1 <u>Title</u>
2 Lineage specific histories of *Mycobacterium tuberculosis* dispersal in Africa and Eurasia
3
4 <u>Author Affiliation</u>
5 Mary B O'Neill[a,b,*], Abigail Shockey[b], Alex Zarley[c], William Aylward[d], Vegard Eldholm[e], Andrew
6 Kitchen[f], Caitlin S Pepperell[b,g]
7
8 [a]Laboratory of Genetics, University of Wisconsin-Madison, Madison, WI 53706, USA
9 [b]Department of Medical Microbiology and Immunology, University of Wisconsin-Madison,
10 Madison, WI 53706, USA
11 [c]Department of Geography, University of Wisconsin-Madison, WI 53706, USA
12 [d]Department of Classical and Ancient Near Eastern Studies, University of Wisconsin-Madison,
13 Madison, WI 53706, USA
14 [e]Infection Control and Environmental Health, Norwegian Institute of Public Health, 0456 Oslo,
15 Norway
16 [f]Department of Anthropology, University of Iowa, Iowa City, IA 52242, USA
17 [g]Department of Medicine, University of Wisconsin-Madison, Madison, WI 53706, USA
18 *Present address: Unit of Human Evolutionary Genetics, Institut Pasteur, 75015 Paris, France
19
20 <u>Corresponding Authors</u>
21 Caitlin S Pepperell
22 1550 Linden Drive
23 5301 Microbial Sciences Building
24 Madison, WI 53706
25 (608) 262-5983
26 cspepper@medicine.wisc.edu
27
28 Andrew Kitchen
29 17 N. Clinton Street
30 114 Macbride Hall
31 Iowa City, Iowa 52242
32 (319) 335-2891
33 andrew-kitchen@uiowa.edu
34
35 <u>Keywords</u>
36 phylogeography, evolution, pathogen, migration, demography
37

38    Abstract

39    *Mycobacterium tuberculosis* (*M.tb*) is a globally distributed, obligate pathogen of humans that

40    can be divided into seven clearly defined lineages.  Identifying how the ancestral clone of *M.tb*

41    spread and differentiated is important for identifying the ecological drivers of the current

42    pandemic.  We reconstructed *M.tb* migration in Africa and Eurasia, and investigated lineage

43    specific patterns of spread.  Applying evolutionary rates inferred with ancient *M.tb* genome

44    calibration, we link *M.tb* dispersal to historical phenomena that altered patterns of connectivity

45    throughout Africa and Eurasia: trans-Indian Ocean trade in spices and other goods, the Silk

46    Road and its predecessors, the expansion of the Roman Empire and, the European Age of

47    Exploration.  We find that Eastern Africa and Southeast Asia have been critical in the dispersal

48    of *M.tb*.  Our results reveal complex relationships between spatial dispersal and expansion of

49    *M.tb* populations, and delineate the independent evolutionary trajectories of bacterial sub-

50    populations underlying the current pandemic.

51

52    Introduction

53    The history of tuberculosis (TB) has been rewritten several times as genetic data accumulate

54    from its causative agent, *Mycobacterium tuberculosis* (*M.tb*).  In the nascent genomic era, these

55    data refuted the long-held hypothesis that human-adapted *M.tb* emerged from an animal

56    adapted genetic background represented among extant bacteria by *Mycobacterium bovis*,

57    another member of the *Mycobacterium tuberculosis* complex (MTBC) (Brosch et al. 2002).

58    Genetic data from bacteria infecting multiple species of hosts revealed that currently known

59    non-primate-adapted strains form a nested clade within the diversity of extant *M.tb* (Behr et al.

60    1999; Brosch et al. 2002; Hershberg et al. 2008).

61

62    *M.tb* can be classified into seven well-differentiated lineages, which differ in their geographic

63    distribution and association with human sub-populations (Hirsh et al. 2004; Gagneux et al.

64    2006).  This observation led to the hypothesis that *M.tb* diversity has been shaped by human

65    migrations out of Africa, and that the most recent common ancestor (MRCA) of extant *M.tb*

66    emerged in Africa approximately 73,000 years ago, coincident with estimated waves of human

67    migration (Comas et al. 2013).  Human out of Africa migrations are a plausible means by which

68    *M.tb* could have spread globally.  However, several features of *M.tb* population genetics suggest

69    it has diversified over relatively short time scales: for example, the species is characterized by

70    low genetic diversity (Eldholm and Balloux 2016) and high rates of non-synonymous

71    polymorphism (Rocha et al. 2006).

72

73    The observation of limited diversity among extant *M.tb* could be reconciled with the out of Africa

74    scenario if *M.tb*'s rate of evolution were orders of magnitude lower than estimates from other

75    bacterial pathogens, or if *M.tb* exhibited dramatic rate decay such that substitution rates varied

76    by several orders of magnitude as a function of temporal sampling window (Comas et al. 2013;

77    Eldholm and Balloux 2016).  There are at least nine published estimates of the rate of *M.tb*

78    molecular evolution, which use a variety of calibration methods (Eldholm et al. 2016).  Rate

79    estimates calibrated with sampling dates, historical events, experimental infection in non-human

80    primates, recent transmission events, and ancient DNA are concordant.  Critically, *M.tb* rate

81    estimates are similar to those of other bacterial species, and are inconsistent with the out of

82    Africa hypothesis.  In addition, a recent meta-analysis of evolutionary rate estimates in bacteria

83    noted that 'reliable rate estimates for *M. tuberculosis*, estimated over sampling frames of 15 and

84    895 years, were nearly identical' (Duchêne et al. 2016), suggesting that *M.tb* is not

85    characterized by the dramatic rate decay that would be needed to reconcile observed data with

86    the out of Africa hypothesis.

87

88    When calibrated with ancient DNA, the estimates of the time to most recent common ancestor

89    (TMRCA) for the MTBC are <6,000 years before present (Bos et al. 2014; Kay et al. 2015).

90    This is not necessarily the time period over which TB first emerged, as it is possible –

91    particularly given the apparent absence of recombination among *M.tb* (Pepperell et al. 2013) –

92    that the global population has undergone clonal replacement events that displaced ancient

93    diversity from the species.

94

95    *M.tb* is an obligate pathogen of humans with a global geographic range.  The finding of a recent

96    origin for the extant *M.tb* population raises the question of how the organism could have spread

97    within this timeframe to occupy its current distribution.  *M.tb* populations in the Americas show

98    the impacts of European colonial movements as well as recent immigration (e.g. Pepperell et al.

99    2011); the role of other historical phenomena in driving TB dispersal is not well understood.

100   Here we sought to reconstruct the migratory history of *M.tb* populations in Africa and Eurasia

101   within the newly established framework of a recent origin and evolutionary rates derived from

102   ancient DNA data (Bos et al. 2014; Kay et al. 2015).  We discovered lineage-specific patterns of

103   migration and a complex relationship between *M.tb* effective population growth and migration.

