

1 Title

2 Lineage specific histories of *Mycobacterium tuberculosis* dispersal in Africa and Eurasia

3

4 Keywords

5 phylogeography, evolution, pathogen, migration, demography

6

7 Abstract

8 *Mycobacterium tuberculosis* (*M.tb*) is a globally distributed, obligate pathogen of humans that  
9 can be divided into seven clearly defined lineages. Identifying how the ancestral clone of *M.tb*  
10 spread and differentiated is important for identifying the ecological drivers of the current  
11 pandemic. We reconstructed *M.tb* migration in Africa and Eurasia, and investigated lineage  
12 specific patterns of spread. Applying evolutionary rates inferred with ancient *M.tb* genome  
13 calibration, we link *M.tb* dispersal to historical phenomena that altered patterns of connectivity  
14 throughout Africa and Eurasia: trans-Indian Ocean trade in spices and other goods, the Silk Road  
15 and its predecessors, the expansion of the Roman Empire and, the European Age of Exploration.  
16 We find that Eastern Africa and Southeast Asia have been critical in the dispersal of *M.tb*. Our  
17 results reveal complex relationships between spatial dispersal and expansion of *M.tb* populations,  
18 and delineate the independent evolutionary trajectories of bacterial sub-populations underlying  
19 the current pandemic.

20

21 Introduction

22 The history of tuberculosis (TB) has been rewritten several times as genetic data accumulate  
23 from its causative agent, *Mycobacterium tuberculosis* (*M.tb*). In the nascent genomic era, these

24 data refuted the long-held hypothesis that human-adapted *M.tb* emerged from an animal adapted  
25 genetic background represented among extant bacteria by *Mycobacterium bovis*, another member  
26 of the *Mycobacterium tuberculosis* complex (MTBC) (Brosch et al. 2002). Genetic data from  
27 bacteria infecting multiple species of hosts revealed that currently known non-primate-adapted  
28 strains form a nested clade within the diversity of extant *M.tb* (Behr et al. 1999; Brosch et al.  
29 2002; Hershberg et al. 2008).

30  
31 *M.tb* can be classified into seven well-differentiated lineages, which differ in their geographic  
32 distribution and association with human sub-populations (Hirsh et al. 2004; Gagneux et al.  
33 2006). This observation led to the hypothesis that *M.tb* diversity has been shaped by human  
34 migrations out of Africa, and that the most recent common ancestor (MRCA) of extant *M.tb*  
35 emerged in Africa approximately 73,000 years ago, coincident with estimated waves of human  
36 migration (Comas et al. 2013). Human out of Africa migrations are a plausible means by which  
37 *M.tb* could have spread globally. However, *M.tb* evolutionary rate estimates based on a variety of  
38 calibration methods are inconsistent with the out of Africa hypothesis (Eldholm et al. 2016;  
39 Brynildsrud et al. 2018).

40  
41 When calibrated with ancient DNA, the estimates of the time to most recent common ancestor  
42 (TMRCA) for the MTBC are <6,000 years before present (Bos et al. 2014; Kay et al. 2015).  
43 This is not necessarily the time period over which TB first emerged, as it is possible –  
44 particularly given the apparent absence of recombination among *M.tb* (Pepperell et al. 2013) –  
45 that the global population has undergone clonal replacement events that displaced ancient  
46 diversity from the species.

47  
48 *M.tb* is an obligate pathogen of humans with a global geographic range. The finding of a recent  
49 origin for the extant *M.tb* population raises the question of how the organism could have spread  
50 within this timeframe to occupy its current distribution. *M.tb* populations in the Americas show  
51 the impacts of European colonial movements as well as recent immigration (Pepperell et al.  
52 2011; Brynildsrud et al. 2018); the role of other historical phenomena in driving TB dispersal is  
53 not well understood. Here we sought to reconstruct the migratory history of *M.tb* populations in  
54 Africa and Eurasia within the framework of a recent origin and evolutionary rates derived from  
55 ancient DNA data (Bos et al. 2014; Kay et al. 2015). We discovered lineage-specific patterns of  
56 migration and a complex relationship between *M.tb* effective population growth and migration.  
57 Our results connect *M.tb* migration to major historical events in human history that altered  
58 patterns of connectivity in Africa and Eurasia. These findings provide context for a recent  
59 evolutionary origin of the MRCA of *M.tb* (Pepperell et al. 2013; Bos et al. 2014; Kay et al.  
60 2015), which represents yet another paradigm shift in our understanding of the history and origin  
61 of this successful pathogen.

62

## 63 Results

### 64 ***Genetic and geographic structures of global *M.tb* populations***

65 In order to establish the contemporary geographic distributions of *M.tb* lineages, we translated  
66 the spoligotypes reported for 42,358 *M.tb* isolates to their corresponding lineage designations  
67 (fig. 1). Geographic patterns in prevalence vary between lineages. Lineage 1 (L1) is prevalent  
68 in regions bordering the Indian Ocean, extending from Eastern Africa to Melanesia. Lineage 2  
69 (L2) is broadly distributed, with a predominance in Eastern Eurasia and South East Asia.

70 Lineage 3 (L3) is similar to L1 in that its distribution rings the Indian Ocean, but it does not  
71 extend into Southeastern Asia, it has a stronger presence in Northern Africa, and a broader  
72 distribution across Southern Asia. Lineage 4 (L4) is strikingly well dispersed, with a  
73 predominance throughout Africa and Europe and the entire region bordering the Mediterranean.  
74 Lineages 5 (L5) and 6 (L6) are found at low frequencies in Western and Northern Africa.  
75 Lineage 7 (L7), as previously described (Blouin et al. 2012; Firdessa et al. 2013; Comas et al.  
76 2015), is limited to Ethiopia.

77  
78 We compiled a diverse collection of *M.tb* genomes for phylogenetic and population genetic  
79 inference of the demographic and migratory history of the extant *M.tb* population (*see Methods*).  
80 Our dataset consists of whole-genome sequences (WGS) from 552 *M.tb* isolates collected from  
81 51 countries (spanning 13 UN geoscheme subregions), which we refer to as the Old World  
82 collection (fig. S1, table S1). We included sites in the alignment where at least half of these  
83 isolates had confident data (60,787 variant sites; 3,838,249 bp) for subsequent analyses, unless  
84 otherwise noted.

85  
86 The inferred maximum likelihood phylogeny reveals the well described *M.tb* lineage structure,  
87 and some associations are evident between lineages and geographic regions (defined here by the  
88 United Nations geoscheme) (fig. S2). The phylogeny has an unbalanced shape, with long  
89 internal branches that define the lineages and feathery tips, suggestive of recent population  
90 expansion.

91

92 Genetic diversity, as measured by the numbers of segregating sites and pairwise differences  
93 (Watterson's  $\theta$  and  $\pi$ ), varied among lineages (table 1). L1 and L4 group together and have the  
94 highest diversity; L2, L3, L5, and L6 have similar levels of diversity and form the middle  
95 grouping; L7 has the lowest diversity. We used an analysis of molecular variance (AMOVA) to  
96 delineate the effects of population sub-division on *M.tb* diversity (table 1). The Old World  
97 collection was highly structured among UN subregions (21% of variation attributable to  
98 between-region comparisons), whereas this structure was less apparent when regions were  
99 defined by the botanical contents outlined by the World geographic scheme for recording plant  
100 distributions (14%). This is consistent with *M.tb*'s niche as an obligate human pathogen, with  
101 bacterial population structure directly shaped by that of its host population (i.e. reflected in UN  
102 subregions) rather than climatic and other environmental features (reflected in botanical  
103 continent definitions). We obtained similar results when the lineages were considered  
104 separately, except for L4, which had little evidence of population structure (4% variation among  
105 UN subregions, 2% among botanical continents).

106

### 107 ***Distinct demographic histories of the *M.tb* lineages***

108 Bayesian inferred trees vary among lineages (fig. 2), likely reflecting their distinct demographic  
109 histories. Branch lengths are relatively even across the phylogenies of L1 and L4, whereas L2  
110 and L3 have a less balanced structure. The long, sparse internal branches and radiating tips of  
111 L2 and L3 phylogenies are consistent with an early history during which the effective population  
112 size remained small (and diversity was lost to drift), followed by more recent population  
113 expansion. L5 has a star-like structure, consistent with rapid population expansion. Jointly  
114 inferred Bayesian skyline plot (BSP) reconstructions of effective population sizes over time

115 suggest that lineages 1-6 have undergone expansion (fig. 3 – top panel, fig. S3). We estimate  
116 that L2 and L3 underwent abrupt expansion at approximately the same time, whereas expansions  
117 of L1 and L4 appeared relatively smooth.

118  
119 We used the methods implemented in  $\partial a \partial i$  to reconstruct the demographic histories of each *M.tb*  
120 population (i.e. lineage) from its synonymous site frequency spectrum (SFS). As demographic  
121 inference with  $\partial a \partial i$  is sensitive to missing data, loci at which any sequence in the individual  
122 lineage alignments had a gap or unknown character were removed for these analyses. Consistent  
123 with the BSP analyses performed in BEAST, instantaneous expansion and exponential growth  
124 models offered an improved fit to the data in comparison with the constant population size model  
125 for each lineage and the entire Old World collection (fig. S4). Parameter estimates varied widely  
126 across runs for the exponential growth model, so we report results only for the instantaneous  
127 expansion model (table 1).

128

### 129 *Major events in M.tb's migratory history*

130 There was evidence of isolation by distance in the global *M.tb* population, as assessed with a  
131 Mantel test of correlations between genetic and geographic distances. We defined geographic  
132 distances using three schemes: great circle distances, great circle distances through waypoints of  
133 human migration as described in (Ramachandran et al. 2005), and distances along historical trade  
134 routes. Waypoints are used to make distance estimates more reflective of presumed human  
135 migration patterns (i.e., when calculating between-continent distances, it is generally thought that  
136 humans did not pass through large bodies of water, and thus a waypoint is used). To allow  
137 comparisons between the schemes, values were centered and standardized (*see Methods*).

138 Values of the Mantel test statistic were similar for great circle distances ( $r = 0.16$ ) and trade  
139 network distances ( $r = 0.16$ ), with distances through waypoints reflective of human migration  
140 patterns having a lower value ( $r = 0.14$ ,  $p = 0.0001$  for all three analyses). In analyses of human  
141 genetic data, adjustment of great circle distances with waypoints results in a higher correlation  
142 between genetic and geographic distances (Ramachandran et al. 2005). Our Mantel test results  
143 therefore do not support a pattern of isolation by distance as expected if out of Africa human  
144 migrations were the primary influence on global diversity of extant *M.tb* (Comas et al. 2013).

