009). The gBGC hypothesis predicts that GC content is higher in regions with high
265 recombination rates, which is observed in mammalian genomes. In bacteria, the role of gBGC is
266 contested. Whether comparative studies find associations between GC content and recombination
267 depends on the method used to infer recombination, and exceptions to general trends are common
268 (Bobay and Ochman, 2017; Lassalle et al., 2015).

269 With its numerous genome sequences that can be placed in a robust phylogenetic framework,
270 the MTBC provides an opportunity to study the evolution of base composition in detail and thus
271 to complement comparative studies. A hypothesis to test is that the MTBC is evolving from the
272 generally GC-rich state of mycobacteria (58 to 70%, *Mycobacterium* sp. genomes on NCBI) to a
273 more AT-rich state characteristic of obligate pathogens (Rocha and Danchin, 2002, Figure 2b),
274 including *Mycobacterium leprae* (58%).

275 The time (in)dependence of evolutionary rates in the MTBC

276 Molecular dating has led to a re-evaluation of the origin and history of the MTBC, as for many
277 other organisms. Earlier studies, assuming a synonymous mutation rate or a co-diversification
278 of humans and the MTBC, located the most recent common ancestor of the existing lineages in
279 Africa and suggested a scenario according to which humans and the MTBC have co-diversified
280 across the globe (Comas et al., 2013; Hughes et al., 2002; Kapur et al., 1994). Recent estimates,

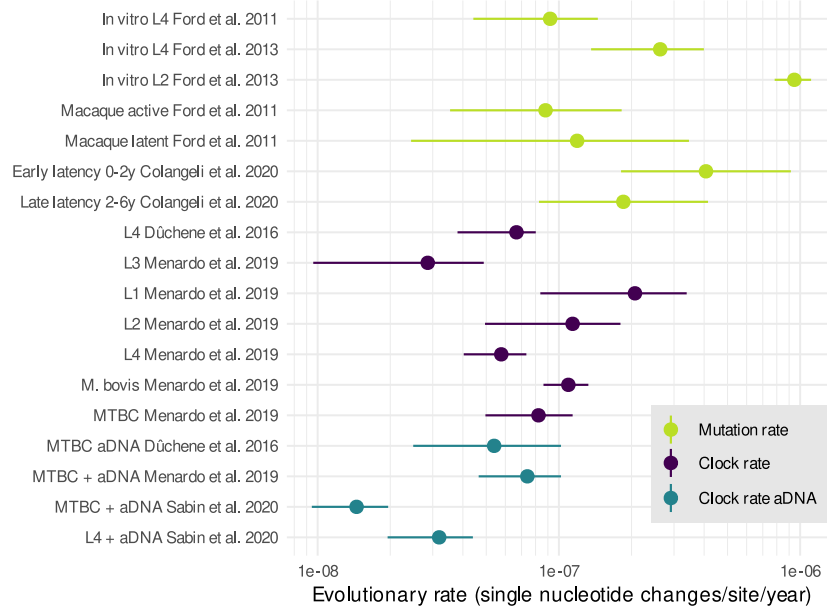


Figure 3: Evolutionary rates in the MTBC. Only studies that report confidence intervals were considered. For the fluctuation assay estimates in Ford et al. (2011, 2013), a generation time $g = 20$ was assumed to translate rates to calendar time. The rates of Colangeli et al. (2020) were translated back to calendar time by assuming $g = 18h$, as reported by the authors. From the molecular clock study of Menardo et al. (2019), BEAST estimates are reported for a $1/x$ clock rate prior and constant population size. For the BEAST analysis of Sabin et al. (2020), results for the birth-death skyline model with an uncorrelated lognormal clock are reported.

281 making use of tip dating, ancient DNA (aDNA) samples, and Bayesian phylogenetics, propose a
 282 more recent common ancestor in the Neolithic ca. 6,000 years ago (Bos et al., 2014; Kay et al., 2015;
 283 Sabin et al., 2020).

284 One caveat regarding these estimates is the poorly understood variability of evolutionary
 285 rates in the MTBC through time. For mitochondrial DNA, viruses, and bacteria, evolutionary
 286 rates usually appear faster when estimated from recent polymorphisms (Ho et al., 2011). For
 287 bacteria, Duchêne et al. (2016) found a clear negative association, described by an exponential
 288 decay curve, between clock rates and sampling time spans in 16 bacterial species, with an order of
 289 magnitude difference between a 10 year and a 100 year sampling period. The delayed effect of
 290 purifying selection is the most prominent explanation for this time dependence of evolutionary
 291 rates, although methodological biases might also contribute (Emerson and Hickerson, 2015; Ho
 292 et al., 2015). Time dependence can have a large effect on molecular dating: Membrane et al. (2019)
 293 showed that accounting for purifying selection by using relaxed clock or epoch models can shift
 294 divergence times one order of magnitude back in time. Could this explain the surprisingly recent
 295 time to the most recent common ancestor estimated by the aDNA studies?

296 In the study of Duchêne et al. (2016), the MTBC does not follow the general pattern of time
 297 dependence: almost identical rates were obtained from samples spanning 15 and 895 years. Similarly,
 298 Menardo et al. (2019) found only marginally lower rates when calibrating the clock with the same
 299 three samples of ancient DNA from Precolumbian human remains and an extensive MTBC dataset

300 covering a sampling period of 30 years. An overview of evolutionary rates estimated for the MTBC
301 illustrates the large variability and uncertainty of rate estimates, but also suggest an overall trend
302 of time dependence (Figure 3). As Menardo et al. (2019) showed in their extensive study of the
303 molecular clock in the MTBC, clock rates vary substantially among lineages and clades of the
304 MTBC and have large confidence intervals. Lineage 1, for instance, seems to have evolved faster
305 than other lineages, and indeed faster than the L4 strain accumulated mutations in the fluctuation
306 assay of Ford et al. (2011). On the slow end of the spectrum is the long-term clock rate estimated
307 by Sabin et al. (2020), for which all six aDNA samples available so far were included (1.4×10^{-8} ,
308 95% HPD 9.46×10^{-9} , 1.96×10^{-8}).

309 A possible methodological bias underlying *in vivo* mutation rate estimates was recently sug-
310 gested in a simulation study of within-host evolution. Morales-Arce et al., 2020 suggested that the
311 genome-wide mutation rate of the MTBC might be two orders of magnitude faster, in the order
312 10^{-8} /bp/generation, if one accounts for progeny skew (Box 2) and the removal of mutations
313 through purifying selection during within-host evolution. The authors simulated a population
314 undergoing a transmission bottleneck, followed by a recovery to a large population size and
315 within-host evolution under purifying selection and with per-generation progeny skew. Com-
316 paring the resulting patterns of diversity with the empirical within-host data of Trauner et al.,
317 2017, they found that mutation rates in the order of 1×10^{-8} to 9×10^{-8} result in similar levels of
318 variation as described by Trauner et al.

319 Box 1: Simulating bacterial populations

320 Simulations are an invaluable tool in evolutionary genetics: they allow to test intuitions and
321 methods, to compare alternative scenarios, and to fit models to data (Hoban et al., 2012; Johri
322 et al., 2022). For bacterial population genetics, the use of simulations was so far rather limited.
323 On the one hand, most simulators are based on the coalescent – the backwards-in-time
324 variant of the Wright-Fisher model. These are fast, but usually limited to neutral scenarios of
325 population size changes and migration. More flexible forward simulators, on the other hand,
326 are much slower because they track the fate of all individuals of the simulated population
327 rather than just of a sample, as in the coalescent (Hoban et al., 2012).

328 Recent advances in forward simulation, however, make it possible to simulate ever more
329 realistic scenarios through improved computational efficiency (Haller et al., 2019) and more
330 flexible non-Wright-Fisher models (Haller and Messer, 2019). The simulation framework of
331 SLiM was recently used to simulate bacteria evolving in a Petri dish in the presence of an
332 antibiotic (Cury et al., 2021). This individual-based forward simulation was spatially explicit
333 and modelled clonal reproduction through binary fission, gene conversion, density-dependent
334 selection, and positive selection for antibiotic resistance. The scriptability of SLiM allows to
335 incorporate more or less arbitrarily complex genetic architectures and life histories, although
336 computational time still sets boundaries.

337 GENETIC DRIFT AND PURIFYING SELECTION

338 Once a mutation appears in a genome, its fate depends on the selective advantage or disadvantage
339 it confers – and on chance. Genetic drift is the "chance factor" in evolution: it describes the
340 undirected, stochastic change of allele frequencies due to sampling effects (Plutynski, 2007). The

341 biological relevance of genetic drift is that it sets limits to natural selection (Kimura, 1983; Lynch,
342 2007). The efficacy of natural selection is inversely related to the strength of drift: when genetic
343 drift is strong, changes in the frequencies of alleles depend less on their effect on fitness, such
344 that, by chance, deleterious alleles can increase and beneficial ones decrease in frequency (Kimura,
345 1983; Ohta, 1992).

346 Genetic drift is frequently invoked as an ad hoc explanation, but actually inferring and
347 quantifying it is difficult. In the standard Wright-Fisher (WF) model with panmixia, discrete
348 generations, and no selection, drift occurs when the alleles to form the next generation are
349 randomly sampled from the parental population (Fisher, 1930; Wright, 1931). In this idealized
350 lottery-like scenario, the strength of drift simply depends on the size of the sample, with less drift
351 in larger samples according to the law of large numbers. Natural populations deviate from the WF
352 model in numerous ways, yet population size remains a useful measure for drift when it is rescaled
353 to account for these deviations (Charlesworth, 2009). The resulting effective population size N_e
354 can be interpreted as the size of an idealized WF population that experiences the same amount of
355 drift as the real population in question (e.g. Gillespie, 2004). In bacteria, population subdivision,
356 linked selection, and demographic changes all imply that sampling effects are stronger than under
357 panmixia (Price and Arkin, 2015), and that effective population sizes are orders of magnitude
358 smaller than census sizes (Bobay and Ochman, 2018).

359 As discussed in this section, arguments about the strength of drift in the MTBC are largely
360 based on indirect evidence in the form of low diversity and overabundant nonsynonymous
361 polymorphisms. Estimates of N_e are sometimes obtained in Bayesian skyline analyses, but their
362 underlying assumptions are problematic. Finally, we discuss transmission bottlenecks in the
363 MTBC, a main mechanism of stochastic sampling whose mid- and long-term consequences go
364 beyond simple reductions in genetic diversity and remain to be understood.