104   Our results connect *M.tb* migration to major historical events in human history that altered

105   patterns of connectivity in Africa and Eurasia.  These findings provide context for a recent

106     evolutionary origin of the MRCA of *M.tb* (Pepperell et al. 2013; Bos et al. 2014; Kay et al. 2015),

107     which represents yet another paradigm shift in our understanding of the history and origin of this

108     successful pathogen.

109

110     <u>Results</u>

111     ***Genetic and geographic structures of global M.tb populations***

112     In order to establish the contemporary geographic distributions of *M.tb* lineages, we translated

113     the spoligotypes reported for 42,358 *M.tb* isolates to their corresponding lineage designations

114     (fig. 1).  Geographic patterns in prevalence vary between lineages.  Lineage 1 (L1) is prevalent

115     in regions bordering the Indian Ocean, extending from Eastern Africa to Melanesia.  Lineage 2

116     (L2) is broadly distributed, with a predominance in Eastern Eurasia and South East Asia.

117     Lineage 3 (L3) is similar to L1 in that its distribution rings the Indian Ocean, but it does not

118     extend into Southeastern Asia, it has a stronger presence in Northern Africa, and a broader

119     distribution across Southern Asia.  Lineage 4 (L4) is strikingly well dispersed, with a

120     predominance throughout Africa and Europe and the entire region bordering the Mediterranean.

121     Lineages 5 (L5) and 6 (L6) are found at low frequencies in Western and Northern Africa.

122     Lineage 7 (L7), as previously described (Blouin et al. 2012; Firdessa et al. 2013; Comas et al.

123     2015), is limited to Ethiopia.

124

125     We compiled a diverse collection of *M.tb* genomes for phylogenetic and population genetic

126     inference of the demographic and migratory history of the extant *M.tb* population (*see Methods*).

127     Our dataset consists of whole-genome sequences (WGS) from 552 *M.tb* isolates collected from

128     51 countries (spanning 13 UN geoscheme subregions), which we refer to as the Old World

129     collection (fig. s1, table s1).  We included sites in the alignment where at least half of these

130     isolates had confident data (60,787 variant sites; 3,838,249 bp) for subsequent analyses, unless

131     otherwise noted.

132

133     The inferred maximum likelihood phylogeny and Bayesian clustering analysis reveals the well

134     described *M.tb* lineage structure, and some associations are evident between lineages and

135     geographic regions (defined here by the United Nations geoscheme) (fig. s2).  The phylogeny

136     has an unbalanced shape, with long internal branches that define the lineages and feathery tips,

137     suggestive of recent population expansion.

138

139   Genetic diversity, as measured by the numbers of segregating sites and pairwise differences

140   (Watterson's Ɵ and π), varied among lineages (table 1).  L1 and L4 group together and have

141   the highest diversity; L2, L3, L5, and L6 have similar levels of diversity and form the middle

142   grouping; L7 has the lowest diversity.  We used an analysis of molecular variance (AMOVA) to

143   delineate the effects of population sub-division on *M.tb* diversity (table 1).  The Old World

144   collection was highly structured among UN subregions (21% of variation attributable to

145   between-region comparisons), whereas this structure was less apparent when regions were

146   defined by the botanical contents outlined by the World geographic scheme for recording plant

147   distributions (14%).  This is consistent with *M.tb*'s niche as an obligate human pathogen, with

148   bacterial population structure directly shaped by that of its host population (i.e. reflected in UN

149   subregions) rather than climatic and other environmental features (reflected in botanical

150   continent definitions).  We obtained similar results when the lineages were considered

151   separately, except for L4, which had little evidence of population structure (4% variation among

152   UN subregions, 2% among botanical continents).

153

154   ***Distinct demographic histories of the M.tb lineages***

155   Bayesian inferred trees vary among lineages (fig. 2), likely reflecting their distinct demographic

156   histories.  Branch lengths are relatively even across the phylogenies of L1 and L4, whereas L2

157   and L3 have a less balanced structure.  The long, sparse internal branches and radiating tips of

158   L2 and L3 phylogenies are consistent with an early history during which the effective population

159   size remained small (and diversity was lost to drift), followed by more recent population

160   expansion.  L5 has a star-like structure, consistent with rapid population expansion.  Jointly

161   inferred Bayesian skyline plot (BSP) reconstructions of effective population sizes over time

162   suggest that lineages 1-6 have undergone expansion (fig. 3 – top panel, fig. S3).  We estimate

163   that L2 and L3 underwent abrupt expansion at approximately the same time, whereas

164   expansions of L1 and L4 appeared relatively smooth.

165

166   We used the methods implemented in ∂a∂i to reconstruct the demographic histories of each

167   *M.tb* population (i.e. lineage) from its synonymous site frequency spectrum (SFS).  As

168   demographic inference with ∂a∂i is sensitive to missing data, loci at which any sequence in the

169   individual lineage alignments had a gap or unknown character were removed for these

170   analyses.  Consistent with the BSP analyses performed in BEAST, instantaneous expansion

171   and exponential growth models offered an improved fit to the data in comparison with the

172   constant population size model for each lineage and the entire Old World collection (fig. S4).

5

173    Parameter estimates varied widely across runs for the exponential growth model, so we report

174    results only for the instantaneous expansion model (table 1).

175

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

177    There was evidence of isolation by distance in the global *M.tb* population, as assessed with a

178    Mantel test of correlations between genetic and geographic distances.  We defined geographic

179    distances using three schemes: great circle distances, great circle distances through waypoints

180    of human migration as described in (Ramachandran et al. 2005), and distances along historical

181    trade routes.  Waypoints are used to make distance estimates more reflective of presumed

182    human migration patterns (i.e., when calculating between-continent distances, it is generally

183    thought that humans did not pass through large bodies of water, and thus a waypoint is used).

184    To allow comparisons between the schemes, values were centered and standardized (*see*

185    *Methods*).  Values of the Mantel test statistic were similar for great circle distances (r = 0.16)

186    and trade network distances (r = 0.16), with distances through waypoints reflective of human

187    migration patterns having a lower value (r = 0.14, *p* = 0.0001 for all three analyses).  In analyses

188    of human genetic data, adjustment of great circle distances with waypoints results in a higher

189    correlation between genetic and geographic distances (Ramachandran et al. 2005).  Our Mantel

190    test results therefore do not support a pattern of isolation by distance as expected if out of Africa

191    human migrations were the primary influence on global diversity of extant *M.tb* (Comas et al.

192    2013).

193

194    To further investigate a potential influence of ancient human migration on *M.tb* evolution, we

195    calculated the correlation between *M.tb* genetic diversity (π) within subregions and their

196    average distances from Addis Ababa, a proxy for a possible origin of anatomically modern

197    human expansion out of Africa.  Contrary to what is observed for human population diversity

198    (Ramachandran et al. 2005), we did not observe a significant decline in *M.tb* diversity as a

199    function of distance in our Old World collection (adjusted R-squared =  -0.1, *p* = 0.88), nor when

200    we included samples from the Americas (adjusted R-squared = $8.9 \times 10^{-4}$, *p* = 0.34, fig. S5,

201    table S2).