145  
146 To further investigate a potential influence of ancient human migration on *M.tb* evolution, we  
147 calculated the correlation between *M.tb* genetic diversity ( $\pi$ ) within subregions and their average  
148 distances from Addis Ababa, a proxy for a possible origin of anatomically modern human  
149 expansion out of Africa. Contrary to what is observed for human population diversity  
150 (Ramachandran et al. 2005), we did not observe a significant decline in *M.tb* diversity as a  
151 function of distance in our Old World collection (adjusted R-squared =  $-0.1$ ,  $p = 0.88$ ), nor when  
152 we included samples from the Americas (adjusted R-squared =  $8.9 \times 10^{-4}$ ,  $p = 0.34$ , note S1, fig.  
153 S5, table S2).

154  
155 We used the methods implemented in BEAST to reconstruct the migratory history of the entire  
156 Old World *M.tb* collection as well as individual lineages within it, modelling geographic origin  
157 of isolates (UN subregion or country) as a discrete trait (fig. 4, figs. S6-S10). Using an  
158 evolutionary rate calibrated with 18<sup>th</sup> century *M.tb* DNA of  $5 \times 10^{-8}$  substitutions/site/year (Kay  
159 et al. 2015), which is similar to the rate inferred with data from 1,000 year old specimens (Bos et  
160 al. 2014), our estimate of the time to most recent common ancestor for extant *M.tb* is between

161 4032 BCE and 2172 BCE (table 1; date ranges are based on the upper and lower limits of the  
162 95% highest posterior density (HPD) for the rate reported in Kay et al. (2015) which is more  
163 conservative than the 95% HPD of our model). We infer an African origin for the MRCA  
164 (Eastern or Western subregion, table 1, fig. 4, fig. S6). Shortly after emergence of the common  
165 ancestor, we infer a migration of the L1-L2-L3-L4-L7 ancestral lineage from Western to Eastern  
166 Africa (we estimate prior to 2683 BCE), with subsequent migrations occurring out of Eastern  
167 Africa.

168  
169 In our phylogeographic reconstruction, emergence of L1 follows migration from Eastern Africa  
170 to Southern Asia at some time between the 3<sup>rd</sup> millennium and 4<sup>th</sup> century BCE (table 1, fig. 4,  
171 fig. S6). L1 has an ‘out of India’ phylogeographic pattern (fig. S7), with diverse Indian lineages  
172 interspersed throughout the phylogeny. This suggests that the current distribution of L1 around  
173 the Indian Ocean (fig. 1) arose from migrations out of India, from a pool of bacterial lineages  
174 that diversified following migration from Eastern Africa.

175  
176 The phylogeographic reconstruction further indicates that following the divergence of L1, *M.tb*  
177 continued to diversify in Eastern Africa, with emergence of L7 there, followed by L4 (table 1,  
178 fig. 4, fig. S6). The contemporary distribution of L4 is extremely broad (fig. 1) and in this  
179 analysis of the Old World collection we infer an East African location for the internal branches  
180 of L4. Notably, in the lineage-specific analyses, we infer a European location for these branches  
181 (fig. S8). The difference is likely due to the fact that inference is informed by deeper as well as  
182 descendant nodes in the Old World collection. Together, these results imply close ties between



183 Europe and Africa during the early history of this lineage that we estimate emerged in the 1<sup>st</sup>  
184 century CE (368 BCE-362 CE, table 1).

185

186 After the emergence of L1 and L7 from Eastern Africa, our analyses suggest that a migration  
187 occurring between 697 BCE and 520 CE established L3 in Southern Asia, with subsequent  
188 dispersal out of Southern Asia into its present distribution, which includes Eastern Africa (i.e., a  
189 back migration of L3 to Africa, fig. 1). We estimate that L2 diversified in South Eastern Asia  
190 following migration from Eastern Africa at some point between 697 BCE and 20 BCE (table 1,  
191 fig. 4, fig. S6). Previously published analyses of L2 phylogeography also inferred a Southeast  
192 Asian origin for the lineage (Luo et al. 2015; Liu et al. 2018).

193

#### 194 ***Lineage and region specific patterns of migration***

195 Our phylogeographic reconstruction indicated that temporal trends in migration varied among  
196 lineages (fig. 3 – bottom panel). We infer that L1 was characterized by high levels of migration  
197 until approximately the 7<sup>th</sup> century CE, when the rate of migration decreased abruptly and  
198 remained stable thereafter. L3, by contrast, exhibited consistently low rates of migration. L2  
199 and L4 had more variable trends in migration, as each underwent punctuated increases in  
200 migration rate. Temporal trends in growth and migration are congruent for L2 and L4, with  
201 increases in migration rate preceding effective population expansions; this is not the case for L1  
202 and L3. Taken together, these results suggest that L1 and L3 populations (as well as L5 and L6,  
203 fig. S3b) grew *in situ*, whereas range expansion may have contributed to the growth of L2 and  
204 L4.

205

206 We employed the Bayesian stochastic search variable selection method (BSSVS) in BEAST  
207 (Lemey et al. 2009) to estimate relative migration rates within the most parsimonious migration  
208 matrix. A map showing inferred patterns of connectivity among UN subregions and relative  
209 rates of *M.tb* migration with strong posterior support is shown in fig. 5. South Eastern Asia was  
210 the most connected region in our analyses, with significant rates of migration connecting it to  
211 eight other regions. Eastern Africa, Eastern Europe, and Southern Asia were also highly  
212 connected, with significant rates with six, six, and five other regions, respectively. Western  
213 Africa, Eastern Asia, and Western Asia were the least connected regions, with just one  
214 significant connection each (to Eastern Africa, South Eastern Asia, and Eastern Europe,  
215 respectively). Our sample from Western Asia is, however, limited (table S1) and migration from  
216 this region may have consequently been underestimated. The highest rates of migration were  
217 seen between Eastern Asia and Southeastern Asia, and between Eastern Africa and Southern  
218 Asia.

219  
220 Lineage specific analyses suggest that migration between Southern Asia, Eastern Africa, and  
221 South Eastern Asia has been important for the dispersal of L1, whereas South Eastern Asia and  
222 Eastern Europe have been important for L2 (fig. S11). L3 is similar to L1 in that there is  
223 evidence of relatively high rates of migration between Southern Asia and Eastern Africa. There  
224 is also evidence of migration within Africa between the eastern and southern subregions. In the  
225 analyses of migration for L4, Eastern Africa appeared highly connected with other regions.

226

227 Discussion

228 Our reconstructions of *M.tb* dispersal throughout the Old World delineate a complex migratory  
229 history that varies substantially between bacterial lineages. Patterns of diversity among extant  
230 *M.tb* suggest that historical pathogen populations were capable of moving fluidly over vast  
231 distances. Using evolutionary rate estimates from ancient DNA calibration, we time the  
232 dispersal of *M.tb* to a historical period of exploration, trade, and increased connectivity among  
233 regions of the Old World.

234  
235 Consistent with prior reports (Comas et al. 2013), we infer an origin of *M.tb* on the African  
236 continent (table 1, fig. 4, fig. S6). There is a modest preference for Western Africa over Eastern  
237 Africa (54% versus 38% inferred probability), likely due to the early branching West African  
238 lineages (i.e. *Mycobacterium africanum*, L5 and L6). Larger samples may allow more precise  
239 localization of the *M.tb* MRCA, and Northern Africa in particular is under-studied.

240  
241 We infer L1 to be the first lineage that emerged out of Africa; L1 is currently concentrated in  
242 regions bordering the Indian Ocean from Eastern Africa to Melanesia (fig. 1). In our  
243 phylogeographic reconstruction, the genesis of this lineage traces to migration from Eastern  
244 Africa to Southern Asia at some point between the 3<sup>rd</sup> millennium and 4<sup>th</sup> century BCE, with  
245 subsequent dispersal occurring out of the Indian subcontinent. Our results suggest that the early  
246 history of L1 was characterized by high levels of migration, particularly between Southern Asia  
247 and Eastern Africa, and between Southern Asia and South Eastern Asia (fig. 3, fig. S11). The  
248 geographic distribution of L1, the timing of its emergence and spread, as well as patterns of  
249 connectivity underlying its dispersal, are all consistent with migration via established trans-  
250 Indian Ocean trade routes linking Eastern Africa to Southern and South Eastern Asia (fig. 6).

251 The interval of our timing estimate for the initial migration overlaps with the so-called Middle  
252 Asian Interaction sphere in The Age of Integration (2600-1900 BCE), which is marked by  
253 increased cultural exchange and trade between civilizations of Egypt, Mesopotamia, the Arabian  
254 peninsula, and the Indus Valley (Vogt 1996; Zarins 1996; Parkin and Barnes 2002; Ray 2003;  
255 Coningham and Young 2015). East-West contact and trade across the Indian Ocean intensified  
256 in the first millennium BCE, when maritime networks expanded to include the eastern  
257 Mediterranean, the Red Sea, and the Black Sea (Dilke 1985; Boussac et al. 1995; Ray et al.  
258 1996; Salles 1996). Historical data from the Roman era indicate that crews on trading ships  
259 crossing the Indian Ocean comprised fluid assemblages of individuals from diverse regions,  
260 brought together under conditions favorable for the transmission of TB (André and Filliozat  
261 1986; Begley and De Puma 1991; Wink 2002; Rauh 2003). These ships would have been an  
262 efficient means of spreading *M.tb* among the distant regions involved in trade.  
263  
264 L2 may similarly have an origin in East-West maritime trade across the Indian Ocean, as we  
265 infer it arose from a migration event from Eastern Africa to South Eastern Asia during the 1<sup>st</sup>  
266 millennium BCE. In this era, increased sophistication in ship technology allowed for longer  
267 voyages (Kent 1979; Blench 1996; Ray et al. 1996; Parkin and Barnes 2002; Wink 2002; Ray  
268 2003). L2 appears to have spread out of Southeast Asia, a highly connected region in our  
269 analyses of *M.tb* migration, and is currently found across Eastern Eurasia and throughout South  
270 Eastern Asia (fig. 1, fig. 4, fig. S6, fig. S11). Interestingly, although L2 is dominant in Eastern  
271 Asia, the region did not appear to have played a prominent role in dispersal of this lineage,  
272 except in its exchanges with South Eastern Asia. A recently published study found that the extant  
273 *M.tb* population in China traces to a limited number of introductions (Liu et al. 2018), which is

274 consistent with our findings of relatively few exchanges of *M. tb* between Eastern Asia and other  
275 regions.

276

277 L3 appears to have had relatively low rates of migration throughout its history (fig. 3). The  
278 contemporary geographic range of L3 is also narrower, extending east from Northern Africa  
279 through Western Asia to the Indian subcontinent (fig. 1). A study of lineage prevalence in  
280 Ethiopia showed that L3 is currently concentrated in the north of the country (Comas et al.  
281 2015), consistent with our observed north to south gradient in its distribution on the African  
282 continent. This is in opposition to L1, which has a southern predominance in Ethiopia and across  
283 Eastern Africa (fig. 1). We estimate L3 emerged in Southern Asia ca. 520 CE (177-739 CE).  
284 Pakistan harbors diverse strains belonging to L3 (fig. S9), and the Southern Asia region was  
285 highly connected with Eastern Africa in our analyses (fig. S11). Trade along the Silk Road  
286 connecting Europe and Asia was very active in the middle of the first millennium, when we  
287 estimate L3 emerged (Hansen 2012; Ball 2016); its distribution suggests it spread primarily  
288 along trading routes connecting Northeast Africa, Western Asia, and South Asia (André and  
289 Filliozat 1986; Sartre 1991; Hansen 2012; Ball 2016) (fig. 6). We speculate that this occurred  
290 *via* overland routes, which may have limited the migration of L3 relative to maritime dispersal of  
291 the other lineages.