365 Do overabundant nonsynonymous polymorphisms indicate strong genetic drift?

366 In the MTBC, the drift-versus-selection discussion has mainly revolved around the large proportion
367 of nonsynonymous polymorphisms observed in the species. The MTBC has a genome-wide ratio
368 of nonsynonymous to synonymous polymorphisms (d_N/d_S) of around 0.5 when diverse strains
369 from across the phylogeny are considered (Figure 2c). This is one third higher than in the closely
370 related *M. canettii* (Supply et al., 2013) and more than six times higher than the median (0.076) of
371 the 153 diverse species studied by (Bobay and Ochman, 2018).

372 Hershberg et al. (2008) have interpreted the high d_N/d_S in the MTBC as evidence for "extremely
373 reduced purifying selection" – in other words strong genetic drift – which would allow the
374 accumulation of deleterious nonsynonymous mutations. The authors refute the alternative
375 explanation that nonsynonymous changes are due to positive selection by pointing out that d_N/d_S
376 does not differ between housekeeping, surface-exposed, and virulence genes, as might be expected
377 if host immunity would drive adaptive diversification. This interpretation of d_N/d_S fits well
378 with the generalization that the intracellular niche of pathogens and symbionts implies smaller
379 population sizes and stronger drift. Kuo et al. (2009) inferred strong drift in human pathogens
380 including the MTBC and reported a strong inverse relationship between drift and genome size.
381 A similar conclusion is reached by Balbi et al. (2009), who compared *E. coli* with the closely
382 related pathogenic *Shigella* and found signs of increased drift in the latter, including an excess of
383 nonsynonymous mutations and of transversions, which are proportionally more nonsynonymous
384 and thus deleterious than transitions.

385 Different studies have challenged the view that purifying selection is "extremely reduced"
386 in the MTBC. Bringing in a temporal perspective on d_N/d_S , Namouchi et al. (2012) found 25%
387 more nonsynonymous SNPs on terminal branches in their tree of 22 globally diverse strains. This
388 suggests that deleterious nonsynonymous mutations are purged through selection over time, such
389 that they become scarce in deeper parts of the phylogeny (Rocha et al., 2006). In general, SNPs
390 are strongly skewed towards rare alleles in the MTBC, be it at the global or the within-host level
391 (O'Neill et al., 2015; Trauner et al., 2017). SNPs are thus not only few in the MTBC, but also to a
392 large proportion singletons (Chiner-Oms et al., 2019b), that is, present in one single strain. While
393 this is consistent with purifying selection preventing variants to rise in frequency, other processes
394 can cause the same pattern, in particular the dynamics of clonal growth. Furthermore, it remains
395 to be understood what biases are introduced by the punctual sampling of highly structured and
396 dynamic within-host populations (Morales-Arce et al., 2021).

397 In the so far only attempt to quantify the strength of purifying selection across the genome,
398 Pepperell et al. (2013) fitted a model including demographic expansion and a fraction of sites
399 under selection to the site frequency spectrum obtained from a global sample of the MTBC. They
400 infer purifying selection at nonsynonymous sites across 95% of the genome, with a selection
401 coefficient s of -9.5×10^{-4} . This value is interpreted as "strong" compared to values in humans
402 and *Drosophila*. The authors used simulations of completely linked genomes to evaluate their
403 models, which assume linkage equilibrium between sites. They find that their best model performs
404 poorly in some scenarios; specifically, strong selection can be misinferred when complete linkage
405 is combined with weak purifying selection, which might thus confound their estimate of s . Other
406 model assumptions were not tested, for example the absence of population subdivision or that the
407 population follows a simple demographic model of exponential growth.

408 Strong genetic drift leaves other signs than an excess of nonsynonymous mutations, includ-
409 ing pseudogenization, proliferation of selfish genetic elements, or an increased proportion of
410 transversions. With strong drift and asexual reproduction, such signatures can accumulate through
411 Muller's ratchet, where lack of recombination and reduced efficacy of purifying selection lead to a
412 build-up of deleterious mutations (Muller, 1964). As pointed out by Namouchi et al. (2012), these
413 signatures are hardly evident in the MTBC. There are 30 pseudogenes in the H37Rv reference
414 genome (Cole et al., 1998), in line with the generally low number of pseudogenes in bacterial
415 genomes (Lawrence et al., 2001) and contrasting with the more than 1,000 pseudogenes described
416 in the genome of *M. leprae* (Gómez-Valero et al., 2007). Also insertion sequences do not thrive in
417 the MTBC: almost all IS activity is due to a single active element, IS6110, which is over-represented
418 in intergenic regions, occurs at low frequencies, and thus probably evolves under strong purifying
419 selection (McEvoy et al., 2007). Finally, transitions occur well in excess of transversions (Payne
420 et al., 2019). Taken together, there is scant evidence for genome erosion driven by Muller's ratchet
421 in the MTBC.

422 Drift is expected to dominate allele frequency changes when $|N_e \times s| \ll 1$ (Kimura, 1983;
423 Ohta, 1992). Thus, rather than small population sizes (N_e), reduced selection coefficients (s), as
424 they might arise when many genes are not required anymore after the transition to an intracellular
425 niche, could explain genome erosion in obligate pathogens. Applied to the MTBC, the absence
426 of genome erosion could indicate that these bacteria still require a large complement of genes,
427 which thus remain under strong purifying selection. Alternatively, the MTBC might be a young
428 pathogen in an early phase of genome degradation, where nonsynonymous mutations are only
429 starting to accumulate (Kuo et al., 2009).

430 Are synonymous sites under selection?

431 How could the high genome-wide d_N/d_S in the MTBC be explained if not by strong drift? An
432 intriguing alternative scenario is purifying selection at synonymous sites (Namouchi et al., 2012).
433 High d_N/d_S can reflect an overabundance of nonsynonymous mutations (numerator), but also
434 a lower number of synonymous mutations (denominator) than in other species. Fitness effects
435 of synonymous mutations can arise when different codons result in variation in RNA stability,
436 protein folding, and translation efficiency and accuracy (reviewed by Hershberg and Petrov, 2008).
437 Already weak selection on synonymous sites can inflate d_N/d_S , as shown in a recent study of
438 codon usage in 13 bacterial genomes (Rahman et al., 2021).

439 In the MTBC, codon frequencies are associated with gene expression (Andersson and Sharp,
440 1996; Pan et al., 1998), but also with the hydrophobicity of proteins and sequence conservation
441 (De Miranda et al., 2000). As suggested in the latter study, a combination of selective pressures may
442 thus act on synonymous sites in the MTBC, including the more efficient and accurate translation
443 of certain codons and constraints on protein folding. Wang and Chen (2013) assessed possible
444 selection on synonymous sites by comparing synonymous (d_S) to intergenic (d_I) diversity across 13
445 diverse MTBC genomes. Diversity varies strongly depending on the genomic position, suggesting
446 variation in mutation rates or selective pressures across the genome. In the majority of windows,
447 however, d_S is higher than d_I . Under the assumption that intergenic regions are free from selection
448 pressures, these results are interpreted as evidence for positive selection on synonymous sites,
449 specifically for increased translational efficiency.

450 The alternative explanation, mentioned but not favored by Wang & Chen, is that purifying
451 selection is stronger in intergenic regions than at synonymous sites. Intergenic regions in bacteria
452 are packed with regulatory motives and can hardly be assumed to evolve neutrally (Molina
453 and Van Nimwegen, 2008; Rocha, 2018). Rather than comparing synonymous against assumed
454 neutral sites, Thorpe et al. (2017) assessed the relative strength of purifying selection by comparing
455 the proportion of singleton mutations among different site categories, reflecting that a higher
456 proportion of singletons indicates stronger purifying selection. In five out of six species, site
457 categories show a clear ranking, with the proportion of singletons increasing from synonymous,
458 intergenic, non-synonymous, to non-sense mutations. In the MTBC, however, no differences
459 between categories are apparent: there are similar proportions of singletons in all four categories.
460 This surprising observation can at least partly be explained by the dataset used by the authors,
461 which includes many near-identical MTBC strains sampled in a single country. Still, that even at
462 short timescales non-sense mutations in the MTBC do not appear to be under stronger selection
463 than synonymous mutations asks for clarification in future studies.

464 In summary, synonymous sites are frequently assumed to be neutral, but studies on codon
465 frequencies and comparisons of synonymous with other sites in the genome suggest a more
466 complex picture. This is a topic deserving a focused study, applying the measures developed in
467 previous work to a large dataset covering different timescales.

468 Bayesian skyline plots and the issue of storytelling

469 Neutral sites are in short supply in prokaryotes (Rocha, 2018). In contrast to eukaryotes, the
470 streamlined genomes of archaea and bacteria do not contain large swaths of decaying repeats and
471 other DNA debris which can be assumed to be non-functional. This poses a particular challenge for
472 the estimation effective population sizes and the quantification of genetic drift, which traditionally

473 relies on the availability of sites not affected by natural selection (Charlesworth, 2009).

474 A popular approach to estimate effective population sizes and their change through time
475 are Bayesian skylines (Ho and Shapiro, 2011). These models are frequently used in Bayesian
476 phylogenetics, where N_e is treated as a nuisance parameter. Many studies, however, interpret N_e
477 literally as historical change in population size and thereby provide instructive examples of how
478 strong assumptions are ignored for the sake of storytelling.

479 Bayesian skyline models assume neutrality in order to translate coalescence times into popula-
480 tion sizes. Several studies have shown that non-neutral processes confound demographic inference
481 and should not simply be assumed away. Recombination (Hedge and Wilson, 2014), population
482 structure (Heller et al., 2013), sampling design, gene conversion, and selection (Lapierre et al.,
483 2016), as well as the skewness of reproductive success (Menardo et al., 2021a) all create spurious
484 signs of population size changes. As observed by Lapierre et al., 2016, such methodological
485 biases might explain why population size trajectories look suspiciously similar for a wide range of
486 species.

