

Mutation rates and adaptive variation among the clinically dominant clusters of *Mycobacterium abscessus*

Nicoletta Commins

Harvard Medical School

Mark Sullivan

Harvard TH Chan School of Public Health

Kerry McGowen

Harvard TH Chan School of Public Health

Evan Koch

Harvard Medical School

Eric Rubin

Harvard TH Chan School of Public Health <https://orcid.org/0000-0001-5120-962X>

Maha Farhat (✉ maha_farhat@hms.harvard.edu)

Harvard Medical School <https://orcid.org/0000-0002-3871-5760>

Article

Keywords:

Posted Date: May 18th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1620528/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Mutation rates and adaptive variation among the clinically dominant clusters of**
2 ***Mycobacterium abscessus***

3

4 Nicoletta Commins¹, Mark R. Sullivan², Kerry McGowen², Evan Koch¹, Eric J. Rubin^{2,3}, Maha
5 Farhat^{1,4*}

6

7 ¹ Department of Biomedical Informatics, Harvard Medical School, Boston, MA, 02115, USA.

8 ² Department of Immunology and Infectious Diseases, Harvard T. H. Chan School of Public
9 Health, Boston, MA, 02115, USA.

10 ³ Department of Microbiology, Harvard Medical School, Boston, MA, 02115, USA.

11 ⁴ Division of Pulmonary and Critical Care Medicine, Massachusetts General Hospital, Boston,
12 MA, 02114, USA.

13

14 * Correspondence: maha_farhat@hms.harvard.edu

15

16 **Abstract**

17 *Mycobacterium abscessus* (*Mab*) is a multi-drug resistant pathogen increasingly responsible for
18 severe pulmonary infections. Analysis of whole genome sequences (WGS) of *Mab* demonstrates
19 dense genetic clustering of clinical isolates collected from disparate geographic locations. This
20 has been interpreted as supporting patient-to-patient transmission, but epidemiological studies
21 have contradicted this interpretation. Here we present evidence for a slowing of the *Mab*
22 molecular clock rate coincident with the emergence of phylogenetic clusters. We find that
23 clustered isolates are enriched in mutations affecting DNA repair machinery and have lower
24 spontaneous mutation rates *in vitro*. We propose that *Mab* adaptation to the host environment
25 through variation in DNA repair genes affects the organism's mutation rate and that this manifests

26 as phylogenetic clustering. These results inform our understanding of niche switching for
27 facultative pathogens, and challenge the model of transmission as the major mode of
28 dissemination of clinically dominant *Mab* clusters.

29

30 **Introduction**

31 *Mycobacterium abscessus* (*Mab*) is an emerging, multi-drug resistant pathogen that is
32 increasingly responsible for opportunistic pulmonary disease, most commonly in patients with
33 underlying structural lung conditions such as cystic fibrosis (CF)¹. *Mab* is divided into three
34 subspecies: *M. a. abscessus*, *M. a. massiliense*, and *M. a. bolletii*, the former two being the most
35 clinically relevant subspecies². Infections with *Mab* are associated with accelerated decline in lung
36 function and can result in severe disseminated disease and/or loss of eligibility for lung
37 transplantation. *Mab* infections are difficult to treat, requiring a prolonged antibiotic regimen with
38 high rates of treatment failure³.

39 Until recently, *Mab* infections were thought to be acquired independently by each patient
40 from environmental reservoirs in soil or water. Large scale analyses of whole genome sequencing
41 (WGS) of *Mab* clinical isolates have revealed that a large proportion of clinical isolates fall into
42 several phylogenetically characterized clusters, often termed dominant circulating clones⁴.
43 Isolates within these clusters have high core genome similarity that is not explained by point
44 source outbreaks as they were sampled from diverse geographic origins. These observations
45 have been interpreted as supporting widespread recent transmission, either through direct contact
46 or indirectly through fomites, but a number of epidemiological studies have refuted person-to-
47 person transmission as a major mode of *Mab* dissemination⁵⁻¹¹. Reconciling these conflicting lines
48 of evidence is important as they dictate public health practice and allocation of resources to
49 protect vulnerable populations from infection. However, it is still unclear what role transmission
50 plays in the emergence of *Mab* clusters.

51 While widespread transmission is one proposed explanation for the phylogenetic structure
52 observed in *Mab*, it is important to consider other factors in the species' evolutionary history that
53 may contribute. For example, ecological niche switching from environmental sources to a human
54 host may introduce bottlenecks resulting in low diversity and effective population size. Another
55 explanation is a change in mutation rates, either through changes in bacterial generation time or
56 through adaptive genetic variation. A third consideration is sampling bias, where particular *Mab*
57 lineages that cause more frequent, severe and/or persistent disease are more likely to be isolated
58 and sequenced. Consistent with the latter hypothesis, previous work has shown that isolates
59 within clusters are more likely to contain point mutations associated with antibiotic resistance and
60 are associated with greater bacterial burden and lung inflammation in mice⁴.

61 Interpreting tree structures in the context of population history is aided by phylogenetic
62 dating, which translates genetic distances into time scale. Previous studies^{4,9,12,13} have estimated
63 molecular clock rates in *Mab* but these have been limited to subpopulations within clusters. Here
64 we report on molecular clock rate estimation for the deep branches of the *M. a. abscessus*
65 phylogeny and for within-host isolates belonging to clusters. We provide evidence for an
66 evolutionary slow-down coincident with the expansion of clinically dominant clusters and propose
67 that the emergence of dense phylogenetic clustering in *Mab* may result from genetically encoded
68 changes to the mutation rate. These results indicate that transmission is not the only plausible
69 explanation for the dense clustering observed in the *Mab* phylogeny, and that more research is
70 needed to clarify the transmission dynamics of this emerging pathogen.

71

72 **Results**

73 Phylogenetic analysis confirms the presence of large, dominant clusters (DCs) spanning multiple
74 geographies. We identified 1,629 *M. abscessus* isolates with publicly available whole genome
75 sequencing (WGS) data and associated dates of collection (Extended Data Table 1). We

76 excluded isolates with likely contamination or low-quality assemblies (Methods). Among the 1,461
77 isolates passing quality filters, the majority (80.1%) were derived from *in vitro* culture of pulmonary
78 samples. Of the 1,181 pulmonary isolates, 98.3% (1,161/1,181) were collected from patients with
79 cystic fibrosis (CF), while the rest were collected from patients with no documented diagnosis.
80 The dataset represented 483 patients with unique patient identifiers, with 59.6% (288/483) of
81 patients contributing a single isolate, and the remaining contributing between 2 to 24 isolates
82 sampled over time. Isolates originated from nine countries and were collected between 1998-
83 2017, with the exception of the reference strain ATCC 19977 which was collected in 1957 (Figure
84 1B).

85 The majority (91.7%) of samples had $\geq 98\%$ average nucleotide identity (ANI) with one of
86 three reference genomes (Supplementary Table 1) representing each of the three *Mab*
87 subspecies, consistent with the commonly used ANI threshold for bacterial subspecies. Isolates
88 with $< 98\%$ ANI with any subspecies were excluded from downstream analysis. The resulting
89 dataset overrepresented clinically relevant subspecies of *M. abscessus*, with 93.9% belonging to
90 either *M. a. abscessus* or *M. a. massiliense*. To better resolve true evolutionary relationships
91 among samples, we defined and restricted analysis to the core genome for each subspecies, then
92 further excluded predicted recombination events (Methods). We performed maximum likelihood
93 (ML) phylogenetic analysis for each of the three *Mab* subspecies, using only the first isolate from
94 each patient to focus on between-host diversity. The resulting subspecies trees confirmed the
95 presence of multiple, globally distributed clades with dense phylogenetic clustering (Methods).
96 We excluded clustering due to potential point source outbreaks by limiting to clusters containing
97 isolates from more than one country and a minimum of three isolates for downstream analysis. In
98 total, we confirmed 16 clusters (Extended Data Figure 1A) including four clusters with > 50
99 isolates, which we refer to here as dominant clusters (DCs). Isolates within clusters comprised
100 76.4% of isolates included in the phylogenetic analysis. The median pairwise distance in *M. a.*

101 *abscessus* DC 1 was 68 SNPs. For comparison, the mean pairwise distance across the
102 subspecies trees were 5537, 4100, and 8034 SNPs for *M. a. abscessus*, *M. a. massiliense*, and
103 *M. a. bolletii*, respectively. All four DCs contained isolates originating from three continents (North
104 America, Europe, and Australia). Clustered isolates also spanned the full range of isolation times,
105 including ATCC 19977, which was placed within *M. a. abscessus* DC 1 despite having been
106 collected at least 40 years before any other isolates in the tree (Figure 1A-B).

107

108 Coalescent analysis reveals a faster long-term evolutionary rate than observed in DCs. To aid in
109 transmission inference and to reconstruct the evolutionary history of *Mab*, we performed a
110 coalescent analysis using our SNP alignment. Previous work has reported molecular clock rates
111 within *Mab* clusters, but to our knowledge no study has successfully dated any of the full *Mab*
112 subspecies phylogenies. To evaluate the presence of temporal signal, we performed root-to-tip
113 regression for each subspecies tree but in each case did not observe a positive correlation
114 between sampling date and root-to-tip distance. We reasoned that the high degree of population
115 structure in the *Mab* phylogeny may obscure the genetic-temporal signal. To test this hypothesis,
116 we pruned each subspecies tree until it contained the minimal number of samples required to
117 capture 95% of the genetic diversity in the original tree (Methods). Following pruning of the *M. a.*
118 *abscessus* subspecies tree, root-to-tip regression demonstrated the presence of temporal signal
119 (Figure 1C, D). Evidence of temporal signal in *M. a. abscessus* was further strengthened using a
120 permutation test and a Bayesian evaluation of temporal signal (Extended Data Figure 1C,
121 Extended Data Table 3).

122 We used a Bayesian implementation of coalescent analysis, BEAST¹⁴, to estimate the
123 molecular clock rate of the pruned *M. a. abscessus* subspecies tree assuming a relaxed clock
124 model, which allows each branch of the tree to have its own clock rate. To explore the effect of
125 the tree prior on molecular clock estimation, we ran analyses using two different tree priors,
126 assuming either a constant coalescent model, which assumes constant population size over time,

127 or a Bayesian skyline coalescent model, which allows the population size to change over time.
128 While the Bayesian skyline coalescent had a slightly better marginal likelihood, we found no
129 significant difference between the clock rates estimated by the two models (Table 1, Figure 2C).
130 Our estimates for the clock rate in the long, internal branches of the tree were approximately 10-
131 fold faster than nearly all previously reported estimates of the clock rate within *Mab* clusters
132 (Figure 2C). We next attempted to date the full (unpruned) *M. a. abscessus* phylogeny using the
133 faster rate inferred from our coalescent analysis of the pruned tree. These chains suffered from
134 poor convergence and sampled trees with a distorted topology, elongating the short branches
135 and eliminating the densely clustered structure observed in the ML phylogeny (data not shown).
136 We therefore hypothesized that the lack of temporal signal in the full *M. a. abscessus* dataset
137 may result from a difference between the long-term clock rate and the more recent evolutionary
138 rate within phylogenetic clusters.