292

293 The geographic distribution of L4 is strikingly broad (fig. 1) and it exhibits minimal population  
294 structure (table 1). This suggests L4 dispersed efficiently and continued to mix fluidly among  
295 regions, a pattern we would expect if it was carried by an exceptionally mobile population of  
296 hosts. L4 is currently concentrated in regions bordering the Mediterranean, and elsewhere

297 throughout Africa and Europe (fig. 1). We estimate the MRCA of L4 emerged in the 1<sup>st</sup> century  
298 CE (range 368 BCE-362 CE), during the peak of Roman Imperial power across the entire  
299 Mediterranean world and expansionist Roman policies into Africa, Europe, and Mesopotamia  
300 (Luttwak 1976; Isaac 2004). The empire reached its greatest territorial extent in the early second  
301 century CE, when all of North Africa, from the Atlantic Ocean to the Red Sea, was under a  
302 single power, with trade on land and sea facilitated by networks of stone-paved roads and  
303 protected maritime routes (Luttwak 1976; Millar 1993; Ball 2016). Primary sources from  
304 Roman civilization attest to trade with China, purposeful expeditions for exploration,  
305 cartography, and trade in the Red Sea and Indian Ocean (Pfister and Bellinger 1945; Dilke 1985;  
306 Begley and De Puma 1991; Erdkamp 2002; Butcher 2003).

307

308 We hypothesize that the broad distribution of L4 reflects rapid diffusion from the Mediterranean  
309 region along trade routes extending throughout Africa, the Middle East, and on to India, China,  
310 and South Eastern Asia. High rates of migration appear to have been maintained for this lineage  
311 over much of its evolutionary history (fig. 3); patterns of connectivity implicate Europe and  
312 Africa in its dispersal (fig. S11). The association of L4 with European migrants is well  
313 described, particularly migrants to the Americas (Gagneux et al. 2006; Pepperell et al. 2011;  
314 Brynildsrud et al. 2018). Here we note bacterial population growth preceded geographic range  
315 expansion in L4 ~ca. 15<sup>th</sup> century (fig. 3), which coincides with the onset of the ‘age of  
316 exploration’ (□Ālam and Subrahmanyam 2009) that would have provided numerous  
317 opportunities for spread of this lineage from Europeans to other populations. We also note the  
318 origin and concentration of this lineage on the African continent. Our sample of L4 isolates

319 includes several deeply rooting African isolates, and African isolates are interspersed throughout  
320 the phylogeny (fig. 4, fig. S6, fig. S8).

321

322 The migratory histories of L5, L6, and L7 are less complicated than those of lineages 1-4.

323 Specifically, L5 and L6 are restricted to Western Africa and L7 is found only in Ethiopia (fig. 4,  
324 fig. S6). The reasons for the restricted distributions of these lineages are not immediately

325 obvious: there is evidence in our analyses that other lineages migrated in and out of Western

326 Africa, and Eastern Africa emerged as highly connected and central to the dispersal of *M.tb* (fig.

327 5). A potential explanation is restriction of the pathogen population to human sub-populations

328 with distinct patterns of mobility and connectivity that did not facilitate dispersal. This is likely

329 the case for L7, which was discovered only recently (Blouin et al. 2012), and is currently largely

330 restricted to the highlands of northern Ethiopia (Firdessa et al. 2013; Comas et al. 2015). In the

331 case of L6 (also known as *Mycobacterium africanum*), there is evidence suggesting infection is

332 less likely to progress to active disease than for *M. tuberculosis sensu stricto* (Jong et al. 2008),

333 which could have played a role in limiting its dispersal.

334

335 Our reconstructions of *M.tb*'s migratory history suggest that patterns of migration were highly

336 dynamic: the pathogen appears to have dispersed efficiently, in complex patterns that

337 nonetheless preserved the distinct structure of each lineage. Some findings, notably inference of

338 population expansion, were consistent across lineages. Though growth of the global *M.tb*

339 population has been described previously (Comas et al. 2013; Pepperell et al. 2013), our results

340 here suggest that the pace and magnitude of expansion, and its apparent relationship to trends in

341 migration, varied among lineages (fig. 3, fig. S3, fig. S11).

342  
343 Our analyses suggest that the expansion of L2 was preceded by an impressive increase in its rate  
344 of migration (fig. 3), implying that growth of the pathogen population was facilitated by  
345 expansion into new niches. Our phylogeographic reconstructions implicate Russia, Central Asia,  
346 and Western Asia in the recent migratory history of L2 (fig. S10, fig. S11), which is consistent  
347 with a published phylogeographic analysis of L2 (Luo et al. 2015). The inferred timing of the  
348 growth and increased migration of L2 (~ca. 13<sup>th</sup> century) is close to the well documented  
349 incursion of *Yersinia pestis* from Central Asia into Europe that resulted in explosive plague  
350 epidemics (Benedictow 2004). The experience with plague suggests that patterns of connectivity  
351 among humans and other disease vectors were shifting at this place and time, which would  
352 potentially open new niches for pathogens including *M.tb*.

353  
354 We estimate that L1 underwent expansion ~ca. 17<sup>th</sup> century (fig. 3) but in this case it appears to  
355 have grown *in situ*, e.g. due to changing environmental conditions such as increased crowding,  
356 and/or growth of local human populations. A study of the molecular epidemiology of TB in  
357 Vietnam identified numerous recent migrations of L2 and L4 into the region, versus a stable  
358 presence of L1 (Holt et al. 2018); this is consistent with our finding of higher recent rates of  
359 migration for L2 and L4 versus L1 (fig. 3). A pattern similar to L1 has been identified  
360 previously, in the delay between dispersal of *M.tb* from European migrants to Canadian First  
361 Nations and later epidemics of TB driven by shifting disease ecology (Pepperell et al. 2011).  
362 These results demonstrate the complex relationship between *M.tb* population growth and  
363 migration, and show that under favorable conditions the pathogen can expand into novel niches  
364 or accommodate growth in an existing niche.



365  
366 In a previous study, analyses of synonymous and non-synonymous SFS have been used to  
367 delineate effects of purifying selection, linkage of sites, and population expansion on global  
368 populations of *M.tb* (Pepperell et al. 2013). Simulation studies have shown that purifying  
369 selection can affect demographic inference with BEAST and SFS-based methods (Ewing and  
370 Jensen 2015; Lapierre et al. 2016). Although our analyses here using  $\hat{d}a\hat{d}i$  were restricted to  
371 synonymous SFS, it is likely that inference of population size changes with this method and with  
372 BEAST were affected by purifying selection on this fully linked genome. The magnitude of  
373 inferred expansions may thus reflect both population size changes and background selection, and  
374 should not be interpreted as direct reflections of historical changes in census population size. We  
375 did not detect an effect of purifying selection on inference of migration in our three population  
376 simulation analyses (note S2, fig. S12, fig. S13), but differences in the strength of purifying  
377 selection could contribute to the lineage-specific differences we observed in the size of inferred  
378 population expansions: i.e., genome-wide patterns of purifying selection could differ among  
379 lineages. Previous evidence has suggested that the fitness trade-offs of drug resistance mutations  
380 vary among lineages (Mortimer et al. 2018), making this intriguing possibility potentially  
381 feasible.

382  
383 This study has some important limitations. We did not attempt to estimate the rate or timescale  
384 of *M.tb* evolution, instead relying on published rates that were calibrated with ancient DNA.  
385 This is an active area of research, and newly discovered ancient *M.tb* DNA samples will likely  
386 refine inference of both the timing and locations of historical migration events, though it is  
387 critical to note that recent substitution rate estimates of *M.tb* have converged on rates around

388  $5 \times 10^{-8}$  substitutions per site per year (Eldholm et al. 2016). Even when substitution rate  
389 estimates can be estimated with confidence, the precision with which individual events can be  
390 dated using genetic data should not be over-stated, as evidenced by broad 95% credible intervals  
391 for internal node date estimates (e.g., Eldholm et al. 2016). Our goal here was to reconstruct  
392 historical migration of *M.tb* throughout Eurasia and Africa and place this evolutionary history  
393 within a broad historical context; the historical phenomena that we connect with the spread of  
394 TB involved vast areas and extended over hundreds and in some cases thousands of years. Our  
395 reconstruction of the global dispersal of TB within a temporal framework provided by ancient  
396 *M.tb* DNA analysis links spread of the disease to the first ~1500y of the common era, a period of  
397 remarkable intensification in the connectedness among peoples of Africa, Asia and Europe  
398 (Green 2018).

399

#### 400 Methods

401 **Lineage Frequencies.** The SITVIT WEB database (Demay et al. 2012), which is an open access  
402 *M.tb* molecular markers database, was accessed on September 5, 2016. Spoligotypes were  
403 translated to lineages based on the following study (Shabbeer et al. 2012). The following  
404 conversions were also included: EAI7-BGD2 for L1, CAS for L3, and LAM7-TUR, LAM12-  
405 Madrid1, T5, T3-OSA, and H4 for L4. Isolates containing ambiguous spoligotypes (denoted  
406 with >1 spoligotype) were inspected manually and assigned to appropriate lineages. Relative  
407 lineage frequencies of lineages 1-6 for each country containing data for >10 isolates were  
408 calculated and plotted with the rworldmap package in R (South 2016).

409

#### 410 **Sample Description.**

411 *Old World collection.* We assembled/aligned publicly available whole genome sequences  
412 (WGS) of thousands of *M.tb* isolates from recently published studies and databases for which  
413 country of origin information were known and fell within regions traditionally defined as the Old  
414 World. Isolates were assembled via reference guided assembly (RGA) when FASTQ data were  
415 available and by multiple genome alignment (MGA) when only draft genome assemblies were  
416 accessible (see below). As we were interested in reconstructing historical migrations of the  
417 pathogen, we excluded countries where the majority of contemporary TB cases are identified in  
418 recent immigrants (Government of Canada 2005; White et al. 2017; Australian Government  
419 Department of Health and Ageing; Centers for Disease Control; Institute of Environmental  
420 Science and Research Limited; Public Health England). Due to computational limitations  
421 (BEAST analyses), we necessarily took measures to limit our dataset to <600 isolates. For  
422 countries with large numbers of available genomes, we implemented a sub-sampling strategy  
423 similar one previously described (Thorpe et al. 2017), whereby phylogenetic lineage diversity  
424 was captured thus minimizing the overrepresentation of clonal complexes (e.g., outbreaks):  
425 phylogenetic inference on all isolates available from a country was performed with Fasttree  
426 (Price et al. 2010) and a random isolate was selected from each clade extending from  $n$  branches,  
427 where  $n$  was the desired number of isolates from the country. Numbers of isolates per country  
428 were selected based on the availability of appropriate genome sequence data as well as relative  
429 TB prevalence (fig. S1) (World Health Organization 2017). All isolates belonging to lineages 5-  
430 7 were retained. As a whole, this dataset reflects a ‘mixed’ sampling scheme (Lapierre et al.  
431 2016), where lineages L5-L7 are overrepresented relative to their contemporary frequencies (fig.  
432 1). At the lineage-specific scale, L1-L4 approximate random sampling of available genomes.  
433 Our final Old World collection consisted of the WGS of 552 previously published *M.tb* isolates

434 collected from 51 countries spanning 13 UN geoscheme subregions. Accession numbers and  
435 pertinent information about each sample can be found in table S1.