487 Despite these caveats, Bayesian skyline plots continue to be used and interpreted liberally in
488 the MTBC literature. Skyline plots were presented as evidence for a Neolithic expansion (Comas
489 et al., 2013), expansions of specific lineages (Merker et al., 2022; Mulholland et al., 2019; O'Neill
490 et al., 2019), or a recent co-expansion with humans in Tibet (Liu et al., 2021). That population
491 size trajectories "make sense" in the historical narratives of these articles does not add to their
492 credibility, but rather puts into question the way results are made sense of (Katz, 2013). Instead of
493 literal interpretations of Bayesian skylines, an improved understanding is required of how far the
494 demographic past can be reconstructed from the genomes of extremely clonal bacteria without
495 taking into account different confounding factors.

496 497 Box 2: Progeny skew in prokaryotes?

498 Recently, progeny skew was brought up as a neglected aspect of MTBC evolution with
499 potentially significant effects on genetic diversity (Morales-Arce et al., 2020) and population
500 genetic inference (Menardo et al., 2021a). Progeny skew refers to the unequal distribution
501 of offspring among parental individuals in a population. Frequently mentioned examples
502 are viruses, where a single parental sequence can give rise to numerous copies, or marine
503 organisms reproducing through broadcast spawning. Wright-Fisher and coalescence models
504 assume the variation in offspring number is small (Tellier and Lemaire, 2014), which leads
505 to mis-inference of population genetic statistics when applied to such organisms (Sackman
506 et al., 2019).

507 While progeny skew in viruses or has a direct interpretation in the way these organisms
508 reproduce, it is less straightforward to apply to prokaryotes. Archea and bacteria reproduce
509 through binary fission, which can be thought of as each parent having two offspring and
510 dying after division (Cury et al., 2021); or, in an age-structured population, as each parent
511 having one offspring and surviving. Progeny skew can arise over multiple generations
512 through rapid adaptation, superspreading events, or repeated bottlenecks, and it is thus a
513 meaningful parameter in population-based models with a continuous timescale (Menardo
514 et al., 2021a). In individual-based, discrete-generation models, it is preferable to simulate the
515 processes giving rise to progeny skew explicitly.

516 How do bottlenecks affect genetic diversity?

517 In the MTBC, genetic drift is often associated with transmission bottlenecks or founder events,
518 when few or even single strains initiate an infection or an outbreak (Pepperell et al., 2010; Smith
519 et al., 2006). TB infections can be initiated by single to few cells (Ryndak and Laal, 2019); each
520 transmission might thus be a massive founder event where, from the millions of cells forming a
521 within-host population, only a few cells are sampled to start a new population. Similar, small-scale
522 colonization dynamics occur during within-host dissemination, as single to few cells "found" new
523 granulomas in the highly structured habitat of the lung (Martin et al., 2017).

524 While genetic bottlenecks entail an immediate loss of genetic diversity, the mid- and long-term
525 effects of periodic bottlenecks on genetic diversity and differentiation in clonal pathogens, where
526 extreme bottlenecks alternate with clonal expansions, are less clear. Periodic bottlenecks have been
527 investigated in the context of experimental evolution, where studies mainly focused on the effects
528 of bottlenecks on the rate of adaptation (e.g. Windels et al., 2021). More general considerations
529 can be found in the population genetics literature. One insight of potential relevance for the
530 evolutionary dynamics of the MTBC is that, under predominant purifying selection, rates of
531 evolution are accelerated when N_e is small because more deleterious mutations fix due to genetic
532 drift (Lanfear et al., 2014). A classic example of this is the increased rate of sequence evolution
533 in aphid endosymbionts versus free-living bacteria of the genus *Buchnera* (Moran, 1996). In the
534 absence of homogenizing gene flow, founder events might thus be expected to increase genetic
535 differentiation and overall diversity among lineages of the MTBC. Following this logic, the low
536 global diversity of the MTBC (Figure 2a) is not evidence for strong bottlenecks. The puzzling
537 observation rather is that there is not more diversity given the repeated bottlenecks during
538 within- and between-host evolution and the absence of gene flow. Low diversity despite frequent
539 bottlenecking could thus indicate purifying selection.

540 The purpose of these considerations is to show that genetic bottlenecks are more complex and
541 interesting than they appear in the literature, where they often serve as *ad hoc* explanation for low
542 diversity. More work on periodic bottlenecks in bacterial pathogens is needed. This work could
543 take into account some real-world complications such as the unclear number of cells actually
544 transmitted, which is most likely larger than the minimum number required to start an infection
545 (Namouchi et al., 2012). Furthermore, infection might not occur at a single time point, but extend
546 through time as hosts are repeatedly exposed to bacteria-laden aerosol droplets (Ryndak and Laal,
547 2019). This situation resembles the source-sink dynamics of metapopulation models with repeated
548 colonization events rather than a single bottleneck.

549 POSITIVE SELECTION

550 As unclear as the role of genetic drift and purifying selection in the evolution of the MTBC is
551 the role of positive selection. Most insights about how the MTBC has adapted to environmental
552 challenges either regard pathoadaptation in the distant past before the MRCA, as revealed through
553 comparative genomics (reviewed by Pepperell, 2022), or the recent evolution of antibiotic resistance
554 (reviewed by Gygli et al., 2017). Much less is known about the genetics underlying adaptation
555 to different mammalian host species, evident in host tropism (Brites et al., 2018; Zwyrer et al.,
556 2021), or about local adaptation to different human populations, as suggested by sympatric
557 patient-pathogen associations observed in cosmopolitan settings (Gagneux et al., 2006).

558 Identifying signatures of positive selection in linked genomes is challenging since most tests

559 rely on the comparison of haplotypes within genomes (Shapiro et al., 2009). Two diversity-based
560 signatures that are not haplotype-based have been used extensively to identify positive selection
561 in MTBC genomes: homoplasy and excess of nonsynonymous polymorphisms. In the following,
562 we discuss the properties of these measures and whether they can be used to elucidate the role of
563 positive selection beyond the case of antibiotic resistance, which so far provides the confirmed
564 cases of adaptive evolution in the MTBC.

565 Homoplasies: how common is convergent adaptation?

566 Molecular homoplasy designates the independent appearance of identical mutations in different
567 parts of a phylogeny through chance, recombination, or convergent selection (Stern, 2013). Chance
568 homoplasy between genomes showing so little overall diversity is rare (Comas et al., 2009, Figure
569 2d), and its probability can be assessed through permutation tests (Farhat et al., 2013). Mutation
570 hotspots can facilitate chance homoplasy (Galtier et al., 2006): in the MTBC, highly mutable
571 tandem repeats frequently cause homoplasy (Outhred et al., 2020), while it is not known how
572 rates of point mutations vary along the genome. Recombination has been argued against as a
573 cause of homoplasies because homoplasies in the MTBC do not occur in clusters, as would be
574 expected when recombination involves diverged DNA (Chiner-Oms et al., 2019b). Non-clustering
575 homoplasies, however, are also expected when recombinant genomes are similar (Bobay et al.,
576 2015). Furthermore, intrachromosomal recombination can generate homoplasies, as suggested by
577 their increased occurrence in homologous *PE/PPE* genes (Tantivitayakul et al., 2020).

578 Clear examples of convergent selection as a cause of homoplasy have been presented for genes
579 involved in antimicrobial resistance (Comas et al., 2012; Farhat et al., 2013). Against a background
580 of low diversity and rare homoplasy, some of these genes show exceptional patterns. In 1,161
581 strains sampled in Russia and South Africa, one specific mutation in the *katG* gene, which confers
582 isoniazid resistance, has originated more than 70 times independently (Mortimer et al., 2017).
583 This is an extreme pattern that arises because *katG* is a "tight target" of selection, that is, only
584 single to few mutations can cause resistance without incurring high fitness costs. In other genes
585 ("sloppy targets"), fewer homoplasies are observed but in more positions. The high incidence of
586 parallelism in resistance evolution, in combination with large datasets, allows the use of genome
587 wide association approaches to identify new drug resistance loci and to elucidate the genetic
588 architecture of resistance phenotypes (e.g. Crook et al., 2022).

589 The basic limitation of homoplasies as a signature of selection is that they only reveal cases of
590 convergent evolution. In the case of antibiotic resistance, convergence is ubiquitous. Thousands
591 of parallel evolutionary experiments are conducted when people around the world are treated
592 with the same antibiotics proposed by the WHO, imposing strong selective pressures with high
593 rewards for resistance mutations in target genes (Walker et al., 2022). For other selective pressures,
594 things are less clear. Recently, two cases of convergent selection were shown in studies of
595 experimental evolution with *M. canettii* and the MTBC. Selecting *M. canettii* strains for *in vivo*
596 persistence in mice, Allen et al. (2021) identified two parallel mutations and demonstrated their
597 effect on persistence through gene knock-out and complementation. Smith et al. (2022) selected for
598 biofilm formation in experimentally evolved MTBC strains and identified two loci that mutated
599 independently and are associated to biofilm-associated traits and fitness proxies. Both studies
600 found that parallel mutations emerged in similar strains, suggesting that the genetic background
601 constrains evolutionary trajectories. These studies also illustrate the rapidity with which mutations
602 otherwise rare or absent can prevail in the presence of new selective pressures; and the significance

603 of structural variation, as convergent evolution involved a large duplication (Smith et al., 2022)
604 and a deletion of two genes (Allen et al., 2021).

605 Convergence might not only be favored by strong selective pressures, but also through human
606 demography and migration. Repeated introductions of sublineages into a region, as described
607 for Tibet (Liu et al., 2021), are natural experiments where genetically highly similar strains are
608 confronted with a new environment. Liu et al. identified several genes that accumulate mutations
609 independently after repeated introductions to the Tibetan Plateau, including *sseA*, a gene involved
610 in the detoxification of reactive oxygen species, and three genes involved in DNA repair (*dnaE2*,
611 *recB*, *mfid*). With the already large and still growing amount of data on MTBC outbreaks, such
612 natural experiments of parallel evolution can provide valuable insights into the dynamics and
613 genes involved in local adaptation.

614 Nonsynonymous polymorphisms: how frequent is positive selection?

615 The second widely used statistic to infer selection and its direction is the ratio of non-synonymous
616 to synonymous polymorphisms d_N/d_S . The intuition behind this measure is that an increased
617 rate of nonsynonymous compared to synonymous changes indicates positive selection. As for
618 homoplasies, genes involved in antibiotic resistance provide the clearest examples (Osório et al.,
619 2013; Wilson et al., 2020), and indeed the two signatures often co-occur.