139

140 Analysis of longitudinally sampled within-host isolates supports a slower clock rate in DCs. To
141 further investigate the differences in clock rate between long-term historical time scales and short-
142 term contemporary time, we used pairs of isolates sampled from the same patient over time to
143 estimate the mutation rate. We focused on estimating a within-host clock rate for *M. a. abscessus*
144 to compare directly with the estimate from our coalescent analysis and because we had the most
145 longitudinal samples for this subspecies. Because polyclonal *Mab* infection is common in CF
146 patients, we first excluded any isolate pairs with a SNP distance > 20. This threshold was selected
147 by comparing the distribution of all possible within-host isolate pairs to the distribution of pairwise
148 SNP distances within *Mab* clusters (Figure 2A, Extended Data Figure 2). Regressing the SNP
149 distance between isolate pairs on the time between sample pairs yielded a mutation rate of 2.21
150 SNPs/year [1.06 – 3.37 95% CI], or 5.17×10^{-7} SNPs/site/year [2.48×10^{-7} – 7.92×10^{-7} 95% CI]. This
151 value is more comparable to previously reported molecular clock rate estimates within *Mab*
152 clusters^{4,9,12,13} than to those estimated in the long branches of the phylogeny (Figure 2C).

153
154 Simulated ancestries support mutation rate slow-down over population size effects as the major
155 driver of phylogenetic clustering. To test whether sampling from clades with low genetic diversity
156 is sufficient to explain the degree of clustering observed in the ML phylogeny, or whether changes
157 in the molecular clock rate are a better explanation for the observed phylogenetic data, we
158 performed coalescent simulations of ancestries with one clustered subpopulation (A) and one
159 unclustered subpopulation (B). Under a range of effective population sizes (N_e) and mutation rates
160 (μ), we simulated ancestries through time with 1000 replicates per simulation. For each replicate
161 we assessed the degree of phylogenetic clustering by measuring d_A/d_B where d_i is the mean
162 pairwise SNP distance within subpopulation i . The total N_e (N_e^T), number of samples, and
163 sampling dates for subpopulations A and B were drawn from the true values for *M. a. abscessus*
164 DC 1 and from all “unclustered” isolates in our dataset, respectively (Figure 3A). We first tested
165 whether a smaller population size for the clustered subpopulation A (N_e^A) could explain the degree
166 of phylogenetic clustering observed. We estimated a range of realistic values for N_e^A using BEAST
167 by assuming a range of possible clock rates (Extended Data Table 4). These rates spanned a
168 published clock rate estimate for *M. a. abscessus* DC1 to a clock rate that is 50-fold higher, far
169 exceeding the upper 95% highest posterior density (HPD) interval for the published value¹³ Under
170 the the simulated conditions of equal mutation rates between subpopulation A and B, we only
171 observed phylogenetic clustering to the extent observed in the clinical isolate phylogeny in the
172 scenario where $N_e^T \gg N_e^A$ (Figure 3B, Extended Data Figure 3). We next tested whether a change
173 in the mutation rate at the root of supopulation A could better explain the observed phylogenetic
174 structure in *Mab*. The simulated ancestries where N_e^A and N_e^B are fixed and the mutation rate
175 slowdown is between 10- and 20-fold recapitulated the observed phylogenetic structure (Figure
176 3C). This result was robust to changes in N_e^A and N_e^T across the 95% HPD for each parameter
177 (Extended Data Figure 4). These results indicate that for person-to-person transmission to drive
178 phylogenetic clustering in *Mab*, the total N_e must be very large compared to the N_e of the clustered

179 subpopulation, an unlikely scenario given the observed genetic data and sampling times. By
180 contrast, a mutation rate change of 10- to 20- fold explains the degree of clustering observed in
181 *Mab* under more realistic population sizes given these data.

182

183 UvrD/Rep family helicases are enriched in DCs in *M. a. abscessus*. To identify genetic variants
184 under positive selection in *Mab* clusters, we searched for evidence of phylogenetic convergence
185 in the core genome of *M. a. abscessus* with regions of recombination excluded. Briefly, we
186 performed ancestral sequence reconstruction for the internal nodes of the subspecies tree and
187 counted the number of independent arisals of each SNP (Methods). Of the 65,202 SNPs identified
188 in *M. a. abscessus*, 34,813 (53.4%) of mutations occurred only once and were excluded. For the
189 remaining SNPs, we ranked each by its enrichment in DCs relative to other regions of the *M. a.*
190 *abscessus* tree (Figure 4A). The most significantly enriched variants in the *M. a. abscessus* DCs
191 included homologs of several known or putative virulence factors as well as a number of cluster-
192 enriched mutations in UvrD-like helicases *MAB_1054*, *MAB_3515c*, and *MAB_3516c* (Table 2,
193 Figure 4A). While the function of these genes is not well characterized in *M. abscessus*, all three
194 share homology with the UvrD/Rep family of helicases. This family of genes is involved in DNA
195 repair and has been associated with growth and persistence *in vivo*^{15,16}.

196

197 Isolates in *Mab* DCs are associated with a lower spontaneous mutation rate. To test the
198 association between phylogenetic clustering and spontaneous mutation rate, we selected *M. a.*
199 *abscessus* and *M. a. massiliense* strains that were grouped within DCs or that were not present
200 in those clusters. Each of the strains present in the DCs has the derived (minor) allele at four sites
201 in the UvrD/Rep family genes *MAB_1054*, *MAB_3515c*, and *MAB_3516c*, and the strains outside
202 of DCs possess the ancestral (major) allele at those sites (Table 2, Figure 4A). The strains used
203 and their genotypes are described in Supplementary Data Tables 2 and 3. The spontaneous
204 mutation rate was estimated in all strains using a fluctuation assay¹⁷ adapted for use in *M.*

205 *abscessus*. The fluctuation assay counts cells gaining resistance to the antibiotic amikacin during
206 growth in culture without antibiotic. By normalizing the number of mutational events to the
207 population size of the culture, the fluctuation assay estimates the spontaneous mutation rate per
208 generation per the number of bases that can be mutated to confer resistance to amikacin. The
209 strains chosen display similar growth rates (Extended Data Figure 5) and very low baseline
210 resistance to amikacin (Extended Data Table 5) which allows them to be compared using a
211 fluctuation assay. Among *M. a. abscessus* strains, those present in DCs with the derived (minor)
212 allele at the four UvrD/Rep loci had lower spontaneous mutation rates than those with the
213 ancestral (major) allele (Figure 4C). A similar trend was observed in *M. a. massiliense* strains
214 (Figure 4D). Further, the average mutation rate was significantly lower in isolates present in DCs
215 across both subspecies (Figure 4E), consistent with the hypothesis that the emergence of dense
216 clustering may be driven in part by changes in mutation rate.

217

218 **Discussion**

219 In this study we have pooled a large global sample of *Mab* genome sequences and
220 examined evolutionary rates and positive selection with the goal of understanding the emergence
221 of *Mab* clusters. We present evidence supporting a slower evolutionary rate within DCs compared
222 to the long internal branches of the phylogeny and further data associating adaptive genetic
223 variation in DNA repair genes with the slower evolutionary rate. We posit that changes to the
224 mutation rate underlie the dense clustering in *Mab*, and that while other effects such as population
225 or transmission bottlenecks and sampling bias may contribute to the low genetic diversity within
226 clusters, they are not sufficient to explain the observed phylogenetic structure.

227 Previous studies have restricted molecular clock rate estimation to *Mab* clusters, most
228 likely because of the lack of temporal signal at the subspecies-level. By pruning the *Mab* tree
229 while preserving the majority of the subspecies diversity, we were able to measure a molecular
230 clock rate for *Mab* spanning the time scale from the common ancestor of the subspecies to the

231 present. Estimates of evolutionary rates are recognized to depend on the time scale of
232 measurement. In such examples, the clock rate usually decreases with increasing time depth¹⁸.
233 We observe the opposite effect where the clock rate increases with time depth. The pruned *Mab*
234 tree most likely reflects an evolutionary process taking place predominantly in environmental
235 reservoirs, in contrast to phylogenetic clusters where terminal branches reflect evolution within
236 the human lung environment. We also note that our measurement of within-host mutation rate in
237 *Mab* was independent of phylogenetic inference and the full temporal signal in our sample.
238 Nevertheless, the measured rate within-host was comparable to the clock rate measured
239 previously using coalescent analysis within clusters (Figure 2C).

240 The molecular clock rate is influenced by both the spontaneous mutation rate and the
241 generation time of an organism. Changes in generation time could result from physiological
242 constraints of the human lung environment, if for example, differences in nutrient availability
243 change the growth rate in the lung compared to those living in environmental reservoirs.
244 Generation times and spontaneous mutation rates can also be modified by genetic changes, for
245 example to DNA replication or repair machinery. The observation that some lineages have
246 adapted to the host environment⁴ raises the possibility that genetic changes that underlie or
247 associate with this adaptation affect mutation rate, either directly by affecting DNA replication or
248 repair machinery, or indirectly by affecting the growth rate. Ruis et al. reported differences in
249 mutational signatures along internal branches of the subspecies phylogeny compared to internal
250 branches of DCs, supporting a change in the mutational process within some clades of the tree¹³.

251 We scanned the core genomes for variants enriched in the DCs in *M. a. abscessus*, the
252 subspecies for which we have the most WGS data and for which there are the highest number of
253 clusters. Notably, we observed several cluster-enriched mutations in UvrD-like helicases. UvrD is
254 a superfamily 1 helicase that is involved in DNA repair and is widely conserved in gram-negative
255 bacteria. UvrD1 in *M. tuberculosis* plays an important role both in nucleotide excision repair (NER)
256 and in pathogenesis and persistence because its DNA repair activity confers tolerance to reactive

257 oxygen intermediates (ROI) and reactive nitrogen intermediates (RNI)¹⁵. Mutations that promote
258 the activity of UvrD1 and its homologs are likely advantageous in the context of lung infection and
259 may have the additional consequence of decreasing the spontaneous mutation rate. We
260 additionally observed mutations in UvrD-like helicases occurring along the cluster-defining branch
261 of the large *M. a. massiliense* cluster, further supporting the role of the NER pathway in *Mab*
262 pathogenesis.

263 Our study has several limitations. Our dataset overrepresented clinical isolates from the
264 UK where routine whole genome sequencing is part of *Mab* surveillance. Oversampling of some
265 lineages and geographies may mean that we have not adequately captured the global population
266 structure and that we may have focused on large phylogenetic clusters because of their
267 geographic distribution and not because of their clinical significance. While our study argues
268 against transmission as the main driver of *Mab* dissemination and the emergence of DCs, we do
269 not examine the geographic distribution of *Mab* within clusters. It is possible that future work
270 considering this data can highlight mechanisms of parallel evolution that explain why clustered
271 *Mab* isolates are detected in distant geographic locations. While our results support a faster
272 molecular clock rate in the long internal branches of the phylogeny compared to clustered
273 lineages, we are unable to distinguish between rate differences due to the spontaneous mutation
274 rate from generation time effects. Differences in growth conditions across environments is difficult
275 to quantify as *Mab* is ubiquitous in both natural and built environments, but nutrient availability
276 likely differs greatly between natural reservoirs such as soil and water compared to man-made
277 reservoirs such as showerheads. A healthy lung supports relatively little bacterial growth, but
278 inflammation increases nutrient availability by promoting mucus production and vascular leakage
279 into the airway ¹⁹. Because there is little publicly available sequencing data from environmental
280 isolates, we were also unable to compare genotypes or mutation rates in environmental samples
281 compared to clinical isolates. Our analysis of variants associated with phylogenetic clustering is

282 also limited in its ability to identify true causal genotypes or those that are enriched in clusters due
283 to hitchhiking to other adaptive variants.