202

203    We used the methods implemented in BEAST to reconstruct the migratory history of the entire

204    Old World *M.tb* collection as well as individual lineages within it, modelling geographic origin of

205    isolates (UN subregion or country) as a discrete trait (fig. 4, figs. S6-S10).  Using an

206    evolutionary rate calibrated with 18[th] century *M.tb* DNA of $5 \times 10^{-8}$ substitutions/site/year (Kay et

6

207    al. 2015), which is similar to the rate inferred with data from 1,000 year old specimens (Bos et

208    al. 2014), our estimate of the time to most recent common ancestor for extant *M.tb* is between

209    4032 BCE and 2172 BCE (table 1; date ranges are based on the upper and lower limits of the

210    95% highest posterior density (HPD) for the rate reported in Kay et al. (2015) which is more

211    conservative than the 95% HPD of our model).  We infer an African origin for the MRCA

212    (Eastern or Western subregion, table 1, fig. 4, fig. S6).  Shortly after emergence of the common

213    ancestor, we infer a migration of the L1-L2-L3-L4-L7 ancestral lineage from Western to Eastern

214    Africa (we estimate prior to 2683 BCE), with subsequent migrations occurring out of Eastern

215    Africa.

216

217    In our phylogeographic reconstruction, emergence of L1 follows migration from Eastern Africa to

218    Southern Asia at some time between the 3rd millennium and 4th century BCE (table 1, fig. 4, fig.

219    S6).  L1 has an 'out of India' phylogeographic pattern (fig. S7), with diverse Indian lineages

220    interspersed throughout the phylogeny.  This suggests that the current distribution of L1 around

221    the Indian Ocean (fig. 1) arose from migrations out of India, from a pool of bacterial lineages

222    that diversified following migration from Eastern Africa.

223

224    The phylogeographic reconstruction further indicates that following the divergence of L1, *M.tb*

225    continued to diversify in Eastern Africa, with emergence of L7 there, followed by L4 (table 1, fig.

226    4, fig. S6).  The contemporary distribution of L4 is extremely broad (fig. 1) and in this analysis of

227    the Old World collection we infer an East African location for the internal branches of L4.

228    Notably, in the lineage-specific analyses, we infer a European location for these branches (fig.

229    S7).  The difference is likely due to the fact that inference is informed by deeper as well as

230    descendant nodes in the Old World collection.  Together, these results imply close ties between

231    Europe and Africa during the early history of this lineage that we estimate emerged in the 1st

232    century CE (368 BCE-362 CE, table 1).

233

234    After the emergence of L1 and L7 from Eastern Africa, our analyses suggest that a migration

235    occurring between 697 BCE and 520 CE established L3 in Southern Asia, with subsequent

236    dispersal out of Southern Asia into its present distribution, which includes Eastern Africa (i.e., a

237    back migration of L3 to Africa, fig. 1).  We estimate that L2 diversified in South Eastern Asia

238    following migration from Eastern Africa at some point between 697 BCE and 20 BCE (table 1,

239    fig. 4, fig. S6).  A previously published analysis of L2 phylogeography also inferred a Southeast

240    Asian origin for the lineage (Luo et al. 2015).

241

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

243  Our phylogeographic reconstruction indicated that temporal trends in migration varied among

244  lineages (fig. 3 – bottom panel).  We infer that L1 was characterized by high levels of migration

245  until approximately the 7$^{th}$ century CE, when the rate of migration decreased abruptly and

246  remained stable thereafter.  L3, by contrast, exhibited consistently low rates of migration.  L2

247  and L4 had more variable trends in migration, as each underwent punctuated increases in

248  migration rate.  Temporal trends in growth and migration are congruent for L2 and L4, with

249  increases in migration rate preceding effective population expansions; this is not the case for L1

250  and L3.  Taken together, these results suggest that L1 and L3 populations (as well as L5 and

251  L6, fig. S3b) grew *in situ*, whereas range expansion may have contributed to the growth of L2

252  and L4.

253

254  We employed the Bayesian stochastic search variable selection method (BSSVS) in BEAST

255  (Lemey et al. 2009) to estimate relative migration rates within the most parsimonious migration

256  matrix.  A map showing inferred patterns of connectivity among UN subregions and relative

257  rates of *M.tb* migration with strong posterior support is shown in fig. 5.  South Eastern Asia was

258  the most connected region in our analyses, with significant rates of migration connecting it to

259  eight other regions.  Eastern Africa, Eastern Europe, and Southern Asia were also highly

260  connected, with significant rates with six, six, and five other regions, respectively.  Western

261  Africa, Eastern Asia, and Western Asia were the least connected regions, with just one

262  significant connection each (to Eastern Africa, South Eastern Asia, and Eastern Europe,

263  respectively).  Our sample from Western Asia is, however, limited (table S1) and migration from

264  this region may have consequently been underestimated.  The highest rates of migration were

265  seen between Eastern Asia and Southeastern Asia, and between Eastern Africa and Southern

266  Asia.

267

268  Lineage specific analyses suggest that migration between Southern Asia, Eastern Africa, and

269  South Eastern Asia has been important for the dispersal of L1, whereas South Eastern Asia and

270  Eastern Europe have been important for L2 (fig. S11).  L3 is similar to L1 in that there is

271  evidence of relatively high rates of migration between Southern Asia and Eastern Africa.  There

272  is also evidence of migration within Africa between the eastern and southern subregions.  In the

273  analyses of migration for L4, Eastern Africa appeared highly connected with other regions.

274

275  ***Phylogeographic reconstruction: limitations and alternatives***

276  These phylogeographic reconstructions are clearly sensitive to sampling, since we cannot

277  identify the roles of unsampled regions in *M.tb*'s migratory history.  We maximized geographic

278  diversity in our sample, but were limited by available data and some regions – notably Middle

279  Africa, Northern Africa, and Western Asia – are absent or underrepresented in our sample (fig.

280  S1).  Defining the contributions of these undersampled regions to *M.tb*'s migratory history awaits

281  more samples and/or further method development.

282

283  De Maio *et al*. (2015) note the sensitivity of discrete trait phylogeographic inference in BEAST to

284  sample selection, as well as overconfidence in the precision of geographic inference, and

285  propose BASTA as an alternative (De Maio et al. 2015).  BASTA is sensitive to the choice of

286  prior and we did not have ancillary data to guide the selection of a prior for the Old World

287  migratory history of *M.tb*, precluding its use here.  We investigated ∂a∂i as an alternative tool for

288  phylogeographic inference but it did not perform well for this application under conditions of

289  complete linkage of sites (Note S1, fig. S12, fig. S13, table S3, table S4).  Lapierre *et al*. (2016)

290  investigated the sensitivity of BEAST inference of demography to sampling regime.  They found

291  that the method performed poorly under a 'uniform' sub-sampling regime, in which populations

292  are randomly sampled to the same size (Lapierre et al. 2016).  Our results are concordant with

293  Lapierre *et al.* (2016), in that we found inferred migration matrices were not consistent across

294  random sub-samples of our dataset (fig. S14).  We also interrogated the relationship between

295  regional sample size and inferred migration rate and did not observe a strong correlation (fig.