620 Compared to homoplasy, which is a fairly intuitive heuristic for convergent selection, d_N/d_S is
621 a more complicated statistic that can be estimated in different ways and whose properties and
622 limitations have been explored in numerous studies (overview in Yang, Ziheng, 2014). Frequently,
623 d_N/d_S is estimated by comparing pairs of sequences (e.g. with the method of Nei and Gojobori,
624 1986). This is e.g. the case for d_N/d_S in Figure 2c or in the study of Hershberg et al. (2008)
625 discussed above, who presented genome-wide average pairwise d_N/d_S as evidence for reduced
626 purifying selection. Although average pairwise d_N/d_S is sometimes used gene-wise in selection
627 scans, it is a coarse measure. The ratio averages over the sites of a locus and the branches in a
628 phylogeny. It thus has low sensitivity, as only in loci with strong signals and multiple sites under
629 selection will the signal not be canceled by sites under purifying selection (Yang and Bielawski,
630 2000). A signal for positive selection may also be canceled if it is only present on a specific branch
631 (Yang and Nielsen, 2002).

632 A family of more versatile maximum likelihood models have been developed that incorporate
633 explicit models of codon evolution and allow to test for increased rates of nonsynonymous
634 changes on particular branches or in particular codons of a gene (implemented in PAML; Yang,
635 2007). These methods are computationally intensive and not suitable for exploratory analyses on
636 large phylogenies, while small MTBC datasets might not contain enough diversity to estimate
637 parameters. They can be used, however, to obtain a more detailed picture of selective pressures in
638 genes of interest and to formally test for selection using model comparisons (Yang, 1998). A recent
639 example of an exploratory selection scan followed by more rigorous statistical testing is the study
640 of Menardo et al. (2021b). In a first step, they identified a hypervariable epitope at the *esxH* locus,
641 which codes for a secreted effector interacting with the human immune system. Codon models
642 were then used to test for site- and branch-specific selection. Significant signatures were found
643 in MTBC lineage 1 but not in other lineages and located to the N-terminal epitope of the gene.
644 Further dissection of these signatures showed that they occur in strains collected in South and
645 Southeast Asia, suggesting that this locus might be involved in adaptation to regional human host
646 populations.

647 Two recent studies have proposed methods to estimate d_N/d_S for large datasets while avoiding
648 site and branch averaging, respectively. Wilson et al. (2020) present a phylogeny-free (and thus fast)
649 method to infer selection at the codon level. Applying their method to more than 10,000 MTBC
650 genomes, they found a d_N/d_S significantly larger than 1 in 2,729 out of 3,979 genes. Chiner-Oms
651 et al. (2022) investigated the temporal trajectories of p_N/p_S in a large phylogeny of 5,000 strains
652 (p_N/p_S is based on simple counts while d_N/d_S includes correction through a substitution model,
653 Yang, Ziheng, 2014, p. 47ff). Focusing on shifts in p_N/p_S along the tree, they found evidence
654 for elevated nonsynonymous changes at some point in time in almost half the genes of the
655 MTBC. While both studies generate long lists of candidate genes, they also lead to the inevitable
656 follow-up question of exploratory selection scans: what to do with these candidates. Considering
657 the difficulty of experimental validation in a human pathogen, further characterization of the
658 candidates with the phylogenetically explicit methods of PAML could be useful.

659 Overall, homoplasies and d_N/d_S tell us little about the big unknown of clonal evolution: the
660 distribution of fitness effects (see introduction). Methods exist to infer the distribution of fitness
661 effects from d_N/d_S using population genetic (reviewed by Eyre-Walker and Keightley, 2007)
662 or phylogenetic (e.g. Tamuri et al., 2012) models. Recently, the relationship between selection
663 coefficients and d_N/d_S under clonal reproduction were explored in the context of somatic evolution
664 (Williams et al., 2020). The model developed in the study relaxes some of the strong assumptions
665 of previous approaches, in particular constant population sizes and evolution over long timescales,
666 by integrating d_N/d_S and the clone size distribution. It would be worthwhile to explore whether
667 this approach can be applied to bacterial within-host populations in order to learn more about the
668 distribution of fitness effects *in vivo*.

669 DISCUSSION

670 What evolutionary processes drive and have driven the evolution of the MTBC? The most robust
671 insight, forming the premise of this review, is that horizontal gene transfer is negligible, al-
672 though recombination more generally is not, considering the mutagenic role of intrachromosomal
673 recombination. Regarding mutation, genetic drift, and natural selection, much remains unclear.

674 The current understanding of evolutionary processes in the MTBC is based on a complex
675 mesh of indirect evidence, intuitions, deductions from general principles, and assumption-rich
676 models. Studies focusing on population genetic processes are few and far apart, their methods
677 and datasets heterogeneous. Some of the hypotheses developed in these articles, for example
678 that the evolution of the MTBC is driven by genetic drift (Hershberg et al., 2008) or purifying
679 selection (Pepperell et al., 2013), have solidified into strong beliefs through repetition, even though
680 the original studies have pointed out caveats and the subtler meanings of "is driven by" remain
681 unexplored. With the large amount of sequencing data now available, covering evolutionary
682 timescales from within-host evolution to global patterns of diversity, it would be a good moment
683 to revisit some past hypotheses. We envisage focused studies that – in contrast to the typically
684 broad scope of studies of the "early" genomics era – address specific hypotheses and pay more
685 attention to methodological limitations.

686 While studying methods is less interesting than studying organisms, the bottleneck in data
687 analysis increasingly lies in the comprehension of complex methods rather than the availability
688 of data (Johri et al., 2022). New tools for evolutionary simulations, such as the versatile forward
689 simulation tool SLiM (Box 1), could provide a long-needed crutch to move forward by allowing to
690 simulate ever more realistic biologies and life histories.

691 To illustrate the utility of simulations, we used SLiM to simulate the within-host dissemination
 692 dynamics of a clonal pathogen (script available on [https://git.scicore.unibas.ch/TBRU/slim_](https://git.scicore.unibas.ch/TBRU/slim_simulations)
 693 [simulations](https://git.scicore.unibas.ch/TBRU/slim_simulations)). The model is inspired by the study of Martin et al., 2017, who used DNA barcoding
 694 and infection mapping to infer the temporal and spatial dynamics of an MTBC infection in
 695 macaques. Populations in the simulation might be thought of as granulomas that grow and give
 696 rise to new granulomas – a metapopulation model with unidirectional migration from "full" to
 697 "empty" populations. Infection begins with a single bacterium giving rise to an exponentially
 698 growing population through clonal reproduction and 19 "empty" populations. Once this popu-
 699 lation reaches carrying capacity $K = 20,000$, it can seed new populations (Figure B1a), which
 700 again grow and can seed new populations when K is reached. More specifically, each generation a
 701 number of n migrants is drawn from a Poisson distribution with mean 1; n individuals are then
 702 drawn from a random population that has reached K and migrated to a random empty population
 703 until all populations are occupied. Exemplary growth dynamics of the simulation are depicted
 704 in Figure B1b. Mutations are simulated at a rate $\mu = 5 \times 10^{-10}$ /bp/gen in a genome of 4 Mb.
 705 Selection is either assumed to be absent ($s = 0$) or purifying ($s = -9.5e - 4$), as proposed by
 706 Pepperell et al., 2013. The simulation ends after 70 generations, which with a generation time
 707 of 24 h corresponds to a 10 week infection. For both selection coefficients, the simulation was
 708 replicated 100 times.

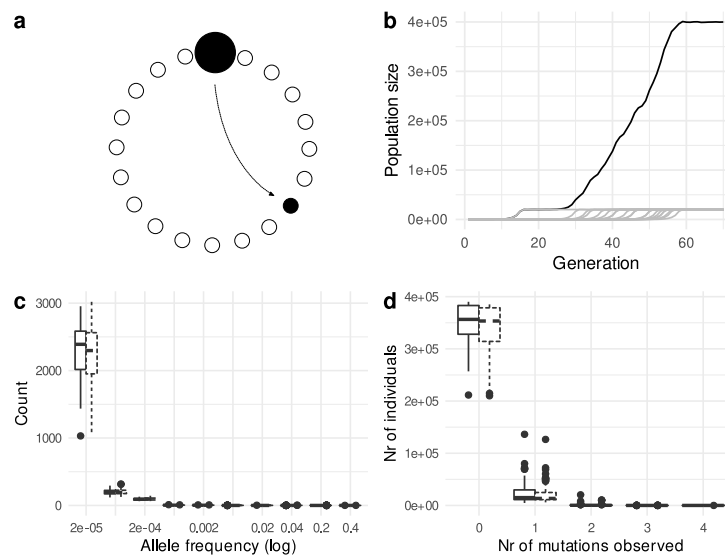


Figure 4: A metapopulation model for within-host evolution. a) Infection begins with a single bacterium giving rise to an exponentially growing population through clonal reproduction. Once this population reaches carrying capacity $K = 20,000$, it can seed new populations which again grow exponentially. b) Exemplary growth dynamics of the model, the solid line showing total population size, dashed lines showing subpopulation sizes. c) Site frequency spectrum at generation 70. d) Number of individuals with 0 to 4 SNPs at generation 70.

709 Independently of purifying selection, the dynamics of clonal growth and dissemination over
 710 70 bacterial generations give rise to an extreme skew towards rare alleles (Figure B1c). A large
 711 proportion of the mutations are in fact singletons, that is, only present in a single individual. At
 712 generation 70, the vast majority of individuals have no mutation, except in few instances where

713 a mutation arose early (Figure B1d). Some simulations produced outlier values because not all
714 populations were "filled" after 70 generations.

715 The purpose of this simulation is to illustrate the simulation approach. Some potential
716 applications of evolutionary simulations for the MTBC are listed in the following.

717 Simulations of structured within-host populations could be used to investigate the implications
718 of sampling and culturing for our understanding of within-host diversity, and to develop new
719 experimental designs and sampling strategies. How representative, for example, can sputum
720 samples possibly be of within-host diversity? Coupling within- and between-host evolution,
721 periodic bottlenecking could be simulated to study how diversity accumulates through time as
722 a function of bottleneck size, purifying selection, or mutation rates. This would lead to a more
723 nuanced understanding of transmission bottlenecks, which have more complex consequences than
724 simple reduction of diversity.