284 In summary, we present new evidence for changes in the molecular clock rate that
285 contribute to the dense phylogenetic clustering observed in *Mab*. These results argue against
286 person-to-person transmission as the primary driver of clustering in the *Mab* phylogeny.
287 Continuing genomic surveillance is needed to understand how *Mab* may be adapting to human
288 infection and to characterize the risk of transmission to at-risk patients. We argue further that the
289 phylogeographic structure of *Mab* is influenced by its complex ecological history as an emerging
290 facultative pathogen that should be considered carefully in future analysis of *Mab* genomic data.
291 Future work should extend these analyses by examining differences in mutation rates and
292 phenotypes between clinical isolates and environmental strains for evidence of host adaptation.

293

294 **Acknowledgements**

295 We thank members of the Farhat lab for feedback and discussion. We are grateful to Luca Freschi
296 for advice on processing WGS data, calling variants, and phylogenetic inference; Maximillian
297 Marin for input on calculating mappability pileup; and Roger Vargas for support in homoplasy
298 counting. We also thank Wilder Wohns for discussions on implementing ancestry simulations with
299 mutation rate changes through time. M.R.S. is a Merck Fellow of the Damon Runyon Cancer
300 Research Foundation, DRG-2415-20. Computational resources and support were provided by the
301 Orchestra High Performance Compute Cluster at Harvard Medical School, which is funded by the
302 NIH (NCRR 1S10RR028832-01).

303

304 **Author Contributions**

305 Conceptualization, N.C. and M.F.; Methodology – Computational analysis, N.C., E.K., and M.F.;
306 Methodology – Experimental studies, M.R.S., K.M.; Investigation – Computational analysis, N.C.;

307 Investigation – Experimental studies, M.R.S., K.M.; Data Curation, N.C.; Writing – Original Draft,
308 N.C. and M.R.S.; Writing – Review & Editing, N.C., M.R.S., E.K., M.F.; Supervision, E.J.R., M.F.

309

310 **Declaration of Interests**

311 The authors declare no competing interests.

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333 **Tables**

334

335 **Table 1.** Coalescent analysis results and model comparison

Model	Marginal Log Likelihood (SD)	UCLD Rate (SNPs/site/year)	UCLD St. Dev.
Coalescent constant population	-7144301.64 (21.42)	2.82×10^{-6} [1.33×10^{-6} , 4.83×10^{-6}]	0.182 [0.135, 0.234]
Bayesian skyline coalescent	-7144118.76 (19.93)	2.63×10^{-6} [1.16×10^{-6} , 4.0×10^{-6}]	0.181 [0.138, 0.232]

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422

Table 2. Summary of top 10 nonsynonymous variants enriched in DCs of *M. a. abscessus*

Gene	Gene Symbol	Description	Comments	Enrichment p-value	Number of clustered taxa containing variant	Number of unclustered taxa containing variant
MAB_3334c	<i>gatB</i> ^c	Aspartyl/glutamyl-tRNA (Asn/Gln) amidotransferase subunit B ^{c,d}	Essential for growth and survival <i>in vivo</i> in <i>M. tuberculosis</i> ¹⁷	1.03x10 ⁻⁶⁵	240	1
MAB_3244		Mercury resistance transport protein/cytochrome c biogenesis protein ^c		1.72x10 ⁻⁶⁴	189	3
MAB_3480		Putative short chain dehydrogenase/ Reductase ^c		7.23x10 ⁻⁶⁴	232	13
MAB_4141 ^a		PE/PPE family protein ^c		1.58x10 ⁻⁶²	189	2
MAB_3242	<i>fni</i> ^c	Isopentenyl-diphosphate delta-isomerase ^{c,d}		8.07x10 ⁻⁶²	189	3
MAB_3481	<i>fadE</i> ^c	Probable acyl-CoA dehydrogenase ^c	FadE mutations in <i>Mtb</i> result in decreased growth and virulence in macrophages and in mice ¹⁸	2.10x10 ⁻⁶¹	232	15
MAB_4057c	<i>mshA</i> ^c	D-inositol 3-phosphate glycosyltransferase ^d	Catalyzes biosynthesis of mycothiol, an important compound in resisting oxidative stress and detoxifying harmful agents such as antibiotics ¹⁹	2.78x10 ⁻⁶⁰	189	2
MAB_3581c		Hypothetical protein		6.09x10 ⁻⁶⁰	169	3
MAB_1054 ^b	<i>pcrA</i> ^c / <i>uvrD1</i> ^d	ATP-dependent DNA helicase ^{c,d}	Involved in nucleotide excision repair (NER), mismatch repair (MMR), homologous recombination (HR), and rolling circle replication ¹⁵	2.72x10 ⁻⁵⁸	216	7
MAB_3515c	<i>rep</i> ^{c,d}	ATP-dependent DNA helicase ^{c,d}	Involved in NER, MMR, HR and rolling circle replication ¹⁵	8.42x10 ⁻⁵⁷	189	7

423 **Table 2 Legend:**

424 ^a Multiple variants in *MAB_4141* have been collapsed and the lowest p-value is shown. See
425 Supplementary Material for full list of homoplastic SNPs ranked by their enrichment in clusters.

426 ^b Gram-negative bacteria possess UvrD and a closely related homolog, Rep, while gram-positive
427 bacteria possess a single homolog called PcrA.²⁰

428 ^c Mycobrowser annotation

429 ^d prokka annotation

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

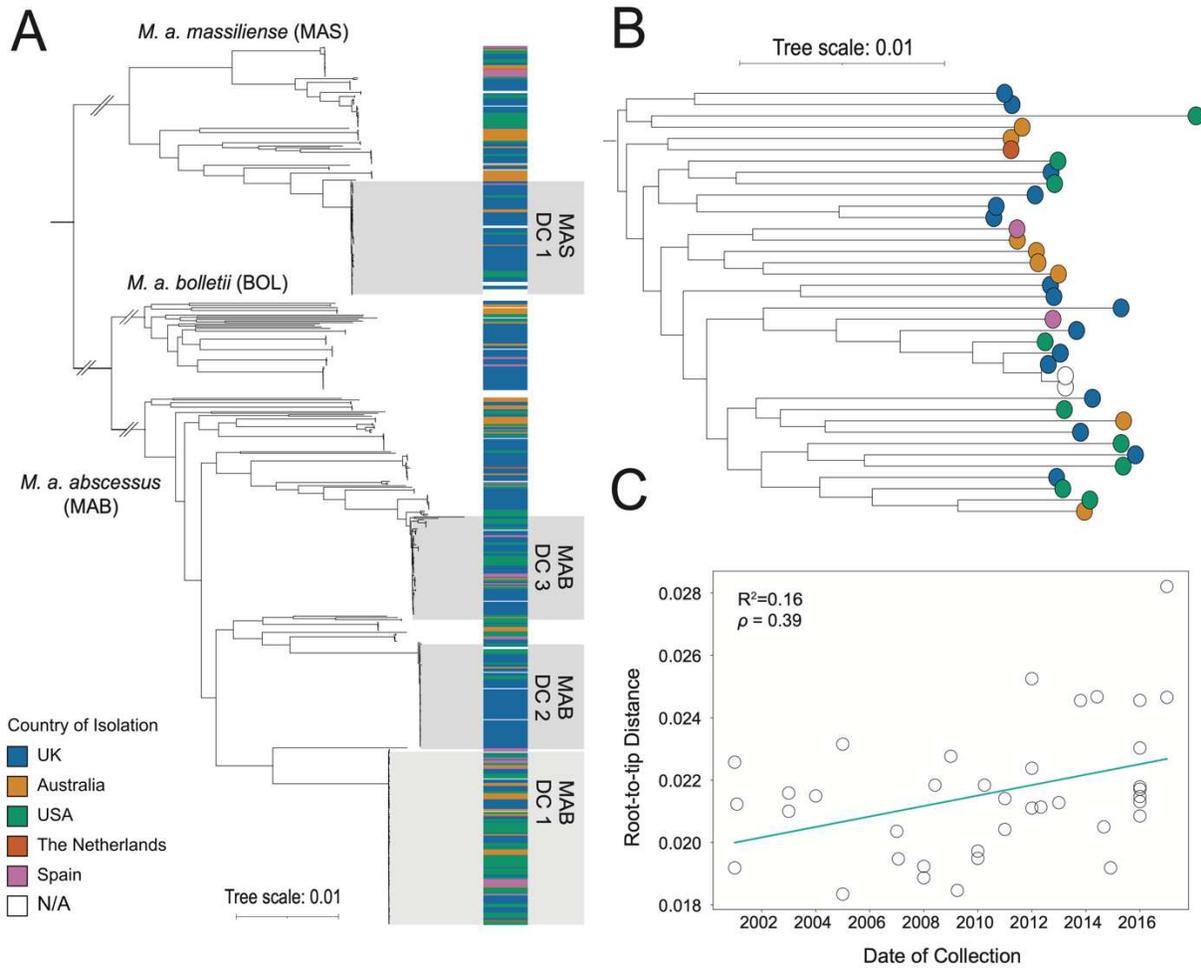
448

449

450

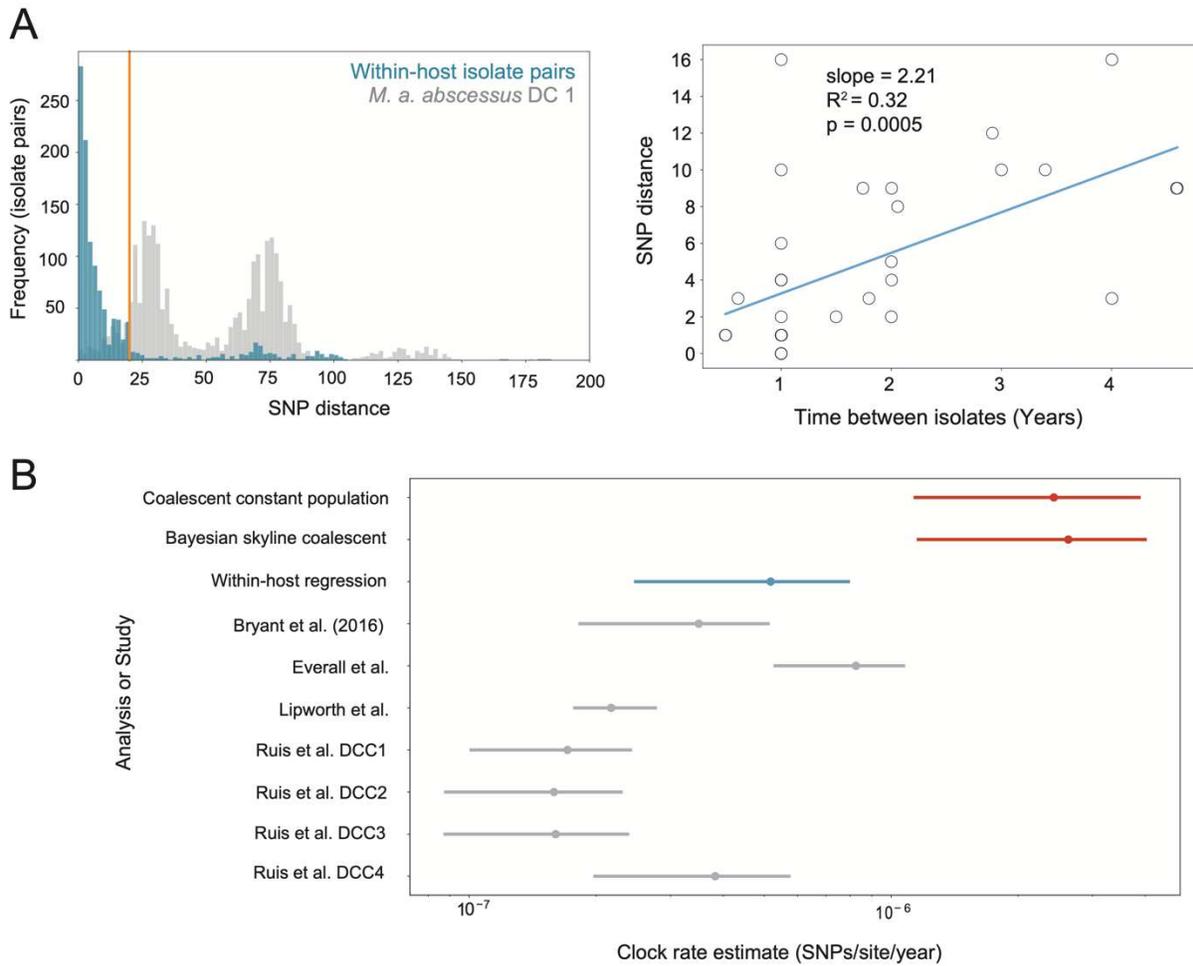
451

452 **Figures**



453

454 **Figure 1. Dominant clusters of *M. abscessus* and the presence of temporal signal in**
 455 **internal branches.** A) Species-wide phylogeny of *M. abscessus* representing one isolate per
 456 patient to avoid within-patient clustering of samples. Each subspecies tree was constructed
 457 independently using a subspecies specific reference genome. The four largest phylogenetic
 458 clusters are highlighted. B) Histogram showing the distribution of collection dates of all clustered
 459 isolates v. all unclustered isolates in *M. a. abscessus*. C) Pruned *M. a. abscessus* phylogeny
 460 representing 95% of the original subspecies tree diversity. C) Root-to-tip regression showing
 461 evidence of positive temporal signal in the pruned tree shown in (B).