436

437 We note that our sample necessarily contains a large number of drug-resistant isolates as these  
438 are more commonly sequenced. We also acknowledge that the studies we draw genomes from  
439 may have been subject to other sampling biases for which we are unaware.

440

441 *Northern and Central American collection.* For one analysis, we included an additional 15  
442 isolates from a previous study (Comas et al. 2015) for which the country of origin was within the  
443 Americas. Isolates were assembled via RGA (see below) and their genotypes at the 3,838,249 bp  
444 considered for all analyses of the Old World collection were extracted.

445

446 **Reference Guided Assembly.** Previously published FASTQ data were retrieved from the  
447 National Center for Biotechnology Information (NCBI) sequence read archive (SRA) (Leinonen  
448 et al. 2011). Low-quality bases were trimmed using a threshold quality of 15, and reads resulting  
449 in less than 20bp length were discarded using Trim Galore!

450 ([http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)), which is a wrapper tool

451 around Cutadapt (Martin 2011) and FastQC

452 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Reads were mapped to H37Rv

453 (NC\_000962.3) (Cole et al. 1998) with the MEM algorithm (Li 2013). Duplicates were removed

454 using Picard Tools (<http://picard.sourceforge.net>), and local realignment was performed with

455 GATK (DePristo et al. 2011). To ensure only high quality sequencing data were included,

456 individual sequencing runs for which <80% of the H37Rv genome was covered by at least 20X

457 coverage were discarded, as were runs for which <70% of the reads mapped as determined by  
458 Qualimap (García-Alcalde et al. 2012). Pilon (Walker et al. 2014) was used to call variants with  
459 the following parameters: --variant --mindepth 10 --minmq 40 --minqual 20.

460

461 **Multiple Genome Alignment.** Draft genome assemblies were aligned to H37Rv  
462 (NC\_000962.3) (Cole et al. 1998) with Mugsy v1.2.3 (Angiuoli and Salzberg 2011). Regions  
463 not present in H37Rv were removed and merged with the reference-guided assembly.

464

465 **SNP alignment.** Variant calls (VCFs) were converted to FASTAs with in-house scripts that  
466 treat ambiguous calls and deletions as missing data (available at <https://github.com>).

467 Transposable elements, phage elements, and repetitive families of genes (PE, PPE, and PE-  
468 PGRS gene families) that are poorly resolved with short read sequencing were masked to  
469 missing data. Isolates with >20% missing sites were excluded from the Old World collection  
470 (table S1). Variant positions with respect to H37Rv were extracted with SNP-sites (Page et al.  
471 2016) resulting in 60,818 variant sites. Only sites where at least half of the isolates had  
472 confident data (i.e., non-missing) were included in the phylogeographic models and population  
473 genetic analyses (60,787 variant sites; 3,838,249 bp). 1.7% of variant sites landed in loci  
474 associated with drug resistance (table S3).

475

476 **Geographic Information.** Geographic locations for each of the 552 samples in the Old World  
477 collection were obtained from NCBI and/or the publications in which the isolates were first  
478 described. When precise geographic information was available (e.g., city, province, etc.),  
479 coordinates were obtained from [www.mapcoordinates.net](http://www.mapcoordinates.net). When only country level geographic

480 information was available, the ‘Create Random Point’ tool in ArcGIS 10.3 was used to randomly  
481 place each isolate without specific latitude and longitude inside its respective country;  
482 inhospitable areas (e.g., deserts and high mountains) and unpopulated areas from each country  
483 using 50m data from Natural Earth (<http://www.naturalearthdata.com/downloads>, accessed  
484 February 17, 2016) were excluded as possible coordinates. The ‘precision’ column of table S1  
485 reflects which method was used.

486

487 **Trade Route Information.** Data for all trade routes active throughout Europe, Africa, and Asia  
488 by 1400 CE were compiled from the Old World Trade Routes (OWTRAD) Project  
489 ([www.ciolek.com/owtrad.html](http://www.ciolek.com/owtrad.html), accessed February 17, 2016). For each route, both node  
490 information (trade cities, oases, and caravanserai) and arc information (the routes between nodes)  
491 were imported into ArcGIS (fig. 6). *M.tb* isolate locations were also imported as points and the  
492 ‘Generate Near Table’ tool was used to assign each isolate to its nearest node in the trade  
493 network and is listed in the ‘NearPost’ column of table S1.

494

495 **Maximum Likelihood Inference.** We used RAxML v8.2.3 (Stamatakis 2014) for maximum  
496 likelihood phylogenetic analysis of the Old World collection (all sites where at least half of  
497 isolates had non-missing data) under the general time reversible model of nucleotide substitution  
498 with a gamma distribution to account for site-specific rate heterogeneity. Rapid bootstrapping of  
499 the corresponding SNP alignment was performed with the -autoMR flag, converging after 50  
500 replicates. Tree visualization was created with the ggtree package in R (Yu et al. 2017).

501

502 **Phylogeographic & Demographic Inference with BEAST.** The Old World collection SNP  
503 alignment and individual lineage SNP alignments were analyzed using the Bayesian Markov  
504 Chain Monte Carlo coalescent method implemented in BEAST v1.8 (Drummond and Rambaut  
505 2007) with the BEAGLE library (Ayres et al. 2012) to facilitate rapid likelihood calculations.  
506 Analyses were performed using the general time reversible model of nucleotide substitution with  
507 a gamma distribution to account for rate heterogeneity between sites, a strict molecular clock,  
508 and both constant and Bayesian skyline plot (BSP) demographic models. Country of origin or  
509 the UN subregion for each isolate was modeled as a discrete phylogenetic trait (Lemey et al.  
510 2009). All Markov chains were run for at least 100 million generations, sampled every 10,000  
511 generations, and with the first 10,000,000 generations discarded as burn-in; replicate runs were  
512 performed for analyses and combined to assess convergence. Estimated sample size (ESS)  
513 values of non-nuisance parameters were >200 for all analyses. Site and substitution model  
514 choice were based on previous analyses of *M.tb* global alignments as opposed to an exhaustive  
515 comparison of models which would require unreasonable computational resources. Strict vs  
516 relaxed molecular clocks did not result in altered trends of migration at the lineage level, and  
517 comparisons between analyses using strict and relaxed clocks show strong correlation between  
518 the estimated height of nodes (e.g.,  $R^2 > 0.97$ ; fig. S14). Table S4 provides a summary of  
519 BEAST analyses presented and the results derived from them. Tree visualizations were created  
520 with FigTree (<http://tree.bio.edu.ac.uk/software/figtree/>) and the ggtree package in R (Yu et al.  
521 2017).

522

523 *Phylogeographic reconstruction: limitations and alternatives*

524 These phylogeographic reconstructions are clearly sensitive to sampling, since we cannot  
525 identify the roles of unsampled regions in *M.tb*'s migratory history. We maximized geographic  
526 diversity in our sample, but were limited by available data and some regions – notably Middle  
527 Africa, Northern Africa, and Western Asia – are absent or underrepresented in our sample (fig.  
528 S1). Defining the contributions of these undersampled regions to *M.tb*'s migratory history  
529 awaits more samples and/or further method development.

530

531 De Maio *et al.* (2015) note the sensitivity of discrete trait phylogeographic inference in BEAST  
532 to sample selection, as well as overconfidence in the precision of geographic inference, and  
533 propose BASTA as an alternative (De Maio *et al.* 2015). BASTA is sensitive to the choice of  
534 prior and we did not have ancillary data to guide the selection of a prior for the Old World  
535 migratory history of *M.tb*, precluding its use here. We investigated  $\hat{\alpha}\hat{\omega}$  as an alternative tool for  
536 phylogeographic inference but it did not perform well for this application under conditions of  
537 complete linkage of sites (note S3, fig. S15, fig. S16, table S5, table S6). The phylogeographic  
538 inference method implemented here relies on the assumption that sample size reflects deme size  
539 (Lemey *et al.* 2009; De Maio *et al.* 2015), and within the constraints of available data, we  
540 attempted to adjust our sample sizes according the regional prevalence of TB (see *Methods* and  
541 fig. S1). We also interrogated the relationship between regional sample size and inferred  
542 migration rate and did not observe a strong correlation (fig. S17). According to the  
543 classifications proposed by Lapierre *et al.* (2015), our Old World collection represents a 'mixed'  
544 sampling scheme (see *Methods*).

545



546 **Demographic inference from the observed site frequency spectrum (SFS).** SNP-sites (Page  
547 et al. 2016) was used to convert the Old World collection alignment to a multi-sample VCF and  
548 SnpEff (Cingolani et al. 2012) was used to annotate variants with respect to H37Rv  
549 (NC\_000962.3) (Cole et al. 1998) as synonymous, non-synonymous, or intergenic. Loci at  
550 which any sequence in the population had a gap or unknown character were removed from the  
551 data set. Demographic inference with the synonymous SFS for each of the seven lineages and  
552 the entire collection was performed using  $\hat{\delta}a\hat{\delta}i$  (Gutenkunst et al. 2009). We modeled constant  
553 population size (standard neutral model), an instantaneous expansion model, and an exponential  
554 growth model, and identified the best-fit model and maximal likelihood parameters of the  
555 demographic model given our observed data. Our parameter estimates,  $\nu$  and  $\tau$ , were optimized  
556 for the instantaneous expansion and exponential growth models. Uncertainty analysis of these  
557 parameters were analyzed using the Godambe Information Matrix (Coffman et al. 2016) on 100  
558 samplings of the observed synonymous SFS with replacement and subsequent model inference.  
559

560 **Population genetic statistics.** Nucleotide diversity ( $\pi$ ) and Watterson's theta ( $\theta$ ) for various  
561 population assignments (e.g., lineage, UN subregion) were calculated with EggLib v2.1.10 (De  
562 Mita and Siol 2012).