296  S14).  The phylogeographic inference method implemented here relies on the assumption that

297  sample size reflects deme size (Lemey et al. 2009; De Maio et al. 2015), and within the

298  constraints of available data, we attempted to adjust our sample sizes according the regional

299  prevalence of TB (see *Methods* and fig. S1).  According to the classifications proposed by

300  Lapierre *et al*. (2015), our Old World collection represents a 'mixed' sampling scheme (see

301  *Methods*).

302

303  We previously demonstrated effects of population expansion, linkage, and purifying selection on

304  *M.tb* genetic diversity (Pepperell et al. 2013).  Given these previous observations, we were

305  curious about a potential impact of purifying selection on inference of migration.  To address this

306  question, we simulated data under demographic models with and without selection and

307  migration, and then analyzed the resulting sequence alignments in BEAST.  Our two population

308  simulation suggests that purifying selection may elevate estimated migration rates, though the

309    distribution of mean rate estimates for simulations with and without purifying selection broadly

310    overlap (fig. S15).  However, analysis of sequence alignments generated under a three

311    population model suggested that selection had a statistically negligible effect on migration rates,

312    which can be observed from plots of the mean relative rates (fig. S16) or of the relative support

313    of migration rates (fig. S17).  We note that the discrete migration model implemented in BEAST

314    was able to capture much of the asymmetry of our three population asymmetrical simulations as

315    evidenced by the distribution of relative migration rates and Bayes factor (BF) support for said

316    rates.  BEAST also consistently produced similar BF support for rates estimated from data

317    simulated under symmetrical migration models (i.e., those with global $M$ = 0.5 or 0.0).  Our

318    simulations thus suggest that consistent purifying selection is unlikely to dramatically affect

319    estimates of, or support for, migration rates between populations in these scenarios.

320

321    Discussion

322    Our reconstructions of *M.tb* dispersal throughout the Old World delineate a complex migratory

323    history that varies substantially between bacterial lineages.  Patterns of diversity among extant

324    *M.tb* suggest that historical pathogen populations were capable of moving fluidly over vast

325    distances.  Using evolutionary rate estimates from ancient DNA calibration, we time the

326    dispersal of *M.tb* to a historical period of exploration, trade, and increased connectivity among

327    regions of the Old World.

328

329    Consistent with prior reports (Comas et al. 2013), we infer an origin of *M.tb* on the African

330    continent (table 1, fig. 4, fig. S6).  There is a modest preference for Western Africa over Eastern

331    Africa (54% versus 38% inferred probability), likely due to the early branching West African

332    lineages (i.e. *Mycobacterium africanum*, L5 and L6).  Larger samples may allow more precise

333    localization of the *M.tb* MRCA, and Northern Africa in particular is under-studied.

334

335    We infer L1 to be the first lineage that emerged out of Africa; L1 is currently concentrated in

336    regions bordering the Indian Ocean from Eastern Africa to Melanesia (fig. 1).  In our

337    phylogeographic reconstruction, the genesis of this lineage traces to migration from Eastern

338    Africa to Southern Asia at some point between the 3$^{rd}$ millennium and 4$^{th}$ century BCE, with

339    subsequent dispersal occurring out of the Indian subcontinent.  Our results suggest that the

340    early history of L1 was characterized by high levels of migration, particularly between Southern

341    Asia and Eastern Africa, and between Southern Asia and South Eastern Asia (fig. 3, fig. S11).

342    The geographic distribution of L1, the timing of its emergence and spread, as well as patterns of

343    connectivity underlying its dispersal, are all consistent with migration via established trans-

344    Indian Ocean trade routes linking Eastern Africa to Southern and South Eastern Asia (fig. 6).

345    The interval of our timing estimate for the initial migration overlaps with the so-called Middle

346    Asian Interaction sphere in The Age of Integration (2600-1900 BCE), which is marked by

347    increased cultural exchange and trade between civilizations of Egypt, Mesopotamia, the Arabian

348    peninsula, and the Indus Valley (Vogt 1996; Zarins 1996; Parkin and Barnes 2002; Ray 2003;

349    Coningham and Young 2015).  East-West contact and trade across the Indian Ocean intensified

350    in the first millennium BCE, when maritime networks expanded to include the eastern

351    Mediterranean, the Red Sea, and the Black Sea (Dilke 1985; Boussac et al. 1995; Ray et al.

352    1996; Salles 1996).  Historical data from the Roman era indicate that crews on trading ships

353    crossing the Indian Ocean comprised fluid assemblages of individuals from diverse regions,

354    brought together under conditions favorable for the transmission of TB (André and Filliozat

355    1986; Begley and De Puma 1991; Wink 2002; Rauh 2003).  These ships would have been an

356    efficient means of spreading *M.tb* among the distant regions involved in trade.

357

358    L2 may similarly have an origin in East-West maritime trade across the Indian Ocean, as we

359    infer it arose from a migration event from Eastern Africa to South Eastern Asia during the 1st

360    millennium BCE.  In this era, increased sophistication in ship technology allowed for longer

361    voyages (Kent 1979; Blench 1996; Ray et al. 1996; Parkin and Barnes 2002; Wink 2002; Ray

362    2003).  L2 appears to have spread out of Southeast Asia, a highly connected region in our

363    analyses of *M.tb* migration, and is currently found across Eastern Eurasia and throughout South

364    Eastern Asia (fig. 1, fig. 4, fig. S6, fig. S11).  Interestingly, although L2 is dominant in Eastern

365    Asia, the region did not appear to have played a prominent role in dispersal of this lineage,

366    except in its exchanges with South Eastern Asia.

367

368    In contrast to L1 and L2, L3 appears to have had relatively low rates of migration throughout its

369    history (fig. 3).  The contemporary geographic range of L3 is also narrower, extending east from

370    Northern Africa through Western Asia to the Indian subcontinent (fig. 1).  A study of lineage

371    prevalence in Ethiopia showed that L3 is currently concentrated in the north of the country

372    (Comas et al. 2015), consistent with our observed north to south gradient in its distribution on

373    the African continent.  This is in opposition to L1, which has a southern predominance in

374    Ethiopia and across Eastern Africa (fig. 1).  We estimate L3 emerged in Southern Asia ca. 520

375    CE (177-739 CE).  Pakistan harbors diverse strains belonging to L3 (fig. S9), and the Southern

376    Asia region was highly connected with Eastern Africa in our analyses (fig. S11).  Trade along

377 the Silk Road connecting Europe and Asia was very active in the middle of the first millennium,

378 when we estimate L3 emerged (Hansen 2012; Ball 2016); its distribution suggests it spread

379 primarily along trading routes connecting Northeast Africa, Western Asia, and South Asia

380 (André and Filliozat 1986; Sartre 1991; Hansen 2012; Ball 2016) (fig. 6).  We speculate that this

381 occurred *via* overland routes, which may have limited the migration of L3 relative to maritime

382 dispersal of the other lineages.