462

463 **Figure 2. Longitudinally sampled within-host isolates and coalescent analysis support a**
 464 **slower clock rate within DCs.** A) Clock rate estimation using isolate pairs sampled from the
 465 same patient over time. Left: distribution of SNP distances between all possible within- host isolate
 466 pairs compared to SNP distances between all possible isolate pairs within *M. a. abscessus* DC 1.
 467 The orange line represents the threshold of 20 SNPs used to exclude pairs representing
 468 independent infections. Right: regression of SNP distance on the time between isolates for each
 469 isolate pair. The slope of the regression line was 2.51 SNPs/year [1.24 – 3.78 95% CI], or roughly
 470 5.93×10^{-7} SNPs/site/year [2.92×10^{-7} – 8.9×10^{-7} 95% CI]. B) UCLD clock rate estimates for the
 471 three models used in this study (two coalescent models shown in red and one regression model
 472 shown in blue) compared to mutation rate estimates from the within-host regression (blue) and
 473 published clock rates (grey).

474

475

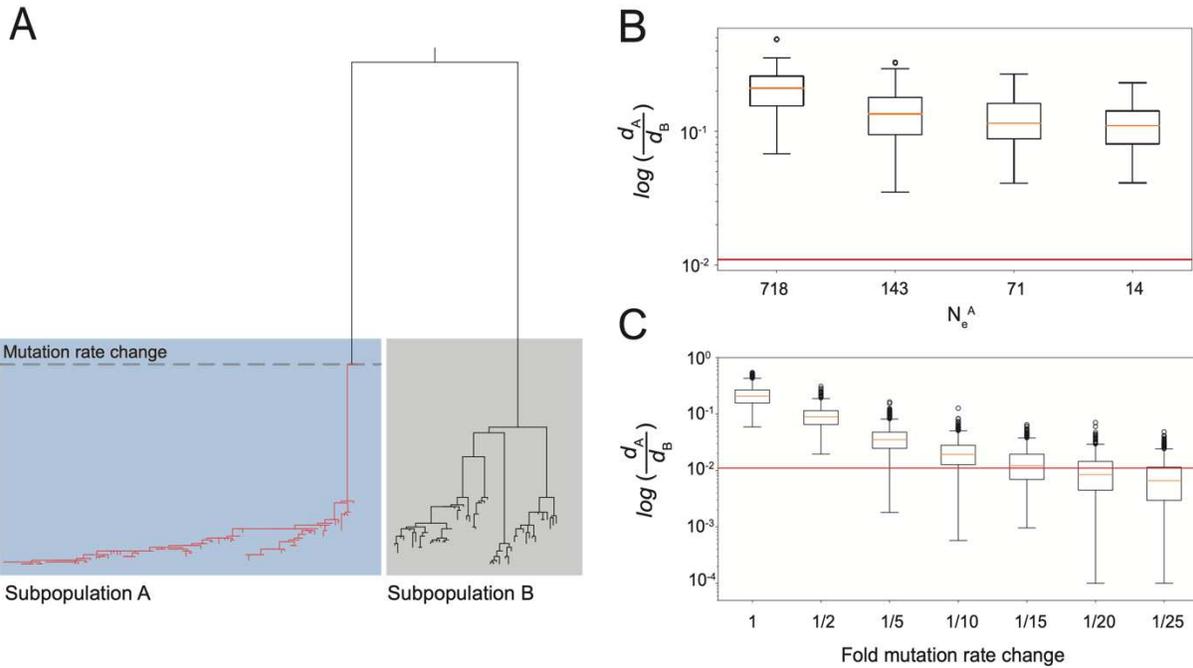
476

477

478

479

480



481

482 **Figure 3. Simulated ancestries support a mutation rate slow-down coincident with DC**
 483 **emergence.** A) Example neighbor joining (NJ) tree from the simulation shown in B. The number
 484 of samples and sampling dates for subpopulations A and B are drawn from real data from *M. a.*
 485 *abscessus* cluster 1 and from all unclustered samples, respectively. The dotted line crosses the
 486 common ancestor node for subpopulation A, and all branches descended from this node (shown
 487 here in red) are subject to the change in mutation rate. B) Boxplot showing the degree of
 488 phylogenetic clustering in subpopulation A relative to subpopulation B over a range of effective
 489 population sizes of subpopulation A (N_e^A). The simulation assumes that $N_e^T=8000$. C) Boxplot
 490 showing the degree of phylogenetic clustering in subpopulation A relative to subpopulation B over
 491 a range of mutation rate changes. The simulation assumes that $N_e^T=8000$ and $N_e^A=718$. In both
 492 A) and B) the degree of clustering is estimated with metric $\log(\frac{d_A}{d_B})$ where d_i is the mean pairwise
 493 SNP distance among all sample pairs in subpopulation i . The red line represents the clustering
 494 metric estimated using the observed ML phylogeny.
 495

496

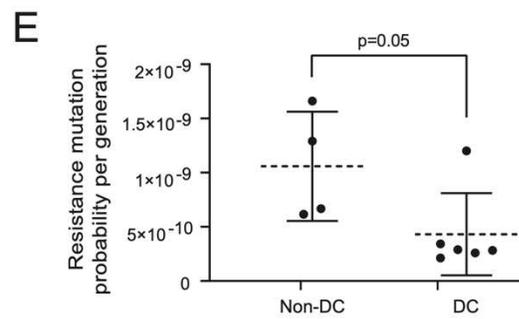
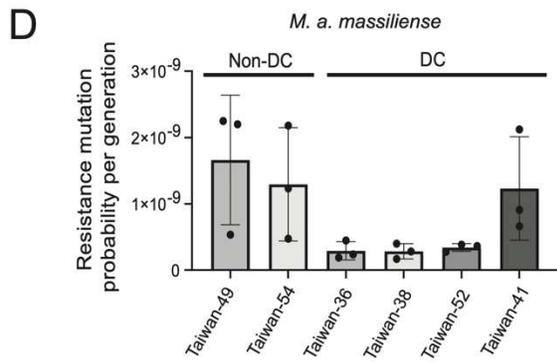
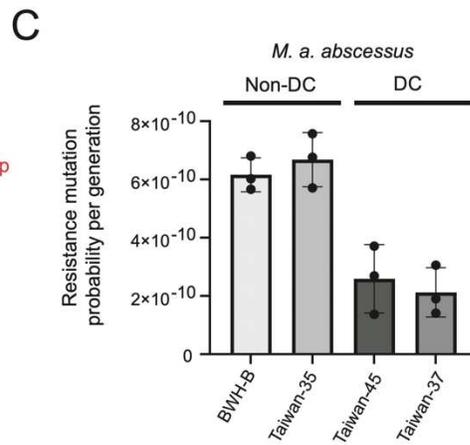
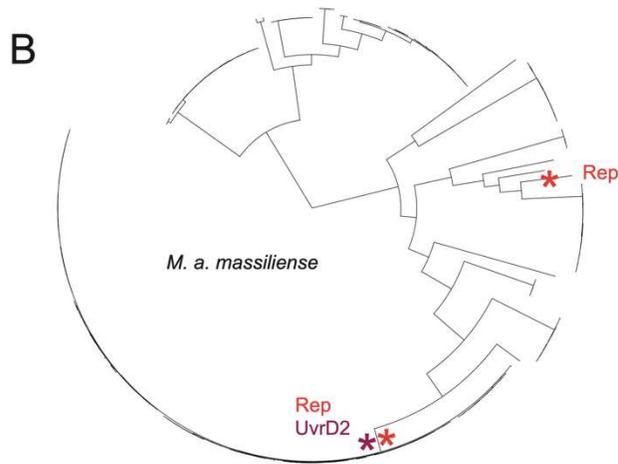
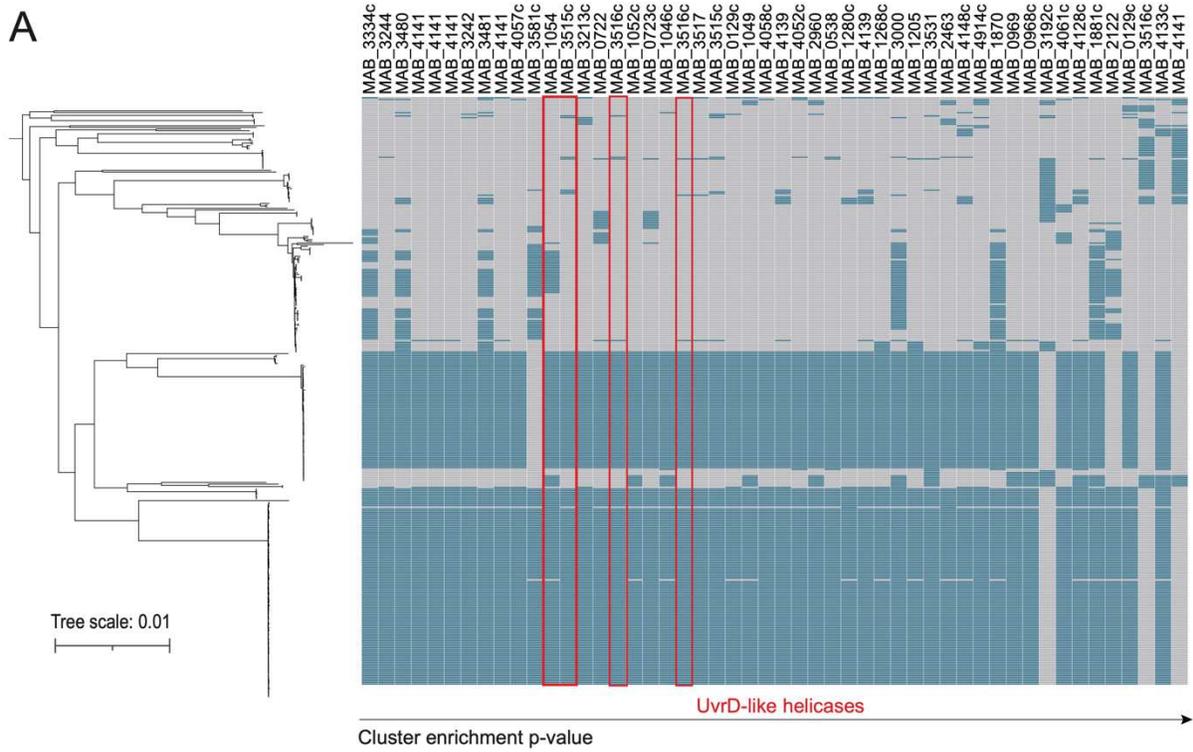
497