563  
564 **Analysis of Molecular Variance (AMOVA).** AMOVAs were performed using the  
565 'poppr.amova' function (a wrapper for the ade4 package (Dray et al. 2007) implementation) in  
566 the poppr package in R (Kamvar et al. 2014). Bins were assigned via the following classification  
567 systems: UN geoscheme subregions and Level 1 ('botanical continents') of the World  
568 geographical scheme for recording plant distributions. Isolate assignment can be found in table

569 S1. Genetic distances between isolates were calculated with the ‘dist.dna’ function of the ape  
570 v4.0 package in R (Paradis et al. 2004) from the SNP alignment of the Old World collection.  
571

572 **Mantel tests.** Great circle distances between *M.tb* isolate locations were calculated with the  
573 ‘distVincentyEllipsoid’ function in the geosphere R package (Hijmans et al. 2016). Geographic  
574 distances between isolate locations along the trade network were calculated by adding the great  
575 circle distances from the isolates to the nearest trade hubs and the shortest distance between trade  
576 hubs along the trade network; the latter was determined using an Origin-Destination Cost Matrix  
577 and the ‘Solve’ tool in the Network Analyst Toolbox of ArcGIS which calculates the shortest  
578 distance from each origin to every destination along the arcs in the trade network. In the event  
579 that two isolates were assigned to the same trade post, the great circle distance between the  
580 isolates was used. To calculate the geographic distance between isolates in a manner that reflects  
581 human migrations, the great circle distance between isolates and waypoints were summed.  
582 These were calculated with a custom R function (available at <https://github.com>) using a series  
583 of rules to define whether or not the path between isolates would have gone through a waypoint.  
584 For all three distance metrics, values were log transformed and standardized. Genetic distances  
585 between isolates were calculated with the ‘dist.dna’ function in the ape v4.0 package in R  
586 (Paradis et al. 2004) from the SNP alignment. The ‘mantel’ function of the vegan package in R  
587 (Oksanen et al. 2017) was used to perform a Mantel test between the genetic distance matrix and  
588 each of the three geographic matrices for both the Old World collection and each individual  
589 lineage. Four of the 552 isolates were excluded from these analyses as they were from Kiribati  
590 and trade networks spanning this region were not compiled.  
591

592 **Relationship between genetic diversity and geographic distance from Addis Ababa.** For this  
593 analysis, we added Northern and Central American datasets, assembled in an identical manner to  
594 those of the Old World collection and masked at sites where less than half of the Old World  
595 collection had confident data (3,838,249 bp). For each UN subregion, the mean latitude and  
596 longitude coordinates for all *M.tb* isolates within the region were calculated. The great circle  
597 distances from these average estimates for regions to Addis Ababa were then calculated, using  
598 waypoints for between-continent distance estimates to make them more reflective of presumed  
599 human migration patterns (Ramachandran et al. 2005). Cairo was used as a waypoint for Eastern  
600 Europe, Central Asia, Western Asia, Southern Asia, Eastern Asia, and South Eastern Asia; Cairo  
601 and Istanbul were used as waypoints for Western Europe and Southern Europe; Cairo, Anadyr,  
602 and Prince Rupert were used as waypoints for Northern and Central America. The distance  
603 between each region and Addis Ababa were the sum of the great circle distances between the two  
604 points (the average coordinates for the UN subregion and Addis Ababa) and the waypoint(s) in  
605 the path connecting them, plus the great circle distance(s) between waypoints if two were used.  
606 Treating each UN subregion as a population, the relationship between genetic diversity (assessed  
607 with  $\pi$ ) and geographic distance from Addis Ababa were explored with linear regression for both  
608 the entire Old World collection and individual lineages in R (R Development Core Team). Code  
609 is available at <https://github.com>.

610

611 **Migration Rate Inference.** Migration rates through time were inferred from the Bayesian  
612 maximum clade credibility trees for the entire Old World collection of *M.tb* isolates ( $n = 552$ ).  
613 Individual lineages that contain isolates from multiple UN subregions (i.e., L1:  $n = 89$ , L2:  $n =$   
614 181, L3:  $n = 65$ , and L4:  $n = 143$ ) were extracted and plotted separately. Only nodes with

615 posterior probabilities greater than or equal to 80% were considered. A migration event was  
616 classified as a change in the most probable reconstructed ancestral geographic region from a  
617 parent to child node. Median heights of the parent and child nodes were treated as a range of  
618 time that the migration event could have occurred. The rate of migration through time for each  
619 lineage or the Old World collection was inferred by summing the number of migration events  
620 occurring across every year of the time-scaled phylogeny, divided by the total number of  
621 branches in existence during each year of the time-scaled phylogeny (both those displaying a  
622 migration event and those that do not). Code for these analyses is available at <https://github.com>.

623  
624 Additionally, relative migration rates between UN subregions were derived from the BEAST  
625 analyses of phylogeography. The Bayesian stochastic search variable selection method (BSSVS)  
626 for identifying the most parsimonious migration matrix implemented in BEAST as part of the  
627 discrete phylogeographic migration model (Lemey et al. 2009) allowed us to use Bayes factors  
628 (BF) to identify the migration rates with the greatest posterior support and provide posterior  
629 estimates for their relative rates. Strongly supported relative rates ( $BF > 5$ ) and connectivity  
630 among subregions were visualized with Cytoscape v3.2.0 (Shannon et al. 2003) and  
631 superimposed onto a map generated with the ‘rworldmap’ package in R (South 2016).

632  
633 **Effect of selection on estimates of migration.** We performed demographic forward-in-time  
634 simulations using the SFS\_CODE package (Hernandez 2008), which allows for demographic  
635 models with arbitrarily complex migration and selection regimes. Our simulations were  
636 performed under a simple two population model or with a more complex three population model.  
637 In all simulations,  $N_e$  for each population was 1000,  $\mu$  was 0.001 (O’Neill et al. 2015), and

638 migration between each pair of populations was symmetrical. As there is substantial evidence  
639 for little to no recombination in the *M.tb* genome, our simulations were performed without  
640 recombination.

641

642 The two population simulations were performed under three scenarios: 1) no migration between  
643 populations after initial divergence; 2) constant migration after divergence (per generation  $M =$   
644  $0.5$ ) without selection; and 3) constant migration ( $M = 0.5$ ) with purifying selection (25% of  
645 alleles of each population have a population selection coefficient of  $-1.0$ , and the rest are neutral)  
646 after divergence.

647

648 The three population simulations were performed under five scenarios: 1) no migration between  
649 populations after simultaneous divergence of the three populations; 2) constant, symmetrical  
650 migration after divergence (per generation  $M = 0.5$  for all population pairs) without selection; 3)  
651 constant, symmetrical migration ( $M = 0.5$ ) with purifying selection (25% of alleles in all  
652 populations have a population selection coefficient of  $-1.0$ , and the rest are neutral); 4) constant,  
653 asymmetrical migration after divergence ( $M = 0.5$  for migration between pop0 and pop1,  $M = 5.0$   
654 for migration between pop1 and pop2, and  $M = 0$  for migration between pop0 and pop2) without  
655 selection; and 5) constant, asymmetrical migration after divergence ( $M = 0.5$  between pop0 and  
656 pop1,  $M = 5.0$  between pop1 and pop2, and  $M = 0$  between pop0 and pop2) with purifying  
657 selection (25% of alleles in all populations have a population selection coefficient of  $-1.0$ , and  
658 the rest are neutral).

659

660 For all simulations, 25 samples were taken from each population, and sequences of 100000 bases  
661 were generated. Twenty simulations were performed under each scenario for both the 2  
662 population (60 simulations) and 3 population (100 simulations) models. Each sequence  
663 alignment was subsequently subjected to migration analysis in  $\hat{\delta}a\hat{\delta}i$  (Gutenkunst et al. 2009, see  
664 note S2) and BEAST v1.8.4 (Drummond and Rambaut 2007). For each Bayesian coalescent  
665 analysis, the HKY+G substitution model, a constant population model, and a strict molecular  
666 clock model were used. A discrete symmetrical migration model (Lemey et al. 2009) was used  
667 to determine migration rates, and BSSVS (Lemey et al. 2009) was used to estimate BF support  
668 for migration rates in the 3 population simulations. All Markov chains were run for 10 million  
669 generations or until convergence, with samples taken every 10,000 steps, and 10% discarded as  
670 burn-in. The package Spread3 v0.96 (Bielejec et al. 2016) was used to calculate BF support for  
671 migration rates.

## 672 **References**

- 673 Australian Government Department of Health and Ageing. Tuberculosis notifications in  
674 Australia, 2008 and 2009. Available from:  
675 <http://www.health.gov.au/internet/main/publishing.nsf/Content/cda-cdi3601c.htm#refs>
- 676 André J, Filliozat J. 1986. *L'Inde vue de Rome: textes latins de l'antiquité, relatifs à l'Inde*.  
677 Paris: Belles Lettres
- 678 Angiuoli SV, Salzberg SL. 2011. Mugsy: fast multiple alignment of closely related whole  
679 genomes. *Bioinformatics* 27:334–342.
- 680 Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO, Huelsenbeck JP, Ronquist F,  
681 Swofford DL, Cummings MP, et al. 2012. BEAGLE: An Application Programming  
682 Interface and High-Performance Computing Library for Statistical Phylogenetics. *Syst*.  
683 *Biol.* 61:170–173.
- 684 Ball W. 2016. *Rome in the East: The Transformation of an Empire*. 2nd ed. London & New  
685 York: Routledge
- 686 Begley V, De Puma RD. 1991. *Rome and India: The Ancient Sea Trade*. Madison: University of  
687 Wisconsin Press
- 688 Behr MA, Wilson MA, Gill WP, Salamon H, Schoolnik GK, Rane S, Small PM. 1999.  
689 Comparative genomics of BCG vaccines by whole-genome DNA microarray. *Science*  
690 284:1520–1523.
- 691 Benedictow OJ. 2004. *The Black Death, 1346-1353: the complete history*. Woodbridge, Suffolk,  
692 UK; Rochester, N.Y., USA: Boydell Press
- 693 Bielejec F, Baele G, Rodrigo AG, Suchard MA, Lemey P. 2016. Identifying predictors of time-  
694 inhomogeneous viral evolutionary processes. *Virus Evol.* [Internet] 2. Available from:  
695 [https://academic.oup.com/ve/article/2/2/vew023/2797617/Identifying-predictors-of-time-](https://academic.oup.com/ve/article/2/2/vew023/2797617/Identifying-predictors-of-time-inhomogeneous-viral)  
696 [inhomogeneous-viral](https://academic.oup.com/ve/article/2/2/vew023/2797617/Identifying-predictors-of-time-inhomogeneous-viral)
- 697 Blench R. 1996. The Ethnographic Evidence for Long-distance Contacts between Oceania and  
698 East Africa. In: Reade J, editor. *The Indian Ocean in antiquity*. Available from:  
699 <http://public.eblib.com/choice/publicfullrecord.aspx?p=1517609>
- 700 Blouin Y, Hauck Y, Soler C, Fabre M, Vong R, Dehan C, Cazajous G, Massoure P-L, Kraemer  
701 P, Jenkins A, et al. 2012. Significance of the Identification in the Horn of Africa of an  
702 Exceptionally Deep Branching Mycobacterium tuberculosis Clade. *PLOS ONE*  
703 7:e52841.
- 704 Bos KI, Harkins KM, Herbig A, Coscolla M, Weber N, Comas I, Forrest SA, Bryant JM, Harris  
705 SR, Schuenemann VJ, et al. 2014. Pre-Columbian mycobacterial genomes reveal seals as  
706 a source of New World human tuberculosis. *Nature* 514:494–497.