725 Gene conversion between closely related strains could be simulated to test different methods to
726 infer recombination. More generally, methods should be tested on simulated data to understand
727 their behavior and make an informed choice, instead of resorting to the typical bioinformatics
728 approach of using multiple methods and reporting intersecting results, which leaves the door
729 open to confirmation bias. Finally, the ultimate challenge would be to try to simultaneously infer
730 demography and selection using approximate Bayesian computation (Johri et al., 2022). It is diffi-
731 cult, however, to conceive what kind of data would be suitable for this. At the microevolutionary
732 scale that is most straightforward to simulate, there is just so little diversity that it is dubious that
733 parameter-rich models could be fitted with any confidence.

734 Simulations are not a panacea, but they allow to raise the debate to a more transparent,
735 quantitative level than achieved by the so far largely verbal arguments. If nothing else, they could
736 allow to better understand what kind of inference is at all possible, given the lack of HGT and the
737 low levels of genetic diversity in monomorphic bacteria.

738 ACKNOWLEDGMENTS

739 Our best thanks for their comments on earlier versions of the manuscripts to Daniela Brites,
740 Etthel Windels, Michaela Zwyrer, Selim Bouaouina, Fabrizio Menardo, Ana Morales-Arce, and
741 the members of the Gagneux group. This work was funded through grants from the European
742 Research Council, grant number 883582, and the Swiss National Science Foundation, grant
743 numbers 310030_188888 and CRSII5_177163.

744 REFERENCES

- 745 Achtman M (2008). Evolution, population structure, and phylogeography of genetically monomor-
746 phic bacterial pathogens. *Annual Review of Microbiology* 62, 53–70. DOI: 10.1146/annurev.micro.
747 62.081307.162832.
- 748 — (2012). Insights from genomic comparisons of genetically monomorphic bacterial pathogens.
749 *Philosophical Transactions of the Royal Society B: Biological Sciences* 367, 860–867. DOI: 10.1098/rstb.
750 2011.0303.
- 751 Allen AC et al. (2021). Parallel *in vivo* experimental evolution reveals that increased stress resistance
752 was key for the emergence of persistent tuberculosis bacilli. *Nature Microbiology* 6, 1082–1093.
753 DOI: 10.1038/s41564-021-00938-4.

- 754 Andersson SG and PM Sharp (1996). Codon usage in the *Mycobacterium tuberculosis* complex.
755 *Microbiology* 142, 915–925. doi: 10.1099/00221287-142-4-915.
- 756 Balbi KJ, EP Rocha, and EJ Feil (2009). The temporal dynamics of slightly deleterious mutations
757 in *Escherichia coli* and *Shigella* spp. *Molecular Biology and Evolution* 26, 345–355. doi: 10.1093/
758 molbev/msn252.
- 759 Bobay LM and H Ochman (2017). Impact of recombination on the base composition of bacteria
760 and archaea. *Molecular Biology and Evolution* 34, 2627–2636. doi: 10.1093/molbev/msx189.
- 761 Bobay LM, CC Traverse, and H Ochman (2015). Impermanence of bacterial clones. *PNAS* 112,
762 8893–8900. doi: 10.1073/pnas.1501724112.
- 763 Bobay LM and H Ochman (2018). Factors driving effective population size and pan-genome
764 evolution in bacteria. *BMC Evolutionary Biology* 18, 1–12. doi: 10.1186/s12862-018-1272-4.
- 765 Boritsch EC et al. (2016). Key experimental evidence of chromosomal DNA transfer among selected
766 tuberculosis-causing mycobacteria. *PNAS* 113, 9876–9881. doi: 10.1073/pnas.1604921113.
- 767 Bos KI et al. (2014). Pre-Columbian mycobacterial genomes reveal seals as a source of New World
768 human tuberculosis. *Nature* 514, 494–497. doi: 10.1038/nature13591.
- 769 Brites D et al. (2018). A new phylogenetic framework for the animal-adapted *Mycobacterium*
770 *tuberculosis* complex. *Frontiers in Microbiology* 9 (NOV), 1–14. doi: 10.3389/fmicb.2018.02820.
- 771 Casadevall A (2008). Evolution of intracellular pathogens. *Annual Review of Microbiology* 62, 19–33.
772 doi: 10.1146/annurev.micro.61.080706.093305.
- 773 Charlesworth B (2009). Effective population size and patterns of molecular evolution and variation.
774 *Nature Reviews Genetics* 10, 195–205. doi: 10.1038/nrg2526.
- 775 — (2012). The effects of deleterious mutations on evolution at linked sites. *Genetics* 190, 5–22. doi:
776 10.1534/genetics.111.134288.
- 777 Chiner-Oms Á, MG López, M Moreno-Molina, V Furió, and I Comas (2022). Gene evolution-
778 ary trajectories in *Mycobacterium tuberculosis* reveal temporal signs of selection. *PNAS* 119,
779 e2113600119–e2113600119. doi: 10.1073/pnas.2113600119.
- 780 Chiner-Oms Á et al. (2019a). Genome-wide mutational biases fuel transcriptional diversity in the
781 *Mycobacterium tuberculosis* complex. *Nature Communications* 10, 1–11. doi: 10.1038/s41467-019-
782 11948-6.
- 783 Chiner-Oms Á et al. (2019b). Genomic determinants of speciation and spread of the *Mycobacterium*
784 *tuberculosis* complex. *Science Advances* (June). doi: 10.1101/314559.
- 785 Clark RR, P Lapierre, E Lasek-Nesselquist, TA Gray, and KM Derbyshire (2022). A polymorphic
786 gene within the *Mycobacterium smegmatis* *esx1* locus determines mycobacterial self-identity and
787 conjugal compatibility. *mBio* 13. doi: 10.1128/mbio.00213-22.
- 788 Colangeli R et al. (2020). *Mycobacterium tuberculosis* progresses through two phases of latent
789 infection in humans. *Nature Communications* 11, 1–10. doi: 10.1038/s41467-020-18699-9.
- 790 Cole S et al. (1998). Deciphering the biology of *Mycobacterium tuberculosis* from the complete
791 genome sequence. *Nature* 393 (NOVEMBER), 537–544. doi: 10.1038/29241.
- 792 Comas I, S Homolka, S Niemann, and S Gagneux (2009). Genotyping of genetically monomorphic
793 bacteria: DNA sequencing in *Mycobacterium tuberculosis* highlights the limitations of current
794 methodologies. *PLoS ONE* 4. doi: 10.1371/journal.pone.0007815.
- 795 Comas I et al. (2012). Whole-genome sequencing of rifampicin-resistant *Mycobacterium tuberculosis*
796 strains identifies compensatory mutations in RNA polymerase genes. *Nature Genetics* 44, 106–
797 110. doi: 10.1038/ng.1038.
- 798 Comas I et al. (2013). Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium tubercu-*
799 *losis* with modern humans. *Nature Genetics* 45, 1176–1182. doi: 10.1038/ng.2744.

- 800 Copin R et al. (2016). Within-host evolution selects for a dominant genotype of *Mycobacterium*
801 *tuberculosis* while T cells increase pathogen genetic diversity. *PLoS Pathogens* 12, 1–16. DOI:
802 10.1371/journal.ppat.1006111.
- 803 Crook DW et al. (2022). Genome-wide association studies of global *Mycobacterium tuberculosis*
804 resistance to 13 antimicrobials in 10,228 genomes identify new resistance mechanisms, 1–27.
805 DOI: 10.1371/journal.pbio.3001755.
- 806 Cury J, BC Haller, G Achaz, and F Jay (2021). Simulation of bacterial populations with SLiM. *Peer*
807 *Community in Evolutionary Biology*, 1–36. DOI: 10.1101/2020.09.28.316869.
- 808 Davies J and D Davis (2010). Origins and evolution of antibiotic resistance. *Microbiología (Madrid,*
809 *Spain)* 12, 9–16. DOI: 10.1128/mmbr.00016-10.
- 810 De Miranda AB, F Alvarez-Valin, K Jabbari, WM Degraeve, and G Bernardi (2000). Gene expression,
811 amino acid conservation, and hydrophobicity are the main factors shaping codon preferences
812 in *Mycobacterium tuberculosis* and *Mycobacterium leprae*. *Journal of Molecular Evolution* 50, 45–55.
813 DOI: 10.1007/s002399910006.
- 814 Denamur E, O Clermont, S Bonacorsi, and D Gordon (2021). The population genetics of pathogenic
815 *Escherichia coli*. *Nature Reviews Microbiology* 19, 37–54. DOI: 10.1038/s41579-020-0416-x.
- 816 Dubos R and J Dubos (1952). *The White Plague*. Little, Brown and Company.
- 817 Duchêne S et al. (2016). Genome-scale rates of evolutionary change in bacteria. *Microbial genomics*
818 2, e000094–e000094. DOI: 10.1099/mgen.0.000094.
- 819 Duret L and N Galtier (2009). Biased gene conversion and the evolution of mammalian genomic
820 landscapes. *Annual Review of Genomics and Human Genetics* 10, 285–311. DOI: 10.1146/annurev-
821 genom-082908-150001.
- 822 Dutta NK and PC Karakousis (Sept. 2014). Latent tuberculosis infection: myths, models, and
823 molecular mechanisms. *Microbiology and Molecular Biology Reviews* 78, 343–371. DOI: 10.1128/
824 MMBR.00010-14.
- 825 Emerson BC and MJ Hickerson (2015). Lack of support for the time-dependent molecular evolution
826 hypothesis. *Molecular Ecology* 24, 702–709. DOI: 10.1111/mec.13070.
- 827 Eyre-Walker A and PD Keightley (2007). The distribution of fitness effects of new mutations.
828 *Nature Reviews Genetics* 8. DOI: 10.1038/nrg2146.
- 829 Farhat MR et al. (2013). Genomic analysis identifies targets of convergent positive selection in
830 drug-resistant *Mycobacterium tuberculosis*. *Nature Genetics* 45, 1183–1189. DOI: 10.1038/ng.2747.
- 831 Fishbein S, N van Wyk, RM Warren, and SL Sampson (2015). Phylogeny to function: PE/PPE
832 protein evolution and impact on *Mycobacterium tuberculosis* pathogenicity. *Molecular Microbiology*
833 96, 901–916. DOI: 10.1111/mmi.12981.
- 834 Fisher RA (1930). *The Genetical Theory of Natural Selection*. Pages: 308. 308 pp.
- 835 Fitzgerald DM and SM Rosenberg (2019). What is mutation? A chapter in the series: how microbes
836 “jeopardize” the modern synthesis. *PLoS Genetics* 15, 1–14. DOI: 10.1371/journal.pgen.1007995.
- 837 Ford CB et al. (2011). Use of whole genome sequencing to estimate the mutation rate of *Mycobac-*
838 *terium tuberculosis* during latent infection. *Nature Genetics* 43, 482–488. DOI: 10.1038/ng.811.
- 839 Ford CB et al. (2013). *Mycobacterium tuberculosis* mutation rate estimates from different lineages
840 predict substantial differences in the emergence of drug-resistant tuberculosis. *Nature Genetics*
841 45, 784–790. DOI: 10.1038/ng.2656.
- 842 Gagneux S (2018). Ecology and evolution of *Mycobacterium tuberculosis*. *Nature Reviews Microbiology*
843 16, 202–213. DOI: 10.1038/nrmicro.2018.8.
- 844 Gagneux S et al. (2006). Variable host-pathogen compatibility in *Mycobacterium tuberculosis*. *PNAS*
845 103, 2869–2873. DOI: 10.1073/pnas.0511240103.