498

499

500

501



503 **Figure 4. Dominant clusters are associated with a decrease in the spontaneous mutation**
504 **rate *in vitro*.** A) Subspecies phylogeny for *M. a. abscessus* (left) and corresponding heatmap
505 representing the presence of the derived allele (blue) or ancestral allele (grey) for each
506 nonsynonymous variant. Variants are ranked in ascending order by their p-value derived from
507 Fisher's tests for enrichment in the three largest clusters. The top 50 hits are shown. The genes
508 where each variant is located are denoted. B) Subspecies phylogeny for *M. a. massiliense*. The
509 branches where mutations have occurred in homologs of Rep and UvrD2 are shown with
510 asterisks. Probability of acquiring an amikacin resistance mutation per generation in C) *M. a.*
511 *abscessus* strains, D) *M. a. massiliense* strains, or D) in all strains tested that are part of a
512 dominant cluster or not. Mean +/- SD is displayed. In C-D, n=3 biological replicates. In E, each
513 point represents the average mutation probability for a single strain. p-value derived from two-
514 tailed, unpaired t-test.
515

516

517 **Methods**

518 Data Sources and Description: WGS data were curated from published studies^{4,7,10,12,21-23} and
519 additional data were downloaded from the NCBI Sequence Read Archive (SRA). A search of
520 "abscessus" was conducted on October 4, 2019 using the filters "DNA" and "Illumina". Sequences
521 were selected that had an associated date of collection with the year at minimum. Collection dates
522 were obtained either from the literature, SRA metadata table, or provided by the authors. A
523 summary of all samples used in this study is provided in Extended Data Table 1. Raw sequencing
524 reads were downloaded using the NCBI SRA toolkit²⁴.

525

526 Quality control of Illumina reads: Reads were trimmed and filtered with trimmomatic²⁵. To detect
527 contaminated samples, the mean GC content of the reads were compared to a modelled normal
528 distribution of GC content using fastqc²⁶. If the sum of the deviations from the normal distribution
529 represented more than 30% of the reads, the sequencing run was excluded from further analysis.

530 Subspecies assignment and genome assembly: Multiple sequencing runs corresponding to one
531 biological sample were combined into one file. Reads were then assembled *de novo* using
532 SPAdes²⁷. Assemblies were compared to a reference sequence for each of the three subspecies
533 (Supplementary Data Table 1) and a whole-genome based average nucleotide identity (gANI)
534 score was calculated using fastANI²⁸. The resulting alignments were assigned to a subspecies

535 based on having gANI of at least 98% with exactly one reference strain. Isolates that did not have
536 at least 98% ANI with any reference strain were excluded. No isolates matched more than one
537 reference genome. After assigning isolates to a subspecies, reads were mapped to the
538 corresponding reference genome using BWA MEM²⁹. After alignment, isolates with <80%
539 coverage of the reference genome at a depth of at least 20x or with missing data at >30% of sites
540 were excluded.

541

542 Variant calling and SNP filtering: Variants were called using Pilon³⁰. We excluded ambiguous calls
543 as well as indels and MNVs. We also excluded calls with low coverage using a minimum depth
544 of either 10% of the mean coverage, or 5, whichever was greater (default settings for Pilon –
545 mindepth). We used bcftools³¹ to exclude calls with mapping quality or base quality scores < 20.

546

547 Core genome inference: Unlike the professional mycobacterial pathogens *M. tuberculosis* or *M.*
548 *leprae*, which have small accessory genomes, *M. abscessus* has a large, open accessory
549 genome³². To define the core genome for phylogenetic analysis and to exclude loci that are highly
550 divergent from the reference sequence, we masked all loci for downstream analysis where read
551 depth was <20x in more than 5% of the isolates in our sample set. The remaining loci were
552 compared to an estimate of the core genome calculated by Roary³³, and the two methods yielded
553 nearly identical results. We further excluded regions with average mapping quality scores or
554 average based quality scores <20. Finally, we used GenMap³⁴ to calculate mappability scores
555 across each reference genome as defined in reference³⁵. We then calculated a mappability pileup
556 score for each base position as described in reference³⁶ and excluded loci where the mappability
557 pileup score was <95% from downstream analysis. In total, we excluded 14.4%, 13.5%, and
558 15.1% of the *M. a. abscessus*, *M. a. massiliense*, and *M. a. bolletii* reference genome length,
559 respectively.

560

561 Phylogenetic tree inference: A full-length sequence alignment was created using a custom script,
562 masking loci excluded by our filtering criteria described above. Gubbins³⁷ was used to predict and
563 mask regions of variation produced by recombination from the full-length alignment. A
564 recombination-free SNP alignment generated by Gubbins was then used to build a phylogeny for
565 each subspecies using IQ-TREE³⁸, using the ModelFinder Plus (-mfp) option and 1000 bootstrap
566 replicates. Clusters were defined as previously described⁴ using a one-dimensional scanning
567 statistic to identify clades with a distribution of branch lengths that is shorter than the distribution
568 of branch lengths in the rest of the tree³⁹.

569

570 Coalescent analysis: Temporal signal was assessed using root-to-tip regression. In both *M. a.*
571 *abscessus* and *M. a. massiliense* subspecies trees and all three major clusters, a slightly negative
572 slope was observed, indicated lack of temporal signal. We recognized that the high degree of
573 clustering and short time span over which the clustered isolates were sampled may confound the
574 genetic temporal signal. To address this, we pruned the trees down to taxa that preserved 95%
575 of the relative tree length (RTL) of the original trees using Treemmer⁴⁰. A threshold of 95%
576 eliminated most of the dense clustering, but preserved the long branches of the original tree
577 (Figure 2B, Figure S2B). Using the pruned tree, we observed evidence of temporal signal in the
578 MAB subspecies using root-to-tip regression. We generated a SNP alignment including only the
579 taxa present in the pruned *M. a. abscessus* tree and used this as input for coalescent analysis
580 using BEAST2¹⁴. ModelTest-NG⁴¹ was used to select the site model GTR+4. MegaX⁴² was used
581 to perform a maximum likelihood test to determine that a relaxed clock model is best suited for
582 our data. This choice was confirmed using the nested sampling algorithm to compare marginal
583 likelihoods for a strict v. relaxed clock model. We then used nested sampling to compare various
584 tree priors. We using Bayesian Evaluation of Temporal Signal (BETS) to confirm the presence of
585 temporal signal by comparing marginal likelihoods when the tip dates were used to those when
586 all samples were assumed to be isochronous (Extended Data Table 3). Based on marginal

587 likelihood values, we ran both a coalescent constant population model and a Bayesian skyline
588 coalescent model, assuming a relaxed clock and including the sampling dates for both models.
589 We used the ML phylogeny of the pruned dataset as a starting tree. For each model we ran two
590 independent chains for at least 100 million states to ensure convergence.

591

592 Phylogenetic dating: To attempt to date the full MAB phylogeny, we ran both coalescent constant
593 population and Bayesian skyline coalescent models assuming a relaxed clock without sampling
594 dates. We set the mean UCLD clock rate to the values estimated from our coalescent analysis of
595 the pruned tree. We ran each chain for at least 100 million states.

596

597 Estimation of mutation rate within-host: We identified 90 patients with two or more *M. a. abscessus*
598 isolates. Because coinfection with multiple *M. abscessus* clones is common in cystic fibrosis
599 patients, we first sought to exclude any isolate pairs that represented distinct clonal infections.
600 We selected a threshold of 20 SNPs by comparing the distribution of all possible within-host
601 isolate pairs to the distribution of pairwise SNP distances within *Mab* clusters. We further excluded
602 any isolates where >20% of the genome is predicted to be rearranged by recombination. For the
603 remaining isolate pairs we regressed the pairwise SNP distance between the first and last isolates
604 sampled from each patient on the time between the collection of each sample. The slope of the
605 regression line was used as an estimate of the number of SNPs accumulated over time.

606

607 Ancestry simulations: We used msprime⁴³ to simulate ancestries to model two subpopulations:
608 (1) a “clustered” subpopulation (A) and (2) an “unclustered” subpopulation (B). We assumed a
609 constant population size and initially set the total effective population size (N_e) to 8000, the N_e
610 estimated in our coalescent analysis of the pruned *Mab* tree. We chose a range of realistic N_e for
611 subpopulation A, N_e^A by running a separate coalescent analysis of *M. a. abscessus* cluster 1 using
612 BEAST, assuming a constant population size, relaxed clock, and a range of UCLD clock rates

613 spanning 1.7×10^{-7} SNPs/site/year, as previously reported for this cluster by Ruis et al. (2021) to
614 8.5×10^{-6} SNPs/site/year, a value well above the 95% HPD interval for our estimate of the clock
615 rate in the pruned tree (Extended Data Table 4). We assumed the N_e of subpopulation B to be:
616 $N_e^B = N_e^{total} - N_e^A$. The number of samples and sampling dates supplied to the model for
617 subpopulations A and B are drawn from the true values for *Mab* cluster 1 and from all
618 “unclustered” isolates in the tree, respectively. We assumed an initial mutation rate of 1×10^{-9}
619 mutations/generation with a sequence length of 5Mb and a divergence time of 30,000 generations
620 between subpopulations A and B. We also assumed the nucleotide substitution model estimated
621 using BEAST.

622 We first simulated a scenario where the mutation rate remains constant, varying N_e^A over
623 the range shown in Extended Data Table 4. For each simulated ancestry we overlaid mutation
624 events, supplying the model with a starting substitution rate of 1×10^{-9} mutations/generation, a
625 nucleotide substitution model estimated using BEAST, and a sequence length of 5Mb. We then
626 extracted the simulated genotypes and constructed a neighbor joining (NJ) tree. The NJ tree was
627 used to exclude any simulations where the two subpopulations did not exhibit reciprocal
628 monophyly. We used the extracted genotypes to estimate the degree of clustering in each tree
629 by comparing the mean pairwise SNP distance in subpopulation A (d_A) to the mean pairwise SNP
630 distance in subpopulation B (d_B). For each simulated condition we performed 1000 replicates. To
631 simulate ancestries with a mutation rate change, we repeated the above procedure conservatively
632 assuming $N_e^A = 718$ (Extended Data Table 4). We modeled a per-generation mutation rate change
633 starting at the common ancestor node for subpopulation A and ranging from 1-fold (no change)
634 to 1/25-fold spanning a range of reasonable rate changes based on the range of clock rate
635 estimates observed in this study and reported in the literature (Figure 2B). For all simulated
636 conditions, we the simulation over a range of values for N_e^{total} spanning the 95% HPD intervals for
637 the N_e estimates using the Treemmer tree under both the constant coalescent and coalescent
638 skyline tree priors.