- 707 Boussac M-F, Salles J-F, France eds. 1995. Athens, Aden, Arikamedu: essays on the  
708 interrelations between India, Arabia, and the eastern Mediterranean. New Delhi:  
709 Manohar: Distributed in South Asia by Foundation Books
- 710 Brosch R, Gordon SV, Marmiesse M, Brodin P, Buchrieser C, Eiglmeier K, Garnier T, Gutierrez  
711 C, Hewinson G, Kremer K, et al. 2002. A new evolutionary scenario for the  
712 Mycobacterium tuberculosis complex. Proc. Natl. Acad. Sci. 99:3684–3689.
- 713 Brynildsrud OB, Pepperell CS, Suffys P, Grandjean L, Monteserin J, Debech N, Bohlin J,  
714 Alfsnes K, Pettersson JO-H, Kirkeleite I, et al. 2018. Global expansion of  
715 Mycobacterium tuberculosis lineage 4 shaped by colonial migration and local adaptation.  
716 Sci. Adv. 4:eaat5869.
- 717 Butcher K. 2003. Roman Syria and the Near East. Los Angeles: J. Paul Getty Museum: Getty  
718 Publications
- 719 Centers for Disease Control. CDC - Reported Tuberculosis in the United States, 2015 - TB.  
720 Available from: <https://www.cdc.gov/tb/statistics/reports/2015/default.htm>
- 721 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012.  
722 A program for annotating and predicting the effects of single nucleotide polymorphisms,  
723 SnpEff. Fly (Austin) 6:80–92.
- 724 Coffman AJ, Hsieh PH, Gravel S, Gutenkunst RN. 2016. Computationally Efficient Composite  
725 Likelihood Statistics for Demographic Inference. Mol. Biol. Evol. 33:591–593.
- 726 Cole ST, Brosch R, Parkhill J, Garnier T, Churcher C, Harris D, Gordon SV, Eiglmeier K, Gas S,  
727 III CEBI, et al. 1998. Erratum: Deciphering the biology of Mycobacterium tuberculosis  
728 from the complete genome sequence. Nat. Lond. 396:190.
- 729 Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, Parkhill J, Malla B, Berg S,  
730 Thwaites G, et al. 2013. Out-of-Africa migration and Neolithic coexpansion of  
731 Mycobacterium tuberculosis with modern humans. Nat. Genet. 45:1176–1182.
- 732 Comas I, Hailu E, Kiros T, Bekele S, Mekonnen W, Gumi B, Tschopp R, Ameni G, Hewinson  
733 RG, Robertson BD, et al. 2015. Population Genomics of Mycobacterium tuberculosis in  
734 Ethiopia Contradicts the Virgin Soil Hypothesis for Human Tuberculosis in Sub-Saharan  
735 Africa. Curr. Biol. 25:3260–3266.
- 736 Coningham R, Young R. 2015. The Archaeology of South Asia: From the Indus to Asoka,  
737 c.6500 BCE–200 CE. Cambridge University Press Available from:  
738 <https://books.google.com/books?id=yaJrCgAAQBAJ>
- 739 De Maio N, Wu C-H, O'Reilly KM, Wilson D. 2015. New Routes to Phylogeography: A  
740 Bayesian Structured Coalescent Approximation. PLoS Genet 11:e1005421.
- 741 De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population  
742 genetics and genomics. BMC Genet. 13:27.



- 743 Demay C, Liens B, Burguière T, Hill V, Couvin D, Millet J, Mokrousov I, Sola C, Zozio T,  
744 Rastogi N. 2012. SITVITWEB – A publicly available international multimarker database  
745 for studying Mycobacterium tuberculosis genetic diversity and molecular epidemiology.  
746 *Infect. Genet. Evol.* 12:755–766.
- 747 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del  
748 Angel G, Rivas MA, Hanna M, et al. 2011. A framework for variation discovery and  
749 genotyping using next-generation DNA sequencing data. *Nat. Genet.* 43:491–498.
- 750 Dilke O. 1985. *Greek and Roman Maps*. Baltimore: Johns Hopkins University Press
- 751 Dray S, Dufour A-B, others. 2007. The ade4 package: implementing the duality diagram for  
752 ecologists. *J. Stat. Softw.* 22:1–20.
- 753 Drummond AJ, Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees.  
754 *BMC Evol. Biol.* 7:214.
- 755 Earl DA, vonHoldt BM. 2012. STRUCTURE HARVESTER: a website and program for  
756 visualizing STRUCTURE output and implementing the Evanno method. *Conserv. Genet.*  
757 *Resour.* 4:359–361.
- 758 Eldholm V, Pettersson JH-O, Brynildsrud OB, Kitchen A, Rasmussen EM, Lillebaek T, Rønning  
759 JO, Crudu V, Mengshoel AT, Debech N, et al. 2016. Armed conflict and population  
760 displacement as drivers of the evolution and dispersal of Mycobacterium tuberculosis.  
761 *Proc. Natl. Acad. Sci.* 113:13881–13886.
- 762 Erdkamp P. 2002. *The Roman Army and the Economy*. Amsterdam: Gieben
- 763 Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of individuals using the  
764 software structure: a simulation study. *Mol. Ecol.* 14:2611–2620.
- 765 Ewing GB, Jensen JD. 2015. The consequences of not accounting for background selection in  
766 demographic inference. *Mol. Ecol.* 25:135–141.
- 767 Firdessa R, Berg S, Hailu E, Schelling E, Gumi B, Erenso G, Gadisa E, Kiros T, Habtamu M,  
768 Hussein J, et al. 2013. Mycobacterial lineages causing pulmonary and extrapulmonary  
769 tuberculosis, Ethiopia. *Emerg. Infect. Dis.* 19:460–463.
- 770 Gagneux S, DeRiemer K, Van T, Kato-Maeda M, de Jong BC, Narayanan S, Nicol M, Niemann  
771 S, Kremer K, Gutierrez MC, et al. 2006. Variable host–pathogen compatibility in  
772 Mycobacterium tuberculosis. *Proc. Natl. Acad. Sci. U. S. A.* 103:2869–2873.
- 773 García-Alcalde F, Okonechnikov K, Carbonell J, Cruz LM, Götz S, Tarazona S, Dopazo J,  
774 Meyer TF, Conesa A. 2012. Qualimap: evaluating next-generation sequencing alignment  
775 data. *Bioinformatics* 28:2678–2679.

- 776 Government of Canada PHA of C. 2005. TUBERCULOSIS PREVENTION AND CONTROL  
777 IN CANADA A FEDERAL FRAMEWORK FOR ACTION. Available from:  
778 <http://www.phac-aspc.gc.ca/index-eng.php>
- 779 Green MH. 2018. Climate and Disease in Medieval Eurasia. *Oxf. Res. Encycl. Asian Hist.*  
780 [Internet]. Available from:  
781 <http://asianhistory.oxfordre.com/view/10.1093/acrefore/9780190277727.001.0001/acrefore-9780190277727-e-6>  
782
- 783 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the Joint  
784 Demographic History of Multiple Populations from Multidimensional SNP Frequency  
785 Data. *PLoS Genet* 5:e1000695.
- 786 Hansen V. 2012. *The Silk Road: A New History*. Oxford: Oxford University Press
- 787 Hernandez RD. 2008. A flexible forward simulator for populations subject to selection and  
788 demography. *Bioinformatics* 24:2786–2787.
- 789 Hershberg R, Lipatov M, Small PM, Sheffer H, Niemann S, Homolka S, Roach JC, Kremer K,  
790 Petrov DA, Feldman MW, et al. 2008. High functional diversity in *Mycobacterium*  
791 tuberculosis driven by genetic drift and human demography. *PLoS Biol.* 6:e311.
- 792 Hijmans RJ, Williams E, Vennes C. 2016. geosphere: Spherical Trigonometry. Available from:  
793 <https://cran.r-project.org/web/packages/geosphere/index.html>
- 794 Hirsh AE, Tsolaki AG, DeRiemer K, Feldman MW, Small PM. 2004. Stable association between  
795 strains of *Mycobacterium tuberculosis* and their human host populations. *Proc Natl Acad*  
796 *Sci U A* 101:4871–4876.
- 797 Holt KE, McAdam P, Thai PVK, Thuong NTT, Ha DTM, Lan NN, Lan NH, Nhu NTQ, Hai HT,  
798 Ha VTN, et al. 2018. Frequent transmission of the *Mycobacterium tuberculosis* Beijing  
799 lineage and positive selection for the EsxW Beijing variant in Vietnam. *Nat. Genet.*  
800 50:849–856.
- 801 Institute of Environmental Science and Research Limited. Tuberculosis in New Zealand: Annual  
802 Report 2014. Available from:  
803 [https://surv.esr.cri.nz/surveillance/AnnualTBReports.php?we\\_objectID=4251](https://surv.esr.cri.nz/surveillance/AnnualTBReports.php?we_objectID=4251)
- 804 Isaac BH. 2004. *The limits of empire: the Roman army in the East*. Oxford: Clarendon Press
- 805 Jong D, C B, Hill PC, Aiken A, Awine T, Martin A, Adetifa IM, Jackson-Sillah DJ, Fox A,  
806 Kathryn D, et al. 2008. Progression to Active Tuberculosis, but Not Transmission, Varies  
807 by *Mycobacterium tuberculosis* Lineage in The Gambia. *J. Infect. Dis.* 198:1037–1043.
- 808 Kamvar ZN, Tabima JF, Grünwald NJ. 2014. Poppr: an R package for genetic analysis of  
809 populations with clonal, partially clonal, and/or sexual reproduction. *PeerJ* 2:e281.