- 846 Galtier N, D Enard, Y Radondy, E Bazin, and K Belkhir (2006). Mutation hot spots in mammalian
847 mitochondrial DNA. *Genome Research* 16, 215–222. doi: 10.1101/gr.4305906.
- 848 Gillespie JH (2004). *Population Genetics – A Concise Guide*. The Johns Hopkins University Press.
- 849 Godfroid M, T Dagan, and A Kupczok (2018). Recombination signal in *Mycobacterium tuberculosis*
850 stems from reference-guided assemblies and alignment artefacts. *Genome Biology and Evolution*
851 10, 1920–1926. doi: 10.1093/gbe/evy143.
- 852 Gómez-Valero L, EP Rocha, A Latorre, and FJ Silva (2007). Reconstructing the ancestor of *Mycobac-*
853 *terium leprae*: the dynamics of gene loss and genome reduction. *Genome Research* 17, 1178–1185.
854 doi: 10.1101/gr.6360207.
- 855 Gray TA and KM Derbyshire (2018). Blending genomes: distributive conjugal transfer in mycobac-
856 teria, a sexier form of HGT. *Molecular Microbiology* 108, 601–613. doi: 10.1111/mmi.13971.
- 857 Guerrini V et al. (2016). Experimental evolution of *Mycobacterium tuberculosis* in human macrophages
858 results in low-frequency mutations not associated with selective advantage. *PLoS ONE* 11, 1–15.
859 doi: 10.1371/journal.pone.0167989.
- 860 Gygli SM, S Borrell, A Trauner, and S Gagneux (2017). Antimicrobial resistance in *Mycobacterium*
861 *tuberculosis*: mechanistic and evolutionary perspectives. *FEMS Microbiology Reviews* 41, 354–373.
862 doi: 10.1093/femsre/fux011.
- 863 Haller BC, J Galloway, J Kelleher, PW Messer, and PL Ralph (2019). Tree-sequence recording in
864 SLiM opens new horizons for forward-time simulation of whole genomes. *Molecular Ecology*
865 *Resources* 19, 552–566. doi: 10.1111/1755-0998.12968.
- 866 Haller BC and PW Messer (2019). SLiM 3: Forward Genetic Simulations Beyond the Wright-Fisher
867 Model. *Molecular Biology and Evolution* 36, 632–637. doi: 10.1093/molbev/msy228.
- 868 Hanage WP (2016). Not so simple after all: bacteria, their population genetics, and recombination.
869 *Cold Spring Harbor Perspectives in Biology* 8. doi: 10.1101/cshperspect.a018069.
- 870 Hedge J and DJ Wilson (2014). Bacterial phylogenetic reconstruction from whole genomes is robust
871 to recombination but demographic inference is not. *mBio* 5, 5–8. doi: 10.1128/mBio.02158-14.
- 872 Heller R, L Chikhi, and HR Siegmund (2013). The confounding effect of population structure on
873 Bayesian skyline plot inferences of demographic history. *PLoS ONE* 8. doi: 10.1371/journal.
874 pone.0062992.
- 875 Hershberg R and DA Petrov (2008). Selection on codon bias. *Annual Review of Genetics* 42, 287–299.
876 doi: 10.1146/annurev.genet.42.110807.091442.
- 877 — (2010). Evidence that mutation is universally biased towards AT in bacteria. *PLoS Genetics* 6.
878 doi: 10.1371/journal.pgen.1001115.
- 879 Hershberg R et al. (2008). High functional diversity in *Mycobacterium tuberculosis* driven by genetic
880 drift and human demography. *PLoS Biology* 6, 2658–2671. doi: 10.1371/journal.pbio.0060311.
- 881 Hildebrand F, A Meyer, and A Eyre-Walker (2010). Evidence of selection upon genomic GC-content
882 in bacteria. *PLoS Genetics* 6. doi: 10.1371/journal.pgen.1001107.
- 883 Ho SY, S Duchêne, M Molak, and B Shapiro (2015). Time-dependent estimates of molecular
884 evolutionary rates: evidence and causes. *Molecular Ecology* 24, 6007–6012. doi: 10.1111/mec.
885 13450.
- 886 Ho SYW and B Shapiro (2011). Skyline-plot methods for estimating demographic history from
887 nucleotide sequences. *Molecular Ecology Resources* 11, 423–434. doi: 10.1111/j.1755-0998.2011.
888 02988.x.
- 889 Ho SYW et al. (2011). Time-dependent rates of molecular evolution. *Molecular Ecology* 20, 3087–3101.
890 doi: 10.1111/j.1365-294X.2011.05178.x.

- 891 Hoban S, G Bertorelle, and OE Gaggiotti (2012). Computer simulations: tools for population and
892 evolutionary genetics. *Nature Reviews Genetics* 13. doi: 10.1038/nrg3130.
- 893 Hughes AL, R Friedman, and M Murray (2002). Genomewide pattern of synonymous nucleotide
894 substitution in two complete genomes of *Mycobacterium tuberculosis*. *Emerging Infectious Diseases*
895 8, 1342–1346. doi: 10.3201/eid0811.020064.
- 896 Johri P et al. (2022). Recommendations for improving statistical inference in population genomics.
897 *PLoS Biology* 20, 1–23. doi: 10.1371/journal.pbio.3001669.
- 898 Kapur V, TS Whittam, and JM Musser (1994). Is *Mycobacterium tuberculosis* 15,000 years old? *Journal*
899 *of Infectious Diseases* 170, 1348–1349. doi: 10.1093/infdis/170.5.1348.
- 900 Karboul A et al. (2008). Frequent homologous recombination events in *Mycobacterium tuberculosis*
901 PE/PPE multigene families: potential role in antigenic variability. *Journal of Bacteriology* 190,
902 7838–7846. doi: 10.1128/JB.00827-08.
- 903 Katju V and U Bergthorsson (2019). Old trade, new tricks: Insights into the spontaneous mu-
904 tation process from the partnering of classical mutation accumulation experiments with
905 high-throughput genomic approaches. *Genome Biology and Evolution* 11, 136–165. doi: 10.1093/
906 gbe/evy252.
- 907 Katz Y (2013). Against storytelling of scientific results. *Nature Methods* 10, 1045–1045. doi: 10.1038/
908 nmeth.2699.
- 909 Kay GL et al. (2015). Eighteenth-century genomes show that mixed infections were common at
910 time of peak tuberculosis in Europe. *Nature Communications* 6. doi: 10.1038/ncomms7717.
- 911 Kimura M (1983). *The Neutral Theory of Molecular Evolution*. Cambridge University Press. doi:
912 10.1017/CBO9781107415324.004.
- 913 Kucukyildirim S et al. (2016). The rate and spectrum of spontaneous mutations in *Mycobacterium*
914 *smegmatis*, a bacterium naturally devoid of the postreplicative mismatch repair pathway. *G3:*
915 *Genes, Genomes, Genetics* 6, 2157–2163. doi: 10.1534/g3.116.030130.
- 916 Kuo CH, NA Moran, and H Ochman (2009). The consequences of genetic drift for bacterial genome
917 complexity. *Genome Research* 19, 1450–1454. doi: 10.1101/gr.091785.109.
- 918 Lanfear R, H Kokko, and A Eyre-Walker (Jan. 2014). Population size and the rate of evolution.
919 *Trends in Ecology & Evolution* 29, 33–41. doi: 10.1016/j.tree.2013.09.009.
- 920 Lapierre M, C Blin, A Lambert, G Achaz, and EP Rocha (2016). The impact of selection, gene
921 conversion, and biased sampling on the assessment of microbial demography. *Molecular biology*
922 *and evolution* 33, 1711–1725. doi: 10.1093/molbev/msw048.
- 923 Lassalle F et al. (2015). GC-content evolution in bacterial genomes: the biased gene conversion
924 hypothesis expands. *PLoS Genetics* 11, 1–20. doi: 10.1371/journal.pgen.1004941.
- 925 Lawrence JG, RW Hendrix, and S Casjens (2001). Where are the pseudogenes in bacterial genomes?
926 *Trends in Microbiology* 9, 535–540. doi: 10.1016/S0966-842X(01)02198-9.
- 927 Liu Q et al. (May 29, 2020). *Mycobacterium tuberculosis* clinical isolates carry mutational signatures
928 of host immune environments. *Science Advances* 6, eaba4901. doi: 10.1126/sciadv.aba4901.
- 929 Liu Q et al. (2021). Local adaptation of *Mycobacterium tuberculosis* on the Tibetan Plateau. *PNAS*
930 118, 1–10. doi: 10.1073/pnas.2017831118.
- 931 Liu X, MM Gutacker, JM Musser, and YX Fu (2006). Evidence for recombination in *Mycobacterium*
932 *tuberculosis*. *Journal of Bacteriology* 188, 8169–8177. doi: 10.1128/JB.01062-06.
- 933 Lynch M (2007). *The Origins of Genome Architecture*. Sinauer Associates Sunderland, MA.
- 934 — (2010). Evolution of the mutation rate. *Trends in Genetics* 26, 345–352. doi: 10.1016/j.tig.2010.05.
935 003.