639

640 Scan for variants enriched in host-adapted lineages: To search for variants associated with host-
641 adapted MAB lineages, we first inferred the number of independent arisals in the MAB phylogeny
642 for each SNP using SNPPar⁴⁴. Briefly, SNPPar uses TreeTime to perform an ancestral sequence
643 reconstruction at each node, then infers mutation events arising on each branch of the tree. We
644 quantified the number of independent mutation events for each SNP. We considered only
645 mutation events where the derived call was the minor allele and the ancestral call was the major
646 allele. Because of the extreme phylogenetic clustering in *Mab*, we defined the major and minor
647 alleles using the pruned alignment generated by Treemmer. For all homoplastic SNPs (occurring
648 two or more times independently in the tree), we used Fisher's exact test to quantify the
649 enrichment of that variant within phylogenetic clusters. We conducted two separate analyses: a)
650 focusing on the three largest clusters in the *M. a. abscessus* tree (described in the main text), and
651 b) considering all clusters identified by this approach. In each case we obtained similar results.
652 For both analyses, full lists of variants ranked by their enrichment scores are provided in
653 Supplementary Data. We also used SNPPar to identify all mutation events occurring on the
654 branch ancestral to *M. a. massiliense* cluster 1.

655

656 *M. abscessus* culture: *M. abscessus* clinical isolates used for fluctuation assays were isolated
657 from patients in Taiwan and in Boston, MA, USA. All *M. abscessus* strains were cultured in
658 Middlebrook 7H9 broth (271310, BD Diagnostics, Franklin Lakes, NJ, USA) with 0.2% (v/v)
659 glycerol, 0.05% (v/v) Tween-80 (P1754, Sigma-Aldrich, St. Louis, MO, USA), and 10% (v/v)
660 OADC (90000-614, VWR, Radnor, PA, USA). Proliferation rates were determined by inoculating
661 5 mL media with 500,000 colony forming units (cfu), then measuring OD600 at 24 and 48 hours
662 post-inoculation. Proliferation rate was calculated as follows: proliferation rate in doublings per
663 day = $\log_2(\text{OD600 at 48 hours} / \text{OD600 at 24 hours})$.

664

665 Fluctuation Assay: Fluctuation assays were performed as described previously¹⁷ with adaptations
666 for *M. abscessus*. *M. abscessus* strains were grown to saturation, then diluted to an OD600 <
667 0.01 and grown overnight until cultures reached an OD600 = 0.6-0.8. These cultures were used
668 for fluctuation assays as well as determination of baseline amikacin sensitivity. Baseline amikacin
669 sensitivity was tested by spreading dilutions of 10 million, 100 million, or 1 billion cfu on agar 7H10
670 (262710, BD Diagnostics) + OADC (90000-614, VWR) + 300 µg/mL amikacin sulfate (A2324,
671 Sigma-Aldrich) plates. For fluctuation assays, 5000 cfu were plated in 1 mL media in each of 21
672 wells of a 24-well culture plate (10861-558, VWR), and initial cultures were plated for cfu on 7H10
673 + OADC to confirm number of cells in initial culture. 24-well plates were incubated shaking at 150
674 rpm for 6 days at 37°C. 18 wells of each strain were plated onto 6-well plates containing 5 mL
675 7H10 agar + OADC + 300 µg/mL amikacin sulfate in each well. Amikacin was used as the
676 antibiotic selection for the assay because all tested strains had low baseline resistance (Extended
677 Data Table 5), which is required for the validity of the fluctuation assay. The remaining 3 wells
678 from each strain were diluted and plated on 7H10 + OADC plates to enumerate average final cfu.
679 Colonies on all agar plates were counted after 5 days incubation at 37°C. Mutation rate per
680 generation per total number of bases that can be mutated to confer amikacin resistance was
681 calculated using the Shinyflan application for the flan R package¹⁷ using the Maximum Likelihood
682 estimation method, assumption of unknown fitness effects of mutations, the Exponential (LD
683 model) distribution of mutant lifetime, and a Winsor parameter of 1024. One strain (BWH-E) was
684 excluded from analysis due to highly variable results across replicates (Figure S6B).

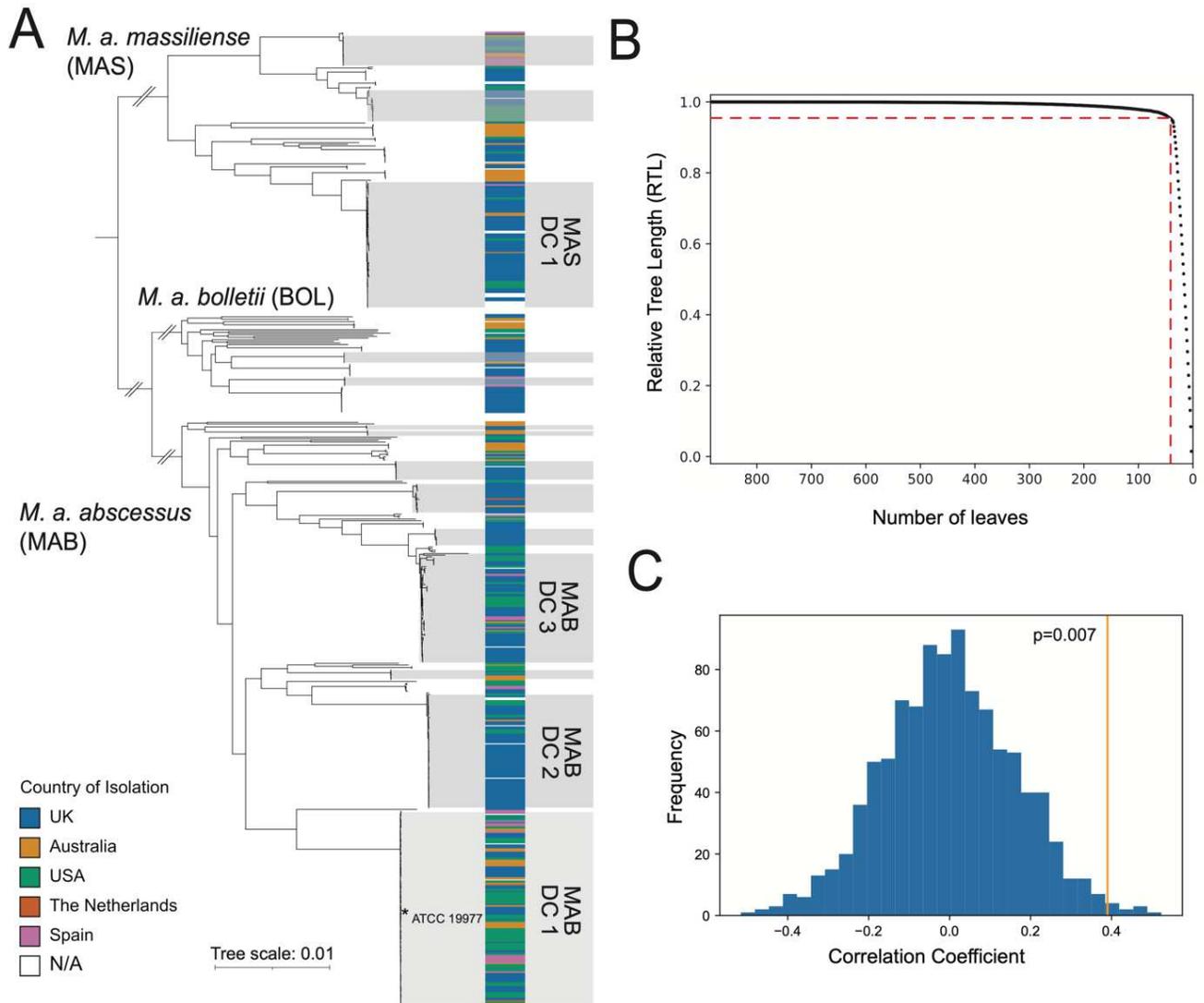
685

686 **Code Availability**

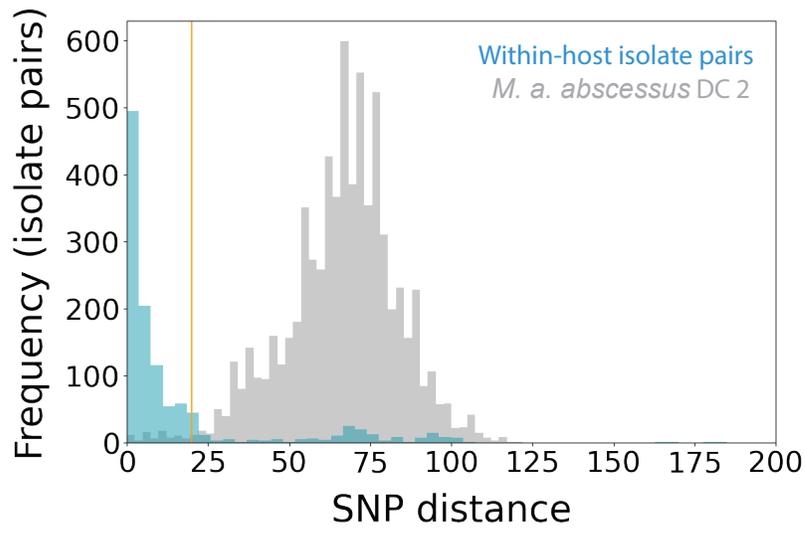
687 All custom scripts are available at:

688 https://github.com/nicolettacommins/mab_mutation_rates_2022.