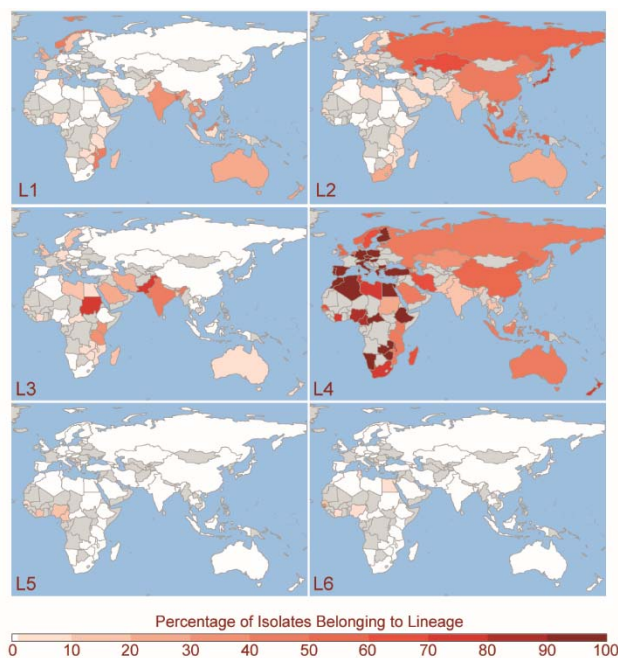
- 810 Kay GL, Sergeant MJ, Zhou Z, Chan JZ-M, Millard A, Quick J, Szikossy I, Pap I, Spigelman M,  
811 Loman NJ, et al. 2015. Eighteenth-century genomes show that mixed infections were  
812 common at time of peak tuberculosis in Europe. *Nat. Commun.* 6:6717.
- 813 Kent RK. 1979. The Possibilities of Indonesian Colonies in Africa with Reference to  
814 Madagascar. In: *Mouvements de Populations dans L’Ocean Indie*. Paris: H. Champion. p.  
815 93–105.
- 816 Lapierre M, Blin C, Lambert A, Achaz G, Rocha EPC. 2016. The Impact of Selection, Gene  
817 Conversion, and Biased Sampling on the Assessment of Microbial Demography. *Mol.*  
818 *Biol. Evol.* 33:1711–1725.
- 819 Leinonen R, Sugawara H, Shumway M. 2011. The Sequence Read Archive. *Nucleic Acids Res.*  
820 39:D19–D21.
- 821 Lemey P, Rambaut A, Drummond AJ, Suchard MA. 2009. Bayesian Phylogeography Finds Its  
822 Roots. *PLOS Comput. Biol.* 5:e1000520.
- 823 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
824 ArXiv13033997 Q-Bio [Internet]. Available from: <http://arxiv.org/abs/1303.3997>
- 825 Liu Q, Ma A, Wei L, Pang Y, Wu B, Luo T, Zhou Y, Zheng H-X, Jiang Q, Gan M, et al. 2018.  
826 China’s tuberculosis epidemic stems from historical expansion of four strains of  
827 *Mycobacterium tuberculosis*. *Nat. Ecol. Evol.* 2:1982.
- 828 Luo T, Comas I, Luo D, Lu B, Wu J, Wei L, Yang C, Liu Q, Gan M, Sun G, et al. 2015.  
829 Southern East Asian origin and coexpansion of *Mycobacterium tuberculosis* Beijing  
830 family with Han Chinese. *Proc. Natl. Acad. Sci.* 112:8136–8141.
- 831 Luttwak EN. 1976. *The grand strategy of the Roman Empire: from the first century A.D. to the*  
832 *third*. London: Weidenweld & Nicholson
- 833 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
834 *EMBnet.journal* 17:10–12.
- 835 Millar F. 1993. *The Roman Near East, 31 B.C.-A.D. 337*. Cambridge, Mass: Harvard University  
836 Press
- 837 Mortimer TD, Weber AM, Pepperell CS. 2018. Signatures of Selection at Drug Resistance Loci  
838 in *Mycobacterium tuberculosis*. *mSystems* 3:e00108-17.
- 839 Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin PR, O’Hara  
840 RB, Simpson GL, Solymos P, et al. 2017. *vegan: Community Ecology Package*.  
841 Available from: <https://cran.r-project.org/web/packages/vegan/index.html>
- 842 O’Neill MB, Mortimer TD, Pepperell CS. 2015. Diversity of *Mycobacterium tuberculosis* across  
843 Evolutionary Scales. *PLOS Pathog.* 11:e1005257.

- 844 Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-sites:  
845 rapid efficient extraction of SNPs from multi-FASTA alignments. Available from:  
846 <http://biorxiv.org/lookup/doi/10.1101/038190>
- 847 Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R  
848 language. *Bioinformatics* 20:289–290.
- 849 Parkin D, Barnes R eds. 2002. Ships and the development of maritime technology in the Indian  
850 Ocean. London: RoutledgeCurzon
- 851 Pepperell CS, Casto AM, Kitchen A, Granka JM, Cornejo OE, Holmes EC, Birren B, Galagan J,  
852 Feldman MW. 2013. The Role of Selection in Shaping Diversity of Natural M.  
853 tuberculosis Populations. *PLoS Pathog* 9:e1003543.
- 854 Pepperell CS, Granka JM, Alexander DC, Behr MA, Chui L, Gordon J, Guthrie JL, Jamieson  
855 FB, Langlois-Klassen D, Long R, et al. 2011. Dispersal of *Mycobacterium tuberculosis*  
856 via the Canadian fur trade. *Proc. Natl. Acad. Sci.* 108:6526–6531.
- 857 Pfister R, Bellinger L. 1945. The Excavations at Dura-Europos. Final Report IV: The Textiles.  
858 New Haven: Yale University Press
- 859 Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees  
860 for Large Alignments. *PLOS ONE* 5:e9490.
- 861 Pritchard JK, Stephens M, Donnelly P. 2000. Inference of Population Structure Using Multilocus  
862 Genotype Data. *Genetics* 155:945–959.
- 863 Public Health England. Tuberculosis in England: annual report - GOV.UK. Available from:  
864 <https://www.gov.uk/government/publications/tuberculosis-in-england-annual-report>
- 865 R Development Core Team. R: A Language and Environment for Statistical Computing. Vienna,  
866 Austria: R Foundation for Statistical Computing Available from: [http://www.R-](http://www.R-project.org/)  
867 [project.org/](http://www.R-project.org/)
- 868 Ramachandran S, Deshpande O, Roseman CC, Rosenberg NA, Feldman MW, Cavalli-Sforza  
869 LL. 2005. Support from the relationship of genetic and geographic distance in human  
870 populations for a serial founder effect originating in Africa. *Proc. Natl. Acad. Sci. U. S.*  
871 *A.* 102:15942–15947.
- 872 Rauh NK. 2003. Merchants, sailors and pirates in the Roman world. Stroud: Tempus
- 873 Ray HP. 2003. The archaeology of seafaring in ancient South Asia. Cambridge; New York:  
874 Cambridge University Press
- 875 Ray HP, Salles J-F, Institute of Southeast Asian Studies, Maison de l’Orient méditerranéen  
876 ancien (Lyon F, National Institute of Science, Technology and Development Studies,  
877 France, Ambassade (India), Centre for Human Sciences. 1996. Tradition and

- 878 archaeology: early maritime contacts in the Indian Ocean. New Delhi: Manohar  
879 Publishers
- 880 Salles J-F. 1996. Achaemenid and Hellenistic Trade in the Indian Ocean. In: The Indian Ocean in  
881 antiquity. p. 251–267. Available from:  
882 <http://public.eblib.com/choice/publicfullrecord.aspx?p=1517609>
- 883 Sartre M. 1991. L’Orient romain: provinces et sociétés provinciales en Méditerranée orientale  
884 d’Auguste aux Sévères (31 avant J.-C-235 après J.-C.). Paris: Seuil
- 885 Shabbeer A, Cowan LS, Ozcaglar C, Rastogi N, Vandenberg SL, Yener B, Bennett KP. 2012.  
886 TB-Lineage: An online tool for classification and analysis of strains of Mycobacterium  
887 tuberculosis complex. *Infect. Genet. Evol.* 12:789–797.
- 888 Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,  
889 Ideker T. 2003. Cytoscape: A Software Environment for Integrated Models of  
890 Biomolecular Interaction Networks. *Genome Res.* 13:2498–2504.
- 891 South A. 2016. rworldmap: Mapping Global Data. Available from: [https://cran.r-](https://cran.r-project.org/web/packages/rworldmap/index.html)  
892 [project.org/web/packages/rworldmap/index.html](https://cran.r-project.org/web/packages/rworldmap/index.html)
- 893 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
894 large phylogenies. *Bioinformatics* 30:1312–1313.
- 895 Thorpe HA, Bayliss SC, Hurst LD, Feil EJ. 2017. Comparative Analyses of Selection Operating  
896 on Nontranslated Intergenic Regions of Diverse Bacterial Species. *Genetics* 206:363–  
897 376.
- 898 Vogt B. 1996. Bronze Age Maritime Trade in the Indian Ocean: Harappan Traits on the Oman  
899 Peninsula. In: Reade J, editor. The Indian Ocean in antiquity. p. 107–132. Available  
900 from: <http://public.eblib.com/choice/publicfullrecord.aspx?p=1517609>
- 901 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,  
902 Wortman J, Young SK, et al. 2014. Pilon: An Integrated Tool for Comprehensive  
903 Microbial Variant Detection and Genome Assembly Improvement. *PLOS ONE*  
904 9:e112963.
- 905 White Z, Painter J, Douglas P, Abubakar I, Njoo H, Archibald C, Halverson J, Robson J, Posey  
906 DL. 2017. Immigrant Arrival and Tuberculosis among Large Immigrant- and Refugee-  
907 Receiving Countries, 2005–2013;2009. *Tuberc. Res. Treat.* [Internet]. Available from:  
908 <https://www.hindawi.com/journals/trt/2017/8567893/>
- 909 Wink A. 2002. From the Mediterranean to the Indian Ocean: Medieval History in Geographic  
910 Perspective. *Comp. Stud. Soc. Hist.* 44:416–445.
- 911 World Health Organization. 2017. Global Tuberculosis Report 2017. World Health Organization

- 912 Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. 2017. ggtree: an r package for visualization and  
913 annotation of phylogenetic trees with their covariates and other associated data. *Methods*  
914 *Ecol. Evol.* 8:28–36.
- 915 Zarins J. 1996. Obsidian in the Larger Context of Predynastic/Archaic Egyptian Red Sea Trade.  
916 In: Reade J, editor. *The Indian Ocean in antiquity*. p. 107–132. Available from:  
917 <http://public.eblib.com/choice/publicfullrecord.aspx?p=1517609>
- 918 □ Ālam M, Subrahmanyam S. 2009. *Indo-Persian travels in the age of discoveries, 1400-1800*.  
919 Digit. pr. Cambridge: Cambridge University Press
- 920

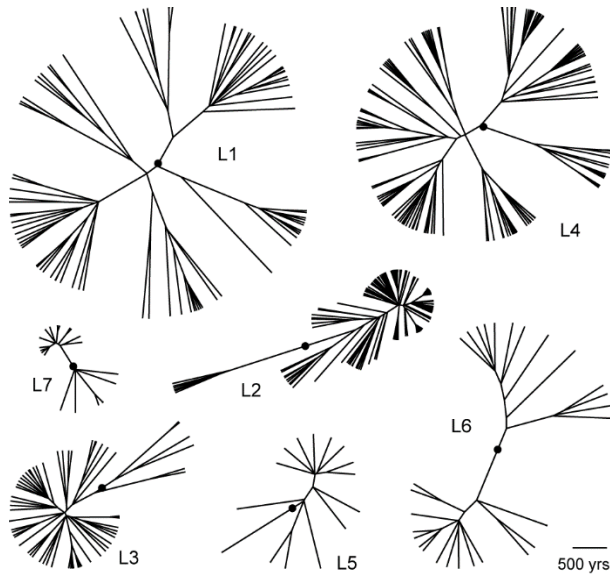
921 Figures  
922



923  
924

925 **Fig. 1.** Geographic distributions of *Mycobacterium tuberculosis* lineages 1-6. Spoligotypes from  
926 the SITVIT WEB database ( $n = 42,358$ ) were assigned to lineages 1-6. Countries are colored  
927 from white to dark-red based on the percentage of isolates from the country belonging to each  
928 lineage. Unsampled countries and those with less than 10 isolates in the database are shown in  
929 grey. Lineage 7 (not pictured) is found exclusively in Ethiopia.