383

384 The geographic distribution of L4 is strikingly broad (fig. 1) and it exhibits minimal population

385 structure (table 1).  This suggests L4 dispersed efficiently and continued to mix fluidly among

386 regions, a pattern we would expect if it was carried by an exceptionally mobile population of

387 hosts.  L4 is currently concentrated in regions bordering the Mediterranean, and elsewhere

388 throughout Africa and Europe (fig. 1).  We estimate the MRCA of L4 emerged in the 1$^{st}$ century

389 CE (range 368 BCE-362 CE), during the peak of Roman Imperial power across the entire

390 Mediterranean world and expansionist Roman policies into Africa, Europe, and Mesopotamia

391 (Luttwak 1976; Isaac 2004).  The empire reached its greatest territorial extent in the early

392 second century CE, when all of North Africa, from the Atlantic Ocean to the Red Sea, was under

393 a single power, with trade on land and sea facilitated by networks of stone-paved roads and

394 protected maritime routes (Luttwak 1976; Millar 1993; Ball 2016).  Primary sources from Roman

395 civilization attest to trade with China, purposeful expeditions for exploration, cartography, and

396 trade in the Red Sea and Indian Ocean (Pfister and Bellinger 1945; Dilke 1985; Begley and De

397 Puma 1991; Erdkamp 2002; Butcher 2003).

398

399 We hypothesize that the broad distribution of L4 reflects rapid diffusion from the Mediterranean

400 region along trade routes extending throughout Africa, the Middle East, and on to India, China,

401 and South Eastern Asia.  High rates of migration appear to have been maintained for this

402 lineage over much of its evolutionary history (fig. 3); patterns of connectivity implicate Europe

403 and Africa in its dispersal (fig. S11).  The association of L4 with European migrants is well

404 described, particularly migrants to the Americas (Gagneux et al. 2006; Pepperell et al. 2011).

405 Here we note bacterial population growth preceded geographic range expansion in L4 ~ca. 15$^{th}$

406 century (fig. 3), which coincides with the onset of the 'age of exploration' (Alam and

407 Subrahmanyam 2009) that would have provided numerous opportunities for spread of this

408 lineage from Europeans to other populations.  We also note the origin and concentration of this

409 lineage on the African continent.  Our sample of L4 isolates includes several deeply rooting

12

410    African isolates, and African isolates are interspersed throughout the phylogeny (fig. 4, fig. S6,

411    fig. S8).

412

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

414    Specifically, L5 and L6 are restricted to Western Africa and L7 is found only in Ethiopia (fig. 4,

415    fig. S6).  The reasons for the restricted distributions of these lineages are not immediately

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

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

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

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

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

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

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

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

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

425

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

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

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

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

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

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

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

433

434    Our analyses suggest that the expansion of L2 was preceded by an impressive increase in its

435    rate of migration (fig. 3), implying that growth of the pathogen population was facilitated by

436    expansion into new niches.  Our phylogeographic reconstructions implicate Russia, Central

437    Asia, and Western Asia in the recent migratory history of L2 (fig. S10, fig. S11), which is

438    consistent with a published phylogeographic analysis of L2 (Luo et al. 2015).  The inferred

439    timing of the growth and increased migration of L2 (~ca. 13th century) is close to the well

440    documented incursion of *Yersinia pestis* from Central Asia into Europe that resulted in explosive

441    plague epidemics (Benedictow 2004).  The experience with plague suggests that patterns of

442    connectivity among humans and other disease vectors were shifting at this place and time,

443    which would potentially open new niches for pathogens including *M.tb*.

444

445  We estimate that L1 underwent expansion ~ca. 17[th] century (fig. 3) but in this case it appears to

446  have grown *in situ*, e.g. due to changing environmental conditions such as increased crowding,

447  and/or growth of local human populations.  A study of the molecular epidemiology of TB in

448  Vietnam identified numerous recent migrations of L2 and L4 into the region, versus a stable

449  presence of L1 (Holt et al. 2018); this is consistent with our finding of higher recent rates of

450  migration for L2 and L4 versus L1 (fig. 3).  A pattern similar to L1 has been identified previously,

451  in the delay between dispersal of *M.tb* from European migrants to Canadian First Nations and

452  later epidemics of TB driven by shifting disease ecology (Pepperell et al. 2011).  These results

453  demonstrate the complex relationship between *M.tb* population growth and migration, and show

454  that under favorable conditions the pathogen can expand into novel niches or accommodate

455  growth in an existing niche.

456

457  In a previous study, we used analyses of synonymous and non-synonymous SFS to delineate

458  effects of purifying selection, linkage of sites, and population expansion on global populations of

459  *M.tb* (Pepperell et al. 2013).  Simulation studies have shown that purifying selection can affect

460  demographic inference with BEAST and SFS-based methods (Ewing and Jensen 2015;

461  Lapierre et al. 2016).  Although our analyses here using $\partial a\partial i$ were restricted to synonymous

462  SFS, it is likely that inference of population size changes with this method and with BEAST were

463  affected by purifying selection on this fully linked genome.  The magnitude of inferred

464  expansions may thus reflect both population size changes and background selection, and

465  should not be interpreted as direct reflections of historical changes in census population size.

466  We did not detect an effect of purifying selection on inference of migration in our three

467  population simulation analyses (fig. S16, fig. S17), but differences in the strength of purifying

468  selection could contribute to the lineage-specific differences we observed in the size of inferred

469  population expansions: i.e., genome-wide patterns of purifying selection could differ among

470  lineages.  We previously found evidence suggesting that the fitness trade-offs of drug resistance

471  mutations vary among lineages (Mortimer et al. 2018), making this intriguing possibility

472  potentially feasible.

473

474  While this study comprises the largest phylogeographic analysis on *M.tb* done to date, with 552

475  isolates collected from 51 countries and all seven described lineages represented, it has some

476  important limitations.  We did not attempt to estimate the rate or timescale of *M.tb* evolution,

477  instead relying on published rates that were calibrated with ancient DNA.  This is an active area

14

478    of research, and newly discovered ancient *M.tb* DNA samples will likely refine inference of both

479    the timing and locations of historical migration events, though it is critical to note that recent

480    substitution rate estimates of *M.tb* have converged on rates around $5 \times 10^{-8}$ substitutions per site

481    per year (Eldholm et al. 2016). Even when substitution rate estimates can be estimated with

482    confidence, the precision with which individual events can be dated using genetic data should

483    not be over-stated, as evidenced by broad 95% credible intervals for internal node date

484    estimates (e.g., Eldholm et al. 2016). Our goal here was to reconstruct historical migration of

485    *M.tb* throughout Eurasia and Africa and place this evolutionary history within a broad historical

486    context; the historical phenomena that we connect with the spread of TB involved vast areas

487    and extended over hundreds and in some cases thousands of years. Our reconstruction of the

488    global dispersal of TB within a temporal framework provided by ancient *M.tb* DNA analysis links

489    spread of the disease to the first ~1500y of the common era, a period of remarkable

490    intensification in the connectedness among peoples of Africa, Asia and Europe (Green 2018).