689

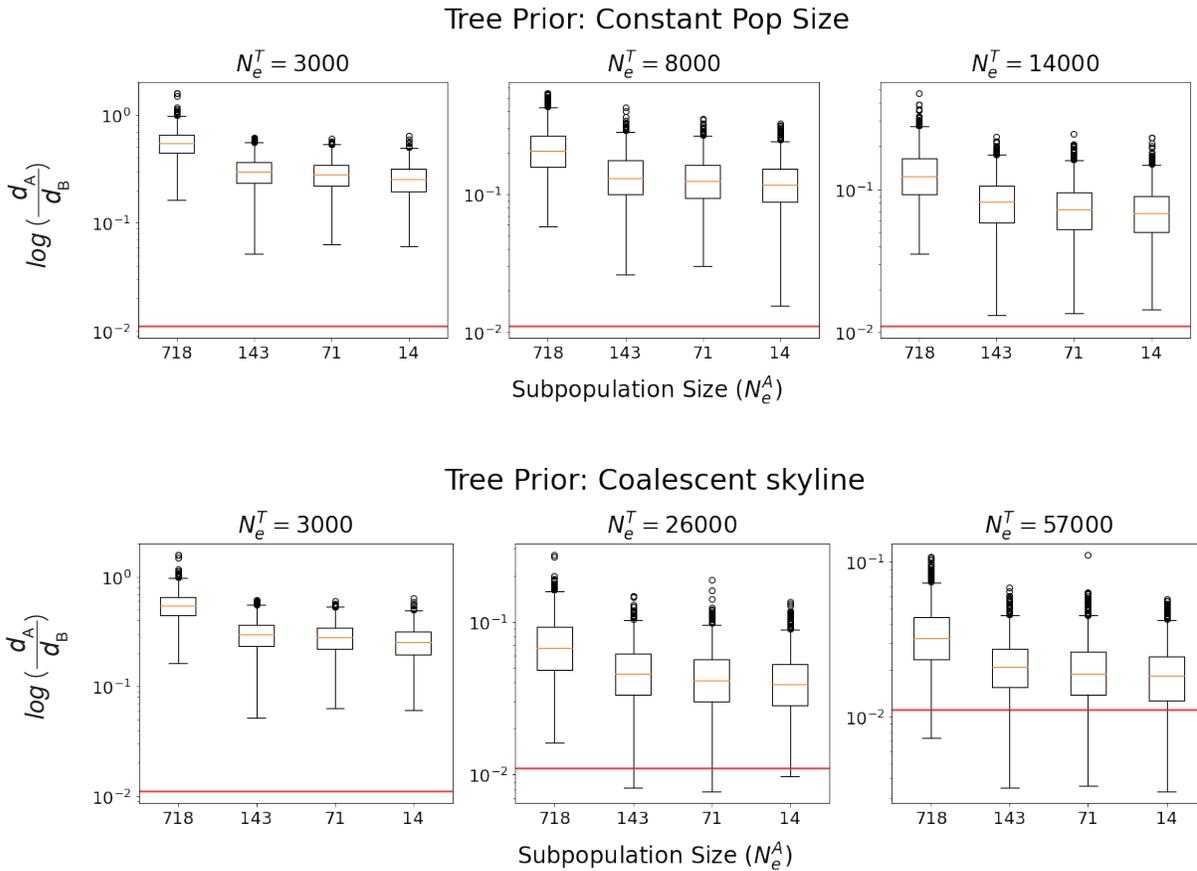


691
 692 **Extended Data Figure 1. *M. a. abscessus* clusters and inference of temporal signal (related**
 693 **to Figure 1).** A) Full species tree showing all clusters identified including DCs and non DC clusters. DCs are labeled. B) The relative tree length (RTL) v. the number of leaves remaining in the tree after each cycle of pruning with Treemmer. RTL is a measure of the diversity in the pruned tree relative to the original tree. The dotted red line indicates where the RTL is 95%, corresponding to 38 isolates. C) Permutation test for significance of temporal signal in the pruned *M. a. abscessus* tree. We randomly shuffled the dates of collection and root-to-tip distances from the pruned subspecies tree 1000 times. For each permutation we calculated Pearson's correlation between the collection dates and root-to-tip distances to estimate an empirical p-value = 0.007.



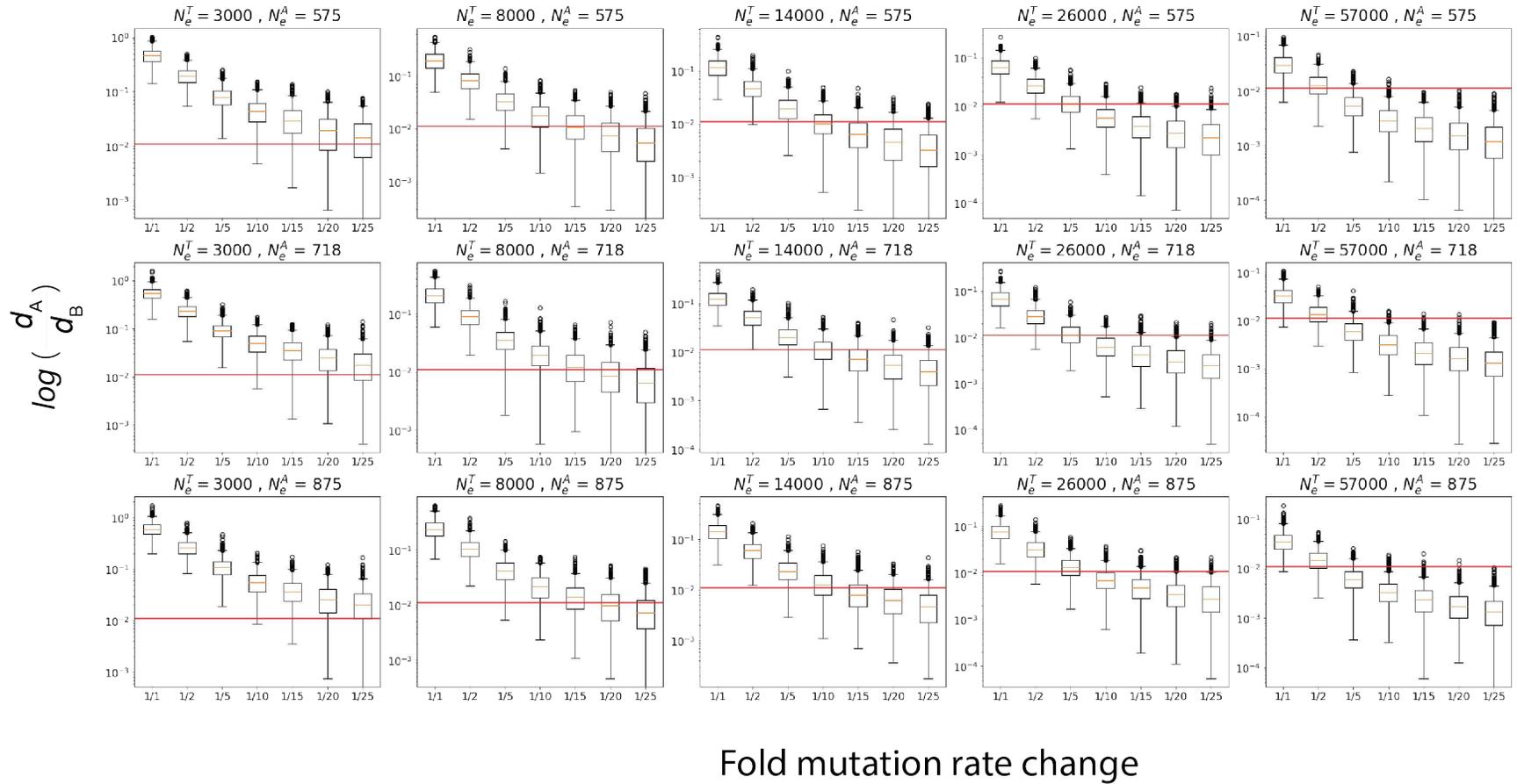
Extended Data Figure 2. SNP distance threshold for within-host isolate pairs (related to Figure 2A). Distribution of SNP distances between all possible within-host isolate pairs in *M. a. abscessus* compared to SNP distances between all possible isolate pairs within *M. a. abscessus* DC 2.

694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719



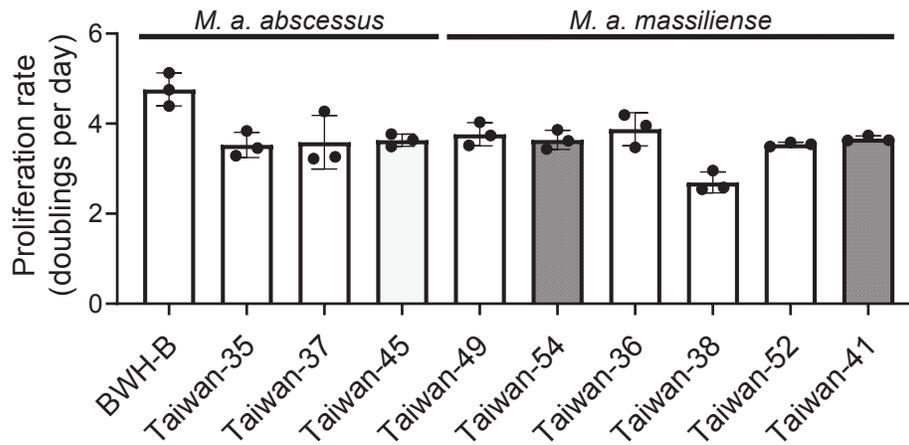
720
 721
 722 **Extended Data Figure 3. Ancestry simulations across a range of possible population sizes**
 723 **(related to Figure 3A).** Boxplots showing the degree of phylogenetic clustering in subpopulation A
 724 relative to subpopulation B (Figure 3A) over a range of effective population sizes of subpopulation A
 725 (N_e^A) and a range of total effective population sizes (N_e^T) spanning the 95% confidence intervals
 726 estimated using BEAST under a constant population size model (top) and a Bayesian coalescent
 727 skyline model (bottom). The top center panel recapitulates figure 3B.

728
 729
 730
 731
 732
 733
 734

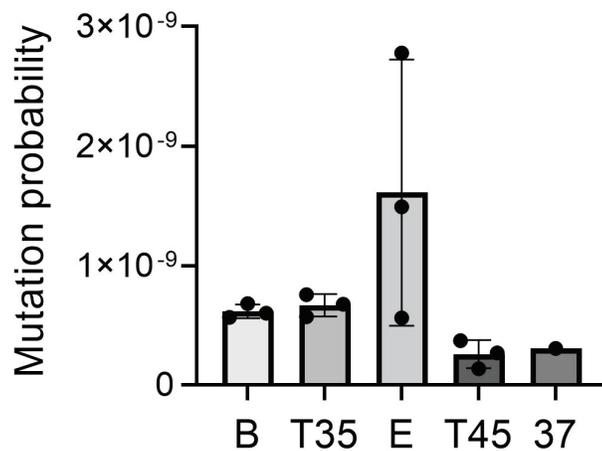


Extended Data Figure 4. Ancestry simulations over a range of population sizes and mutation rate changes (Related to Figure 3B,C). Boxplot showing the degree of phylogenetic clustering in subpopulation A relative to subpopulation B over a range of mutation rate changes. Each set of simulations is repeated over a range of effective population sizes of subpopulation A (N_e^A) and a range of total effective population sizes (N_e^T) spanning the 95% confidence intervals estimated using BEAST assuming a constant population size tree prior. The panel in which $N_e^T=8000$ and $N_e^A=718$ recapitulates Figure 3B.

A



B



736
737
738
739
740
741
742
743
744

Extended Data Figure 5. Growth rates and mutation probability of *M. abscessus* strains. (Related to Figure 4C-E). A) Growth rates of all isolates used in fluctuation assays. B) Probability of acquiring an amikacin resistance mutation per generation in *M. a. abscessus* strains including BWH-E. Data for BWH-B, Taiwan-35, Taiwan-37, and Taiwan-45 are reproduced from Figure 4C. Mean +/- SD for 3 biological replicates is displayed.

745
746
747
748

749 **Extended Data Table 1.** Isolates used in this study

Number of samples	BioProject ID	Citation
854	PRJEB2779	Bryant, J et al. 2016. "Emergence and Spread of a Human-Transmissible Multidrug-Resistant Nontuberculous Mycobacterium." <i>Science</i> 354 (6313): 751–57.
190	PRJNA319839	Hasan, NA et al. 2019. "Population Genomics of Nontuberculous Mycobacteria Recovered from United States Cystic Fibrosis Patients." <i>bioRxiv</i> . https://doi.org/10.1101/663559 .
173	PRJEB7058	Everall, I et al. 2017. "Genomic Epidemiology of a National Outbreak of Post-Surgical Mycobacterium Abscessus Wound Infections in Brazil." <i>Microbial Genomics</i> 3 (5): e000111.
141	PRJEB31559	Doyle, RM et al. 2019. "Cross-Transmission Is Not the Source of New Mycobacterium Abscessus Infections in a Multi-Centre Cohort of Cystic Fibrosis Patients." <i>Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America</i> , June. https://doi.org/10.1093/cid/ciz526 .
69	PRJNA420644	Lipworth, S. <i>et al.</i> Whole-Genome Sequencing for Predicting Clarithromycin Resistance in Mycobacterium abscessus. <i>Antimicrob. Agents Chemother.</i> 63 , (2019)
17	PRJNA439313	Yan, J. <i>et al.</i> Investigating transmission of Mycobacterium abscessus amongst children in an Australian cystic fibrosis centre. <i>J. Cyst. Fibros.</i> 19 , 219–224 (2020)
11	PRJNA297030	Unpublished
1	PRJNA447908	Unpublished
1	PRJNA495001	Chhotaray, C. <i>et al.</i> Comparative Analysis of Whole-Genome and Methylome Profiles of a Smooth and a Rough Mycobacterium abscessus Clinical Strain. <i>G3</i> 10 , 13–22 (2020)
1	PRJEB1520	Unpublished
1	PRJNA347845	Fogelson, S. B. <i>et al.</i> Variation among human, veterinary and environmental Mycobacterium chelonae-abscessus complex isolates observed using core genome phylogenomic analysis, targeted gene comparison, and anti-microbial susceptibility patterns. <i>PLoS One</i> 14 , e0214274 (2019)
1	PRJNA566387	Pearce, C. <i>et al.</i> Inhaled tigecycline is effective against Mycobacterium abscessus in vitro and in vivo. <i>J. Antimicrob. Chemother.</i> 75 , 1889–1894 (2020)
1	PRJNA401495	Unpublished