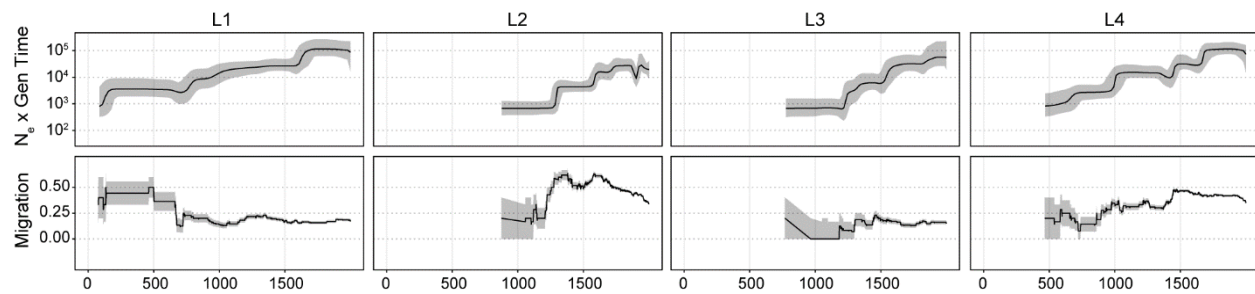




930  
931

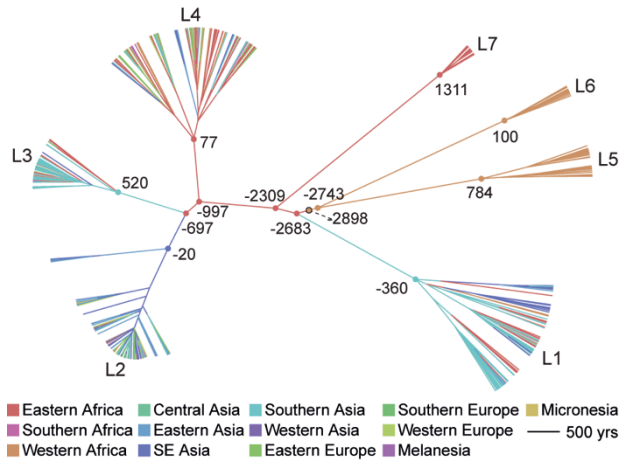
932 **Fig. 2.** Maximum clade credibility phylogenies of *Mycobacterium tuberculosis* lineages 1-6.  
933 Bayesian analyses were performed on each lineage alignment with the general time reversible  
934 model of nucleotide substitution with a gamma distribution to account for rate heterogeneity  
935 between sites, a strict molecular clock, and Bayesian skyline plot demographic models. The  
936 most recent common ancestor (MRCA) of each lineage is indicated with a black circle; the  
937 MRCA of individual lineage phylogenies were informed by the phylogeny of the entire Old  
938 World collection, which was dated using a substitution rate of  $5 \times 10^{-8}$  substitutions/site/year  
939 (Kay et al. 2015).



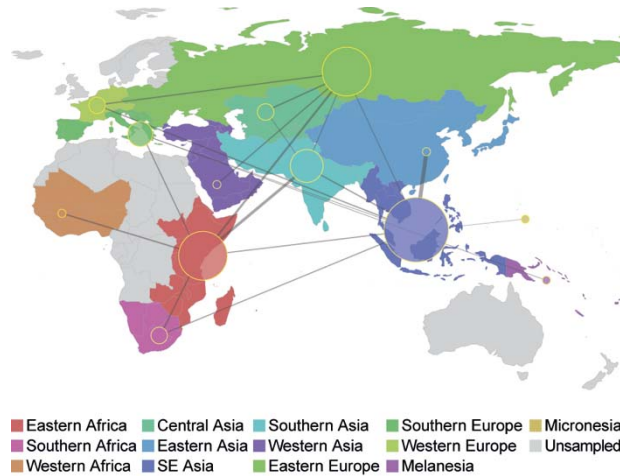


940  
941

942 **Fig. 3.** Patterns of effective population size and migration through time of *Mycobacterium*  
943 *tuberculosis* lineages 1-4. Bayesian skyline plots (top panels) show inferred changes in effective  
944 population size ( $N_e$ ) through time deduced from lineage specific analyses. Black lines denote  
945 median  $N_e$  and gray shading the 95% highest posterior density. Estimated migration through  
946 time (see *Methods*) for each lineage is shown in the bottom panels (see *Methods*). Grey shading  
947 depicts the rates inferred after the addition or subtraction of a single migration event, and  
948 demonstrate the uncertainty of rate estimates, particularly from the early history of each lineage.  
949 Dates are shown in calendar years and are based on scaling the phylogeny of the Old World  
950 collection with a substitution rate of  $5 \times 10^{-8}$  substitutions/site/year (Kay et al. 2015).

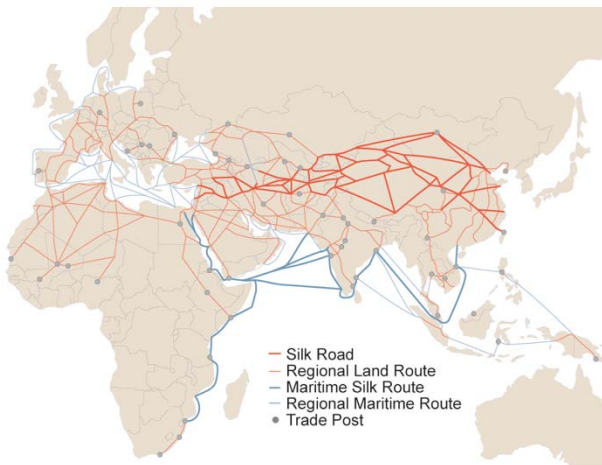


**Fig. 4.** Maximum clade credibility tree of the Old World Collection. Estimated divergence dates are shown in calendar years based on median heights and a substitution rate of  $5 \times 10^{-8}$  substitutions/site/year (Kay et al. 2015). Branches are colored according to the inferred most probable geographic origin. Nodes corresponding to the most recent common ancestors (MRCA) of each lineage, lineage splits, and the MRCA of *M. tuberculosis* (outlined black) are marked with circles and colored to reflect their most probable geographic origin.



959  
960

961 **Fig. 5.** Connectivity of UN subregions during dispersal of *Mycobacterium tuberculosis*. The  
962 Bayesian stochastic search variable selection method was used to identify and quantify  
963 migrations with strong support in discrete phylogeographic analysis of the Old World collection.  
964 Node sizes reflect the number of significant migrations emanating from the region observed in  
965 the phylogeny, whereas the thickness of lines connecting regions reflects the estimated relative  
966 rate between regions.



967  
968 **Fig. 6.** Trade routes active throughout Europe, Africa and Asia by 1400 CE. Nodes (trade cities,  
969 oases, and caravanserai) and arcs (the routes between nodes) are from the Old World Trade  
970 Routes Project ([www.ciolek.com/owtrad.html](http://www.ciolek.com/owtrad.html), accessed February 17, 2016) and are visualized  
971 with ArcGIS.

972 **Table 1.** Genetic diversity of Old World *M.tb* across lineages 1-7. TMRCA estimates reflect  
 973 scaling of results to evolutionary rates calibrated from ancient DNA [median  $5.00 \times 10^{-8}$   
 974 substitutions/ site/ year (Kay et al. 2015) and are written as calendar years. To account for  
 975 uncertainty in this rate estimate, our lower and upper TMRCA estimates reflect scaling of our  
 976 results with the low and high bounds of the 95% highest posterior density estimates of the rate  
 977 reported from ancient DNA analysis (i.e.  $4.06 \times 10^{-8}$  and  $5.87 \times 10^{-8}$ , respectively).

		<i>MTBC</i>	<i>L1</i>	<i>L4</i>	<i>L2</i>	<i>L3</i>	<i>L5</i>	<i>L6</i>	<i>L7</i>
<i>Sample</i>	<i>n</i>	552	89	143	181	65	15	31	28
<i>Diversity</i>	□	2.13E-03	7.56E-04	7.80E-04	4.49E-04	3.88E-04	1.72E-04	3.04E-04	7.99E-05
	π	2.80E-04	1.92E-04	1.54E-04	7.46E-05	9.16E-05	8.77E-05	1.41E-04	4.52E-05
<i>Demographic Inference</i>	<i>N/Nanc</i>	91 ± 4	71 ± 5	55 ± 22	112 ± 102	148 ± 2	504 ± 111	50 ± 5	17 ± 4
	<i>Generations (Nanc)</i>	0.16 ± 0.01	0.80 ± 0.06	0.65 ± 0.35	0.41 ± 0.94	3.54 ± 0.04	3.94 ± 0.73	1.10 ± 0.09	2.45 ± 0.89
	<i>LL expansion</i>	-1788.4	-424.2	-492.8	-467.1	-108.2	-42.4	-151.9	-64.5
	<i>LL neutral</i>	-10549.2	-3246.6	-3474.6	-2378.9	-1717.0	-520.7	-912.3	-159.4
	<i>p-value</i>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
<i>Structure UN subregions</i>	<i>Var. Between</i>	21	19	4	20	16	NA	NA	NA
	<i>Var. Within</i>	79	81	96	80	84	NA	NA	NA
	<i>p-value</i>	<0.001	<0.001	0.001	<0.001	0.004	NA	NA	NA
<i>Structure Botanical Continents</i>	<i>Var. Between</i>	14	5	2	9	13	NA	NA	NA
	<i>Var. Within</i>	86	95	98	91	87	NA	NA	NA
	<i>p-value</i>	<0.001	0.02	0.05	<0.001	<0.001	NA	NA	NA
<i>TMRCA</i>	<i>median</i>	-2898	-360	77	-20	520	784	100	1311
	<i>lower</i>	-4032	-906	-368	-488	177	502	-339	1152
	<i>upper</i>	-2172	-10	362	279	739	964	382	1413
<i>Geographic origin</i>	<i>1st region probability</i>	W Africa	S Asia	E Africa	SE Asia	S Asia	W Africa	W Africa	E Africa
		54.2%	75.6%	98.9%	81.0%	63.5%	99.9%	99.8%	99.8%
	<i>2nd region probability</i>	E Africa	E Africa	E Europe	E Asia	E Africa	E Africa	E Africa	S Africa
	37.5%	24.1%	0.7%	9.2%	36.2%	0.1%	0.2%	0.0%	

978