- 936 Madacki J et al. (2021). ESX-1-Independent horizontal gene transfer by *Mycobacterium tuberculosis*
937 complex strains. *mBio* 12, 1–19. doi: 10.1128/mbio.00965-21.
- 938 Martin CJ et al. (2017). Digitally barcoding *Mycobacterium tuberculosis* reveals *in vivo* infection
939 dynamics in the macaque model of tuberculosis. *mBio* 8, 1–12. doi: 10.1128/mBio.00312-17.
- 940 Martin DP, P Lemey, and D Posada (2011). Analysing recombination in nucleotide sequences.
941 *Molecular Ecology Resources* 11, 943–955. doi: 10.1111/j.1755-0998.2011.03026.x.
- 942 Maynard Smith J (1995). Do bacteria have population genetics? In: *Population genetics of bacteria:*
943 *Symposium 52*. Cambridge University Press, pp. 1–12.
- 944 Maynard Smith J, NH Smith, M O'Rourke, and BG Spratt (1993). How clonal are bacteria? *PNAS*
945 90, 4384–4388. doi: 10.1073/pnas.90.10.4384.
- 946 McEvoy CR et al. (2007). The role of IS6110 in the evolution of *Mycobacterium tuberculosis*. *Tubercu-*
947 *losis* 87, 393–404. doi: 10.1016/j.tube.2007.05.010.
- 948 Mcgrath M, NC Gey van pittius, PD Van helden, RM Warren, and DF Warner (2014). Mutation
949 rate and the emergence of drug resistance in *Mycobacterium tuberculosis*. *Journal of Antimicrobial*
950 *Chemotherapy* 69, 292–302. doi: 10.1093/jac/dkt364.
- 951 Membrebe JV et al. (2019). Bayesian inference of evolutionary histories under time-dependent
952 substitution rates. *Molecular Biology and Evolution* 36, 1793–1803. doi: 10.1093/molbev/msz094.
- 953 Menardo F, S Gagneux, and F Freund (2021a). Multiple merger genealogies in outbreaks of
954 *Mycobacterium tuberculosis*. *Molecular Biology and Evolution* 38, 290–306. doi: 10.1101/2019.12.21.
955 885723.
- 956 Menardo F, S Duchêne, D Brites, and S Gagneux (2019). The molecular clock of *Mycobacterium*
957 *tuberculosis*. *PLoS Pathogens* 15, 1–24. doi: 10.1371/journal.ppat.1008067.
- 958 Menardo F et al. (2021b). Local adaptation in populations of *Mycobacterium tuberculosis* endemic to
959 the Indian Ocean Rim. *F1000Research*. doi: 10.1101/2020.10.20.346866.
- 960 Merker M et al. (Aug. 30, 2022). Transcontinental spread and evolution of *Mycobacterium tuber-*
961 *culosis* W148 European/Russian clade toward extensively drug resistant tuberculosis. *Nature*
962 *Communications* 13, 5105. doi: 10.1038/s41467-022-32455-1.
- 963 Merrikh CN and H Merrikh (2018). Gene inversion potentiates bacterial evolvability and virulence.
964 *Nature Communications* 9. doi: 10.1038/s41467-018-07110-3.
- 965 Molina N and E Van Nimwegen (2008). Universal patterns of purifying selection at noncoding
966 positions in bacteria. *Genome Research* 18, 148–160. doi: 10.1101/gr.6759507.
- 967 Morales-Arce AY, RB Harris, AC Stone, and JD Jensen (2020). Evaluating the contributions
968 of purifying selection and progeny-skew in dictating within-host *Mycobacterium tuberculosis*
969 evolution. *Evolution* 74, 992–1001. doi: 10.1111/evo.13954.
- 970 Morales-Arce AY, SJ Sabin, AC Stone, and JD Jensen (2021). The population genomics of within-
971 host *Mycobacterium tuberculosis*. *Heredity* 126, 1–9. doi: 10.1038/s41437-020-00377-7.
- 972 Moran NA (1996). Accelerated evolution and Muller's ratchet in endosymbiotic bacteria. *PNAS* 93,
973 2873–2878. doi: 10.1073/pnas.93.7.2873.
- 974 Moreno-Molina M et al. (2021). Genomic analyses of *Mycobacterium tuberculosis* from human lung
975 resections reveal a high frequency of polyclonal infections. *Nature Communications* 12, 1–11. doi:
976 10.1038/s41467-021-22705-z.
- 977 Mortimer TD, AM Weber, and CS Pepperell (2017). Signatures of selection at drug resistance loci
978 in *Mycobacterium tuberculosis*. *mSystems* 8, 1–11. doi: 10.1101/173229.
- 979 Mulholland CV et al. (2019). Dispersal of *Mycobacterium tuberculosis* driven by historical European
980 trade in the South Pacific. *Frontiers in Microbiology* 10 (December), 1–13. doi: 10.3389/fmicb.
981 2019.02778.

- 982 Muller HJ (1964). The relation of recombination to mutational advance. *Mutation Research - Funda-*
983 *mental and Molecular Mechanisms of Mutagenesis* 1, 2–9. DOI: 10.1016/0027-5107(64)90047-8.
- 984 Namouchi A, X Didelot, U Schöck, B Gicquel, and EP Rocha (2012). After the bottleneck: genome-
985 wide diversification of the *Mycobacterium tuberculosis* complex by mutation, recombination, and
986 natural selection. *Genome Research* 22, 721–734. DOI: 10.1101/gr.129544.111.
- 987 Neher RA (2013). Genetic draft, selective interference, and population genetics of rapid adaptation.
988 *Annual Review of Ecology, Evolution, and Systematics* 44, 195–215. DOI: 10.1146/annurev-ecolsys-
989 110512-135920.
- 990 Nei M and T Gojobori (1986). Simple methods for estimating the numbers of synonymous and
991 nonsynonymous nucleotide substitutions. *Molecular Biology and Evolution* 3, 418–426. DOI:
992 10.1093/oxfordjournals.molbev.a040410.
- 993 O’Neill MB et al. (2019). Lineage specific histories of *Mycobacterium tuberculosis* dispersal in Africa
994 and Eurasia. *Molecular Ecology* 28, 3241–3256. DOI: 10.1111/mec.15120.
- 995 O’Neill MB, TD Mortimer, and CS Pepperell (2015). Diversity of *Mycobacterium tuberculosis* across
996 evolutionary scales. *PLoS Pathogens* 11, 1–29. DOI: 10.1371/journal.ppat.1005257.
- 997 Ohta T (1992). The nearly neutral theory of molecular evolution. *Annual Review of Ecology and*
998 *Systematics* 23, 263–286. DOI: 10.1146/annurev.es.23.110192.001403.
- 999 Orme IM (2014). A new unifying theory of the pathogenesis of tuberculosis. *Tuberculosis* 94, 8–14.
1000 DOI: 10.1016/j.tube.2013.07.004.
- 1001 Osório NS et al. (2013). Evidence for diversifying selection in a set of *Mycobacterium tuberculosis*
1002 genes in response to antibiotic- and nonantibiotic-related pressure. *Molecular Biology and*
1003 *Evolution* 30, 1326–1336. DOI: 10.1093/molbev/mst038.
- 1004 Outhred AC et al. (2020). Extensive homoplasy but no evidence of convergent evolution of repeat
1005 numbers at MIRU loci in modern *Mycobacterium tuberculosis* lineages. *Frontiers in Public Health*
1006 8 (August), 1–12. DOI: 10.3389/fpubh.2020.00455.
- 1007 Pan A, C Dutta, and J Das (1998). Codon usage in highly expressed genes of *Haemophilus influenzae*
1008 and *Mycobacterium tuberculosis*: translational selection versus mutational bias. *Gene* 215, 405–413.
1009 DOI: 10.1016/S0378-1119(98)00257-1.
- 1010 Panda A, M Drancourt, T Tuller, and P Pontarotti (2018). Genome-wide analysis of horizontally
1011 acquired genes in the genus *Mycobacterium*. *Scientific Reports* 8, 1–13. DOI: 10.1038/s41598-018-
1012 33261-w.
- 1013 Payne JL et al. (2019). Transition bias influences the evolution of antibiotic resistance in *Mycobac-*
1014 *terium tuberculosis*. *PLoS Biology* 17, 1–23. DOI: 10.1371/journal.pbio.3000265.
- 1015 Pepperell C et al. (2010). Bacterial genetic signatures of human social phenomena among *M.*
1016 *tuberculosis* from an aboriginal canadian population. *Molecular Biology and Evolution* 27, 427–440.
1017 DOI: 10.1093/molbev/msp261.
- 1018 Pepperell CS (2022). Evolution of tuberculosis pathogenesis. *Annual Review of Microbiology* 76,
1019 661–680.
- 1020 Pepperell CS et al. (2013). The role of selection in shaping diversity of natural *M. tuberculosis*
1021 populations. *PLoS Pathogens* 9. DOI: 10.1371/journal.ppat.1003543.
- 1022 Phelan JE et al. (2016). Recombination in pe/ppe genes contributes to genetic variation in *Mycobac-*
1023 *terium tuberculosis* lineages. *BMC Genomics* 17, 1–12. DOI: 10.1186/s12864-016-2467-y.
- 1024 Plutynski A (2007). Drift: a historical and conceptual overview. *Biological Theory* 2, 156–167. DOI:
1025 10.1162/biot.2007.2.2.156.
- 1026 Price MN and AP Arkin (2015). Weakly deleterious mutations and low rates of recombination
1027 limit the impact of natural selection on bacterial genomes. *mBio* 6. DOI: 10.1128/mBio.01302-15.