750

751

752

753 **Extended Data Table 2.** Description of Isolation Sources

Isolation Source	Number of samples (%)
Pulmonary	1181 (80.8%)
Skin and soft tissue swab or biopsy	185 (12.7%)
Lymph node	8 (0.5%)
Environmental	4 (0.2%)
CSF	3 (0.2%)
Blood	2 (0.1%)
Feces	2 (0.1%)
Clinical isolates from unknown source	76 (5.2%)

754

755

756 **Extended Data Table 3.** Results from Bayesian Evaluation of Temporal Signal

Model	Marginal likelihood (SD)
Relaxed clock, constant population size (with dates)	-7144301.64 (21.42)
Relaxed clock, constant population size (without dates)	-7144579.04 (21.01)

757

758 Comparison of the marginal likelihoods with sampling dates and without sampling dates yields
 759 the following Bayes Factor (BF) calculation: $\text{Log}_{10}(\text{BF}) = 277.4$

760

761

762

763

764

765

766

767

768

769

770

771

772

773 **Extended Data Table 4.** BEAST estimated N_e for the 'clustered' subpopulation A under a range
 774 of assumed clock rates
 775

Clock rate (SNPs/site/year)	N_e^A [95% HPD]
1.7×10^{-7} ^a	718.99 [575.06, 875.16]
8.5×10^{-7}	142.79 [116.38, 175.52]
1.7×10^{-6}	71.75 [57.57, 87.07]
8.5×10^{-6}	14.27 [11.47, 17.35]

776 ^a Mutation rate estimated for *M. a. abscessus* cluster 1 by Ruis et al.¹³

777
 778
 779
 780
 781

Extended Data Table 5. Baseline Rates of Resistance to Amikacin in Clinical Isolates

Strain	Fraction of cells resistant to amikacin
Mas-Taiwan-36	n.d.
Mas-Taiwan-38	8.77×10^{-8}
Mas-Taiwan-41	n.d.
Mas-Taiwan-52	1.59×10^{-8}
Mas-Taiwan-49	n.d.
Mas-Taiwan-54	5×10^{-8}
Mab-BWH-B	1×10^{-8}
Mab-BWH-E	1×10^{-8}
Mab-Taiwan-35	1×10^{-8}
Mab-Taiwan-37	n.d.
Mab-Taiwan-45	n.d.

782
 783 n.d. = no resistant colonies were detected

784
 785
 786
 787
 788
 789
 790
 791
 792

793 **References**

- 794 1. Adjemian, J., Olivier, K. N. & Prevots, D. R. Epidemiology of Pulmonary Nontuberculous
795 Mycobacterial Sputum Positivity in Patients with Cystic Fibrosis in the United States, 2010-
796 2014. *Ann. Am. Thorac. Soc.* **15**, 817–826 (2018).
- 797 2. Lee, M.-R. *et al.* Mycobacterium abscessus Complex Infections in Humans. *Emerg. Infect.*
798 *Dis.* **21**, 1638–1646 (2015).
- 799 3. Nessar, R., Cambau, E., Reytrat, J. M., Murray, A. & Gicquel, B. Mycobacterium abscessus:
800 a new antibiotic nightmare. *J. Antimicrob. Chemother.* **67**, 810–818 (2012).
- 801 4. Bryant, J. M. *et al.* Emergence and spread of a human-transmissible multidrug-resistant
802 nontuberculous mycobacterium. *Science* **354**, 751–757 (2016).
- 803 5. Bange, F. C., Brown, B. A., Smaczny, C., Wallace, R. J., Jr & Böttger, E. C. Lack of
804 transmission of mycobacterium abscessus among patients with cystic fibrosis attending a
805 single clinic. *Clin. Infect. Dis.* **32**, 1648–1650 (2001).
- 806 6. Olivier, K. N. *et al.* Nontuberculous mycobacteria. I: multicenter prevalence study in cystic
807 fibrosis. *Am. J. Respir. Crit. Care Med.* **167**, 828–834 (2003).
- 808 7. Doyle, R. M. *et al.* Cross-transmission is not the source of new Mycobacterium abscessus
809 infections in a multi-centre cohort of cystic fibrosis patients. *Clin. Infect. Dis.* (2019)
810 doi:10.1093/cid/ciz526.
- 811 8. Sermet-Gaudelus, I. *et al.* Mycobacterium abscessus and children with cystic fibrosis. *Emerg.*
812 *Infect. Dis.* **9**, 1587–1591 (2003).
- 813 9. Lipworth, S. *et al.* Epidemiology of Mycobacterium abscessus in England: an observational
814 study. *The Lancet Microbe* (2021) doi:10.1016/S2666-5247(21)00128-2.
- 815 10. Harris, K. A. *et al.* Whole-genome sequencing and epidemiological analysis do not provide
816 evidence for cross-transmission of mycobacterium abscessus in a cohort of pediatric cystic
817 fibrosis patients. *Clin. Infect. Dis.* **60**, 1007–1016 (2015).

- 818 11. Tortoli, E. *et al.* Mycobacterium abscessus in patients with cystic fibrosis: low impact of inter-
819 human transmission in Italy. *Eur. Respir. J.* **50**, (2017).
- 820 12. Everall, I. *et al.* Genomic epidemiology of a national outbreak of post-surgical Mycobacterium
821 abscessus wound infections in Brazil. *Microb Genom* **3**, e000111 (2017).
- 822 13. Ruis, C. *et al.* Dissemination of Mycobacterium abscessus via global transmission networks.
823 *Nature Microbiology* 1–10 (2021) doi:10.1038/s41564-021-00963-3.
- 824 14. Bouckaert, R. *et al.* BEAST 2.5: An advanced software platform for Bayesian evolutionary
825 analysis. *PLoS Comput. Biol.* **15**, e1006650 (2019).
- 826 15. Houghton, J. *et al.* Important role for Mycobacterium tuberculosis UvrD1 in pathogenesis and
827 persistence apart from its function in nucleotide excision repair. *J. Bacteriol.* **194**, 2916–2923
828 (2012).
- 829 16. Sinha, K. M., Unciuleac, M.-C., Glickman, M. S. & Shuman, S. AdnAB: a new DSB-resecting
830 motor-nuclease from mycobacteria. *Genes Dev.* **23**, 1423–1437 (2009).
- 831 17. Krašovec, R. *et al.* Measuring Microbial Mutation Rates with the Fluctuation Assay. *J. Vis.*
832 *Exp.* (2019) doi:10.3791/60406.
- 833 18. Ho, S. Y. W., Shapiro, B., Phillips, M. J., Cooper, A. & Drummond, A. J. Evidence for time
834 dependency of molecular rate estimates. *Syst. Biol.* **56**, 515–522 (2007).
- 835 19. Huffnagle, G. B., Dickson, R. P. & Lukacs, N. W. The respiratory tract microbiome and lung
836 inflammation: a two-way street. *Mucosal Immunol.* **10**, 299–306 (2017).
- 837 20. Curti, E., Smerdon, S. J. & Davis, E. O. Characterization of the helicase activity and substrate
838 specificity of Mycobacterium tuberculosis UvrD. *J. Bacteriol.* **189**, 1542–1555 (2007).
- 839 21. Bryant, J. M. *et al.* Whole-genome sequencing to identify transmission of Mycobacterium
840 abscessus between patients with cystic fibrosis: a retrospective cohort study. *Lancet* **381**,
841 1551–1560 (2013).
- 842 22. Hasan, N. A. *et al.* Population Genomics of Nontuberculous Mycobacteria Recovered from
843 United States Cystic Fibrosis Patients. *bioRxiv* 663559 (2019) doi:10.1101/663559.

- 844 23. Shaw, L. P. *et al.* Children with cystic fibrosis are infected with multiple subpopulations of
845 *Mycobacterium abscessus* with different antimicrobial resistance profiles. *Clin. Infect. Dis.*
846 (2019) doi:10.1093/cid/ciz069.
- 847 24. Download : Software : Sequence Read Archive : NCBI/NLM/NIH.
848 <https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software>.
- 849 25. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
850 data. *Bioinformatics* **30**, 2114–2120 (2014).
- 851 26. Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput Sequence
852 Data. <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- 853 27. Prjibelski, A., Antipov, D., Meleshko, D., Lapidus, A. & Korobeynikov, A. Using SPAdes De
854 Novo Assembler. *Curr. Protoc. Bioinformatics* **70**, e102 (2020).
- 855 28. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput
856 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.*
857 **9**, 5114 (2018).
- 858 29. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
859 *arXiv [q-bio.GN]* (2013).
- 860 30. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
861 and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 862 31. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and
863 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–
864 2993 (2011).
- 865 32. Choo, S. W. *et al.* Genomic reconnaissance of clinical isolates of emerging human pathogen
866 *Mycobacterium abscessus* reveals high evolutionary potential. *Sci. Rep.* **4**, 4061 (2014).
- 867 33. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*
868 **31**, 3691–3693 (2015).

- 869 34. Pockrandt, C., Alzamel, M., Iliopoulos, C. S. & Reinert, K. GenMap: ultra-fast computation of
870 genome mappability. *Bioinformatics* **36**, 3687–3692 (2020).
- 871 35. Derrien, T. *et al.* Fast computation and applications of genome mappability. *PLoS One* **7**,
872 e30377 (2012).
- 873 36. Marin, M. *et al.* Benchmarking the empirical accuracy of short-read sequencing across the
874 *M. tuberculosis* genome. *Bioinformatics* (2022) doi:10.1093/bioinformatics/btac023.
- 875 37. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial
876 whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
- 877 38. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference
878 in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 879 39. Harris, S. TreeGubbins. (2016).
- 880 40. Menardo, F. *et al.* Treemmer: a tool to reduce large phylogenetic datasets with minimal loss
881 of diversity. *BMC Bioinformatics* **19**, 164 (2018).
- 882 41. Darriba, D. *et al.* ModelTest-NG: A New and Scalable Tool for the Selection of DNA and
883 Protein Evolutionary Models. *Mol. Biol. Evol.* **37**, 291–294 (2020).
- 884 42. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary
885 Genetics Analysis across Computing Platforms. *Mol. Biol. Evol.* **35**, 1547–1549 (2018).
- 886 43. Kelleher, J. & Lohse, K. Coalescent Simulation with msprime. in *Statistical Population*
887 *Genomics* (ed. Dutheil, J. Y.) 191–230 (Springer US, 2020). doi:10.1007/978-1-0716-0199-
888 0_9.
- 889 44. Edwards, D. J., Duchêne, S., Pope, B. & Holt, K. E. SNPPar: identifying convergent evolution
890 and other homoplasies from microbial whole-genome alignments. *bioRxiv*
891 2020.07.08.194480 (2020) doi:10.1101/2020.07.08.194480.
- 892
- 893
- 894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [mabmutRatenatureMicrosupplement.pdf](#)