491

492    Methods

493    **Lineage Frequencies.** The SITVIT WEB database (Demay et al. 2012), which is an open

494    access *M.tb* molecular markers database, was accessed on September 5, 2016. Spoligotypes

495    were translated to lineages based on the following study (Shabbeer et al. 2012). The following

496    conversions were also included: EAI7-BGD2 for L1, CAS for L3, and LAM7-TUR, LAM12-

497    Madrid1, T5, T3-OSA, and H4 for L4. Isolates containing ambiguous spoligotypes (denoted with

498    >1 spoligotype) were inspected manually and assigned to appropriate lineages. Relative

499    lineage frequencies of lineages 1-6 for each country containing data for >10 isolates were

500    calculated and plotted with the rworldmap package in R (South 2016).

501

502    **Sample Description.**

503    *Old World collection*. We assembled/aligned publicly available whole genome sequences

504    (WGS) of thousands of *M.tb* isolates from recently published studies and databases for which

505    country of origin information were known and corresponded to traditional definitions of the Old

506    World. Isolates were assembled via reference guided assembly (RGA) when FASTQ data were

507    available and by multiple genome alignment (MGA) when only draft genome assemblies were

508    accessible (see below). As we were interested in reconstructing historical migrations of the

509    pathogen, we excluded countries where the majority of contemporary TB cases are identified in

510    recent immigrants (Barry et al. 2012; CCDIC 2014; ESR 2015; CDC 2016; PHE 2016; White et

511    al. 2017). Due to computational limitations (BEAST analyses), we necessarily took measures to

512     limit our dataset to <600 isolates. For countries with large numbers of available genomes, we

513     implemented a sub-sampling strategy similar one previously described (Thorpe et al. 2017),

514     whereby phylogenetic lineage diversity was captured thus minimizing the overrepresentation of

515     clonal complexes (e.g., outbreaks): phylogenetic inference on all isolates available from a

516     country was performed with Fasttree (Price et al. 2010) and a random isolate was selected from

517     each clade extending from $n$ branches, where $n$ was the desired number of isolates from the

518     country. Numbers of isolates per country were selected based on the availability of appropriate

519     genome sequence data as well as relative TB prevalence (fig. S1) (WHO 2017). All isolates

520     belonging to lineages 5-7 were retained. As a whole, this dataset reflects a 'mixed' sampling

521     scheme (Lapierre et al. 2016), where lineages L5-L7 are overrepresented relative to their

522     contemporary frequencies (fig. 1). At the lineage-specific scale, L1-L4 approximate random

523     sampling of available genomes. Our final Old World collection consisted of the WGS of 552

524     previously published *M.tb* isolates collected from 51 countries spanning 13 UN geoscheme

525     subregions. Accession numbers and pertinent information about each sample can be found in

526     table S1.

527

528     We note that our sample necessarily contains a large number of drug-resistant isolates as these

529     are more commonly sequenced. We also acknowledge that the studies we draw genomes from

530     may have been subject to other sampling biases for which we are unaware.

531

532     *Northern and Central American collection*. For one analysis, we included an additional 15

533     isolates from a previous study (Comas et al. 2015) for which country of origin information were

534     known and corresponded to the Americas. Isolates were assembled via RGA (see below) and

535     their genotypes at the 3,838,249 bp considered for all analyses of the Old World collection were

536     extracted.

537

538     **Reference Guided Assembly.** Previously published FASTQ data were retrieved from the

539     National Center for Biotechnology Information (NCBI) sequence read archive (SRA) (Leinonen

540     et al. 2011). Low-quality bases were trimmed using a threshold quality of 15, and reads

541     resulting in less than 20bp length were discarded using Trim Galore!

542     (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/), which is a wrapper tool

543     around Cutadapt (Martin 2011) and FastQC

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

545     (NC_000962.3) (Cole et al. 1998) with the MEM algorithm (Li 2013). Duplicates were removed

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

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

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

549    coverage were discarded, as were runs for which <70% of the reads mapped as determined by

550    Qualimap (García-Alcalde et al. 2012).  Pilon (Walker et al. 2014) was used to call variants with

551    the following parameters: --variant --mindepth 10 --minmq 40 --minqual 20.

552

553    **Multiple Genome Alignment.**  Draft genome assemblies were aligned to H37Rv

554    (NC_000962.3) (Cole et al. 1998) with Mugsy v1.2.3 (Angiuoli and Salzberg 2011).  Regions not

555    present in H37Rv were removed and merged with the reference-guided assembly.

556

557    **SNP alignment.**  Variant calls (VCFs) were converted to FASTAs with in-house scripts that

558    treat ambiguous calls and deletions as missing data (available at https://github.com/pepperell-

559    lab/RGAPepPipe).  Transposable elements, phage elements, and repetitive families of genes

560    (PE, PPE, and PE-PGRS gene families) that are poorly resolved with short read sequencing

561    were masked to missing data.  Isolates with >20% missing sites were excluded from the Old

562    World collection (table S1).  Variant positions with respect to H37Rv were extracted with SNP-

563    sites (Page et al. 2016) resulting in 60,818 variant sites.  Only sites where at least half of the

564    isolates had confident data (i.e., non-missing) were included in the phylogeographic models and

565    population genetic analyses (60,787 variant sites; 3,838,249 bp).  1.7% of variant sites landed in

566    loci associated with drug resistance (table S5).

567

568    **Geographic Information.**  Geographic locations for each of the 552 samples in the Old World

569    collection were obtained from NCBI and/or the respective publications from which the isolates

570    were first described.  When precise geographic information was available (e.g., city, province,

571    etc.), coordinates were obtained from www.mapcoordinates.net.  When only country level

572    geographic information was found, the 'Create Random Point' tool in ArcGIS 10.3 was used to

573    randomly place each isolate without specific latitude and longitude inside its respective country;

574    inhospitable areas (e.g., deserts and high mountains) and unpopulated areas from each country

575    using 50m data from Natural Earth (http://www.naturalearthdata.com/downloads, accessed

576    February 17, 2016) were excluded as possible coordinates.  The 'precision' column of table S1

577    reflects which method was used.

578

579     **Trade Route Information.** Data for all trade routes active throughout Europe, Africa, and Asia

580     by 1400 CE were compiled from the Old World Trade Routes (OWTRAD) Project

581     (www.ciolek.com/owtrad.html, accessed February 17, 2016). For each route, both node

582     information (trade cities, oases, and caravanserai) and arc information (the routes between

583     nodes) were imported into ArcGIS (fig. 6). *M.tb* isolate locations were also imported as points

584     and the 'Generate Near Table' tool was used to assign each isolate to its nearest node in the

585     trade network and is listed in the 'NearPost' column of table S1.