- 1028 Rahman S, SL Kosakovsky Pond, A Webb, and J Hey (May 18, 2021). Weak selection on synonymous
1029 codons substantially inflates dN/dS estimates in bacteria. *PNAS* 118, e2023575118. doi: 10.1073/
1030 pnas.2023575118.
- 1031 Reichenberger ER, G Rosen, U Hershberg, and R Hershberg (2015). Prokaryotic nucleotide com-
1032 position is shaped by both phylogeny and the environment. *Genome Biology and Evolution* 7,
1033 1380–1389. doi: 10.1093/gbe/evv063.
- 1034 Rocha C and A Danchin (2002). Base composition bias might result from competition for metabolic
1035 resources. *TRENDS in genetics* 18, 291–294.
- 1036 Rocha EP (2018). Neutral theory, microbial practice: challenges in bacterial population genetics.
1037 *Molecular Biology and Evolution* 35, 1338–1347. doi: 10.1093/molbev/msy078.
- 1038 Rocha EP and EJ Feil (2010). Mutational patterns cannot explain genome composition: are there any
1039 neutral sites in the genomes of bacteria? *PLoS Genetics* 6, 1–4. doi: 10.1371/journal.pgen.1001104.
- 1040 Rocha EP et al. (2006). Comparisons of dN/dS are time dependent for closely related bacterial
1041 genomes. *Journal of Theoretical Biology* 239, 226–235. doi: 10.1016/j.jtbi.2005.08.037.
- 1042 Ryndak MB and S Laal (2019). *Mycobacterium tuberculosis* primary infection and dissemination:
1043 a critical role for alveolar epithelial cells. *Frontiers in Cellular and Infection Microbiology* 9. doi:
1044 10.3389/fcimb.2019.00299.
- 1045 Sabin S et al. (2020). A seventeenth-century *Mycobacterium tuberculosis* genome supports a Neolithic
1046 emergence of the *Mycobacterium tuberculosis* complex. *Genome Biology* 21, 1–24. doi: 10.1186/
1047 s13059-020-02112-1.
- 1048 Sackman AM, RB Harris, and JD Jensen (2019). Inferring demography and selection in organisms
1049 characterized by skewed offspring distributions. *Genetics* 211, 1019–1028. doi: 10.1534/genetics.
1050 118.301684.
- 1051 Selander RK, JM Musser, DA Caugant, MN Gilmour, and TS Whittam (1987). Population genetics
1052 of pathogenic bacteria. *Microbial Pathogenesis* 3, 1–7. doi: 10.1016/0882-4010(87)90032-5.
- 1053 Shapiro BJ, LA David, J Friedman, and EJ Alm (2009). Looking for Darwin’s footprints in the
1054 microbial world. *Trends in Microbiology* 17, 196–204. doi: 10.1016/j.tim.2009.02.002.
- 1055 Smith JM, CG Dowson, and BG Spratt (1991). Localized sex in bacteria. *Nature* 349, 29–31. doi:
1056 10.1038/349029a0.
- 1057 Smith NH, SV Gordon, R de la Rúa-Domenech, RS Clifton-Hadley, and RG Hewinson (2006).
1058 Bottlenecks and broomsticks: the molecular evolution of *Mycobacterium bovis*. *Nature Reviews*
1059 *Microbiology* 4, 670–681. doi: 10.1038/nrmicro1472.
- 1060 Smith TM et al. (2022). Rapid adaptation of a complex trait during experimental evolution of
1061 *Mycobacterium tuberculosis*. *Elife* 11, e78454.
- 1062 Sniegowski PD and PJ Gerrish (2010). Beneficial mutations and the dynamics of adaptation in
1063 asexual populations. *Philosophical Transactions of the Royal Society B: Biological Sciences* 365,
1064 1255–1263. doi: 10.1098/rstb.2009.0290.
- 1065 Stern DL (2013). The genetic causes of convergent evolution. *Nature Reviews Genetics* 14, 751–764.
1066 doi: 10.1038/nrg3483.
- 1067 Supply P et al. (2013). Genomic analysis of smooth tubercle bacilli provides insights into ancestry
1068 and pathoadaptation of *Mycobacterium tuberculosis*. *Nature Genetics* 45, 172–179. doi: 10.1038/
1069 ng.2517.
- 1070 Tamuri AU, M dos Reis, and RA Goldstein (Mar. 1, 2012). Estimating the distribution of selection
1071 coefficients from phylogenetic data using sitewise mutation-selection models. *Genetics* 190,
1072 1101–1115. doi: 10.1534/genetics.111.136432.

- 1073 Tantivitayakul P et al. (2020). Homoplastic single nucleotide polymorphisms contributed to
1074 phenotypic diversity in *Mycobacterium tuberculosis*. *Scientific Reports* 10, 1–10. doi: 10.1038/
1075 s41598-020-64895-4.
- 1076 Tarashi S, A Fateh, M Mirsaeidi, SD Siadat, and F Vaziri (2017). Mixed infections in tuberculosis:
1077 the missing part in a puzzle. *Tuberculosis* 107, 168–174. doi: 10.1016/j.tube.2017.09.004.
- 1078 Tellier A and C Lemaire (2014). Coalescence 2.0: a multiple branching of recent theoretical
1079 developments and their applications. *Molecular Ecology* 23, 2637–2652. doi: 10.1111/mec.12755.
- 1080 Templeton AR (2021). *Population genetics and microevolutionary theory*. Second edition. Hoboken, NJ:
1081 John Wiley & Sons.
- 1082 Thorpe HA, SC Bayliss, LD Hurst, and EJ Feil (2017). Comparative analyses of selection operating
1083 on nontranslated intergenic regions of diverse bacterial species. *Genetics* 206, 363–376. doi:
1084 10.1534/genetics.116.195784.
- 1085 Tibayrenc M and FJ Ayala (2017). Is predominant clonal evolution a common evolutionary adapta-
1086 tion to parasitism in pathogenic parasitic protozoa, fungi, bacteria, and viruses? *Advances in*
1087 *Parasitology* 97, 243–325. doi: 10.1016/bs.apar.2016.08.007.
- 1088 Trauner A et al. (2017). The within-host population dynamics of *Mycobacterium tuberculosis* vary
1089 with treatment efficacy. *Genome Biology* 18, 1–17. doi: 10.1186/s13059-017-1196-0.
- 1090 Uplekar S, B Heym, V Friocourt, J Rougemont, and ST Cole (2011). Comparative genomics of ESX
1091 genes from clinical isolates of *Mycobacterium tuberculosis* provides evidence for gene conversion
1092 and epitope variation. *Infection and Immunity* 79, 4042–4049. doi: 10.1128/IAI.05344-11.
- 1093 Vos M and X Didelot (2009). A comparison of homologous recombination rates in bacteria and
1094 archaea. *ISME Journal* 3, 199–208. doi: 10.1038/ismej.2008.93.
- 1095 Walker TM et al. (2022). The 2021 WHO catalogue of *Mycobacterium tuberculosis* complex mutations
1096 associated with drug resistance: a genotypic analysis. *The Lancet Microbe* 3, e265–e273. doi:
1097 10.1016/s2666-5247(21)00301-3.
- 1098 Wang L et al. (2022). Multiple genetic paths including massive gene amplification allow *My-*
1099 *cobacterium tuberculosis* to overcome loss of ESX-3 secretion system substrates. *PNAS* 119. doi:
1100 10.1073/pnas.2112608119.
- 1101 Wang TC and FC Chen (2013). The evolutionary landscape of the *Mycobacterium tuberculosis*
1102 genome. *Gene* 518, 187–193. doi: 10.1016/j.gene.2012.11.033.
- 1103 Weissman JL, WF Fagan, and PL Johnson (2019). Linking high GC content to the repair of double
1104 strand breaks in prokaryotic genomes. *PLoS Genetics* 15, 1–19. doi: 10.1371/journal.pgen.
1105 1008493.
- 1106 Weller C and M Wu (2015). A generation-time effect on the rate of molecular evolution in bacteria.
1107 *Evolution* 69, 643–652. doi: 10.1111/evo.12597.
- 1108 Williams MJ et al. (2020). Measuring the distribution of fitness effects in somatic evolution by
1109 combining clonal dynamics with dN/dS ratios. *eLife* 9, 1–19. doi: 10.7554/eLife.48714.
- 1110 Wilson DJ et al. (2020). GenomeMap: within-species genome-wide dN/dS estimation from over
1111 10,000 genomes. *Molecular Biology and Evolution*, 1–11. doi: 10.1093/molbev/msaa069.
- 1112 Windels EM et al. (2021). Population bottlenecks strongly affect the evolutionary dynamics of
1113 antibiotic persistence. *Molecular Biology and Evolution* 38, 3345–3357. doi: 10.1093/molbev/
1114 msab107.
- 1115 Woese CR and N Goldenfeld (2009). How the microbial world saved evolution from the Scylla
1116 of molecular biology and the charybdis of the Modern Synthesis. *Microbiology and Molecular*
1117 *Biology Reviews* 73, 14–21. doi: 10.1128/mnbr.00002-09.

- 1118 World Health Organization (2022). *Global tuberculosis report 2022*. Section: xiii, 51 p. Geneva: World
1119 Health Organization.
- 1120 Wright S (1931). Evolution in mendelian populations. *Genetics* 16. doi: 10.4161/hv.21408.
- 1121 Yang Z (1998). Likelihood ratio tests for detecting positive selection and application to primate
1122 lysozyme evolution. *Molecular Biology and Evolution* 15, 568–573. doi: 10.1093/oxfordjournals.
1123 molbev.a025957.
- 1124 — (2007). PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution*
1125 24, 1586–1591. doi: 10.1093/molbev/msm088.
- 1126 Yang Z and JR Bielawski (2000). Statistical methods for detecting molecular adaptation. *Trends in*
1127 *Ecology and Evolution* 15, 496–503. doi: 10.1016/S0169-5347(00)01994-7.
- 1128 Yang Z and R Nielsent (2002). Codon-substitution models for detecting molecular adaptation
1129 at individual sites along specific lineages. *Molecular Biology and Evolution* 19, 908–917. doi:
1130 10.1093/oxfordjournals.molbev.a004148.
- 1131 Yang, Ziheng (2014). *Molecular Evolution – A Statistical Approach*. Oxford University Press.
- 1132 Zwyer M et al. (2021). A new nomenclature for the livestock-associated *Mycobacterium tuber-*
1133 *culosis* complex based on phylogenomics. *Open Research Europe* 1, 100–100. doi: 10.12688/
1134 openreseurope.14029.1.