586

587     **Maximum Likelihood Inference.** We used RAxML v8.2.3 (Stamatakis 2014) for maximum

588     likelihood phylogenetic analysis of the Old World collection (all sites where at least half of

589     isolates had non-missing data) under the general time reversible model of nucleotide

590     substitution with a gamma distribution to account for site-specific rate heterogeneity. Rapid

591     bootstrapping of the corresponding SNP alignment was performed with the -autoMR flag,

592     converging after 50 replicates. Tree visualization was created with the ggtree package in R (Yu

593     et al. 2017).

594

595     **Structure Analyses.** Unsupervised clustering analysis of Old World isolates was performed

596     with STRUCTURE v2.3.4 (Pritchard et al. 2000). For $K$ values between 2-12, ten replicate runs

597     consisting of 10,000 burn-in iterations followed by 50,000 iterations were performed with default

598     settings on a subset of the Old World SNP alignment (4053 SNPs occurring at a minor allele

599     frequency > 0.01 with no missing data). StructureHarvester (Earl and vonHoldt 2012) was used

600     to collate results and determine the most suitable value of $K$ following the "Evanno" method

601     (Evanno et al. 2005). Replicate runs with the lowest log likelihood for each value of $K$ were

602     used for visualization of results.

603

604     **Phylogeographic & Demographic Inference with BEAST.** The Old World collection SNP

605     alignment and individual lineage SNP alignments were analyzed using the Bayesian Markov

606     Chain Monte Carlo coalescent method implemented in BEAST v1.8 (Drummond and Rambaut

607     2007) with the BEAGLE library (Ayres et al. 2012) to facilitate rapid likelihood calculations.

608     Analyses were performed using the general time reversible model of nucleotide substitution with

609     a gamma distribution to account for rate heterogeneity between sites, a strict molecular clock,

610     and both constant and Bayesian skyline plot (BSP) demographic models. Country of origin or

611     the UN subregion for each isolate was modeled as a discrete phylogenetic trait (Lemey et al.

612     2009). All Markov chains were run for at least 100 million generations, sampled every 10,000

613     generations, and with the first 10,000,000 generations discarded as burn-in; replicate runs were

614     performed for analyses and combined to assess convergence.  Estimated sample size (ESS)

615     values of non-nuisance parameters were >200 for all analyses.  Site and substitution model

616     choice were based on previous analyses of *M.tb* global alignments as opposed to an exhaustive

617     comparison of models which would require unreasonable computational resources.  Strict vs

618     relaxed molecular clocks did not result in altered trends of migration at the lineage level, and

619     comparisons between analyses using strict and relaxed clocks show strong correlation between

620     the estimated height of nodes (e.g., $R^2$ > 0.97; fig S18).  Table S6 provides a summary of

621     BEAST analyses presented and the results derived from them.  Tree visualizations were

622     created with FigTree (http://tree.bio.edu.ac.uk/software/figtree/) and the ggtree package in R

623     (Yu et al. 2017).

624

625     We note that phylogeographic inference methods are an active area of research and

626     increasingly sophisticated models are continuously being developed [e.g. (Lemey et al. 2010;

627     De Maio et al. 2015)].  We found alternative methods unsuitable and/or intractable for our large

628     dataset.  As methods improve, comparison of the results inferred herein to other

629     phylogeographic models will be important to investigate the sensitivity of our results to the

630     method of phylogeographic inference.

631

632     **Demographic inference from the observed site frequency spectrum (SFS).**  SNP-sites

633     (Page et al. 2016) was used to convert the Old World collection alignment to a multi-sample

634     VCF and SnpEff (Cingolani et al. 2012) was used to annotate variants with respect to H37Rv

635     (NC_000962.3) (Cole et al. 1998) as synonymous, non-synonymous, or intergenic.  Loci at

636     which any sequence in the population had a gap or unknown character were removed from the

637     data set.  Demographic inference with the synonymous SFS for each of the seven lineages and

638     the entire collection was performed using ∂a∂i (Gutenkunst et al. 2009).  We modeled constant

639     population size (standard neutral model), an instantaneous expansion model, and an

640     exponential growth model, and identified the best-fit model and maximal likelihood parameters

641     of the demographic model given our observed data.  Our parameter estimates, ν and τ, were

642     optimized for the instantaneous expansion and exponential growth models.  Uncertainty

643     analysis of these parameters were analyzed using the Godambe Information Matrix (Coffman et

644     al. 2016) on 100 samplings of the observed synonymous SFS with replacement and subsequent

645     model inference.

646

647 **Population genetic statistics.** Nucleotide diversity (π) and Watterson's theta (Ө) for various

648 population assignments (e.g., lineage, UN subregion) were calculated with EggLib v2.1.10 (De

649 Mita and Siol 2012).

650

651 **Analysis of Molecular Variance (AMOVA).** AMOVAs were performed using the 'poppr.amova'

652 function (a wrapper for the ade4 package (Dray et al. 2007) implementation) in the poppr

653 package in R (Kamvar et al. 2014). Bins were assigned via the following classification systems:

654 UN geoscheme subregions and Level 1 ('botanical continents') of the World geographical

655 scheme for recording plant distributions. Isolate assignation can be found in table S1. Genetic

656 distances between isolates were calculated with the 'dist.dna' function of the ape v4.0 package

657 in R (Paradis et al. 2004) from the SNP alignment of the Old World collection.

658

659 **Mantel tests.** Great circle distances between *M.tb* isolate locations were calculated with the

660 'distVincentyEllipsoid' function in the geosphere R package (Hijmans et al. 2016). Geographic

661 distances between isolate locations along the trade network were calculated by adding the great

662 circle distances from the isolates to the nearest trade hubs and the shortest distance between

663 trade hubs along the trade network; the latter was determined using an Origin-Destination Cost

664 Matrix and the 'Solve' tool in the Network Analyst Toolbox of ArcGIS which calculates the

665 shortest distance from each origin to every destination along the arcs in the trade network. In

666 the event that two isolates were assigned to the same trade post, the great circle distance

667 between the isolates was used. To calculate the geographic distance between isolates in a

668 manner that reflects human migrations, the great circle distance between isolates and

669 waypoints were summed. These were calculated with a custom R function (available at

670 https://github.com/ONeillMB1/Mtb_Phylogeography_v2) using a series of rules to define

671 whether or not the path between isolates would have gone through a waypoint. For all three

672 distance metrics, values were log transformed and standardized. Genetic distances between

673 isolates were calculated with the 'dist.dna' function in the ape v4.0 package in R (Paradis et al.

674 2004) from the SNP alignment. The 'mantel' function of the vegan package in R (Oksanen et al.

675 2017) was used to perform a Mantel test between the genetic distance matrix and each of the

676 three geographic matrices for both the Old World collection and each individual lineage. Four of

677 the 552 isolates were excluded from these analyses as they were from Kiribati and trade

678 networks spanning this region were not compiled.

679

680 **Relationship between genetic diversity and geographic distance from Addis Ababa.** For

681 this analysis, we added Northern and Central American datasets, assembled in an identical

682 manner to those of the Old World collection and masked at sites where less than half of the Old

683 World collection had confident data (3,838,249 bp). For each UN subregion, the mean latitude

684 and longitude coordinates for all *M.tb* isolates within the region were calculated. The great

685 circle distances from these average estimates for regions to Addis Ababa were then calculated,

686 using waypoints for between-continent distance estimates to make them more reflective of

687 presumed human migration patterns (Ramachandran et al. 2005). Cairo was used as a

688 waypoint for Eastern Europe, Central Asia, Western Asia, Southern Asia, Eastern Asia, and

689 South Eastern Asia; Cairo and Istanbul were used as waypoints for Western Europe and

690 Southern Europe; Cairo, Anadyr, and Prince Rupert were used as waypoints for Northern and

691 Central America. The distance between each region and Addis Ababa were the sum of the

692 great circle distances between the two points (the average coordinates for the UN subregion

693 and Addis Ababa) and the waypoint(s) in the path connecting them, plus the great circle

694 distance(s) between waypoints if two were used. Treating each UN subregion as a population,

695 the relationship between genetic diversity (assessed with π) and geographic distance from

696 Addis Ababa were explored with linear regression for both the entire Old World collection and

697 individual lineages in R (R Development Core Team). Code is available at

698 https://github.com/ONeillMB1/Mtb_Phylogeography_v2.

699

700 **Migration Rate Inference.** Migration rates through time were inferred from the Bayesian

701 maximum clade credibility trees for the entire Old World collection of *M.tb* isolates ($n$ = 552).

702 Individual lineages that contain isolates from multiple UN subregions (i.e., L1: $n$ = 89, L2: $n$ =

703 181, L3: $n = 65$, and L4: $n$ = 143) were extracted and plotted separately. Only nodes with

704 posterior probabilities greater than or equal to 80% were considered. A migration event was

705 classified as a change in the most probable reconstructed ancestral geographic region from a

706 parent to child node. Median heights of the parent and child nodes were treated as a range of

707 time that the migration event could have occurred. The rate of migration through time for each

708 lineage or the Old World collection was inferred by summing the number of migration events

709 occurring across every year of the time-scaled phylogeny, divided by the total number of

710 branches in existence during each year of the time-scaled phylogeny (both those displaying a

711 migration event and those that do not). Code for these analyses is available at

712 https://github.com/ONeillMB1/Mtb_Phylogeography_v2.

713

21

714    Additionally, relative migration rates between UN subregions were derived from the BEAST

715    analyses of phylogeography.  The Bayesian stochastic search variable selection method

716    (BSSVS) for identifying the most parsimonious migration matrix implemented in BEAST as part

717    of the discrete phylogeographic migration model (Lemey et al. 2009) allowed us to use Bayes

718    factors (BF) to identify the migration rates with the greatest posterior support and provide

719    posterior estimates for their relative rates.  Strongly supported relative rates (BF > 5) and

720    connectivity among subregions were visualized with Cytoscape v3.2.0 (Shannon et al. 2003)

721    and superimposed onto a map generated with the 'rworldmap' package in R (South 2016). To

722    assess the effect of sampling on migration rate inference, UN regions harboring greater than 10

723    and 20 isolates were randomly subsampled to even numbers and subject to the same analysis.

724    This was done 10 times each for $n$ = 10 and $n$ = 20.

725

726    **Effect of selection on estimates of migration.** We performed demographic forward-in-time

727    simulations using the SFS_CODE package (Hernandez 2008), which allows for demographic

728    models with arbitrarily complex migration and selection regimes.  Our simulations were

729    performed under a simple two population model or with a more complex three population model.

730    In all simulations, $N_e$ for each population was 1000, θ was 0.001 (O'Neill et al. 2015), and

731    migration between each pair of populations was symmetrical.  As there is substantial evidence

732    for little to no recombination in the *M.tb* genome, our simulations were performed without

733    recombination.

734

735    The two population simulations were performed under three scenarios:  1) no migration between

736    populations after initial divergence; 2) constant migration after divergence (per generation $M$ =

737    0.5) without selection; and 3) constant migration ($M$ = 0.5) with purifying selection (25% of

738    alleles of each population have a population selection coefficient of -1.0, and the rest are

739    neutral) after divergence.

740

741    The three population simulations were performed under five scenarios: 1) no migration between

742    populations after simultaneous divergence of the three populations; 2) constant, symmetrical

743    migration after divergence (per generation $M$ = 0.5 for all population pairs) without selection; 3)

744    constant, symmetrical migration ($M$ = 0.5) with purifying selection (25% of alleles in all

745    populations have a population selection coefficient of -1.0, and the rest are neutral); 4) constant,

746    asymmetrical migration after divergence ($M$ = 0.5 for migration between pop0 and pop1, $M$ = 5.0

747    for migration between pop1 and pop2, and $M$ = 0 for migration between pop0 and pop2) without

748    selection; and 5) constant, asymmetrical migration after divergence ($M = 0.5$ between pop0 and

749    pop1, $M = 5.0$ between pop1 and pop2, and $M = 0$ between pop0 and pop2) with purifying

750    selection (25% of alleles in all populations have a population selection coefficient of -1.0, and

751    the rest are neutral).

752

753    For all simulations, 25 samples were taken from each population, and sequences of 100000

754    bases were generated.  Twenty simulations were performed under each scenario for both the 2

755    population (60 simulations) and 3 population (100 simulations) models.  Each sequence

756    alignment was subsequently subjected to migration analysis in ∂a∂i (Gutenkunst et al. 2009, see

757    Note S2) and BEAST v1.8.4 (Drummond and Rambaut 2007).  For each Bayesian coalescent

758    analysis, the HKY+G substitution model, a constant population model, and a strict molecular

759    clock model were used.  A discrete symmetrical migration model (Lemey et al. 2009) was used

760    to determine migration rates, and BSSVS (Lemey et al. 2009) was used to estimate BF support

761    for migration rates in the 3 population simulations.  All Markov chains were run for 10 million

762    generations or until convergence, with samples taken every 10,000 steps, and 10% discarded

763    as burn-in.  The package SpreaD3 v0.96 (Bielejec et al. 2016) was used to calculate BF support

764    for migration rates.

765

766    Acknowledgments

## **References**

773 Alam M, Subrahmanyam S. 2009. Indo-Persian travels in the age of discoveries, 1400-1800.
775      Digit. pr. Cambridge: Cambridge University Press.

776 André J, Filliozat J. 1986. L'Inde vue de Rome: textes latins de l'antiquité, relatifs à l'Inde.
777      Paris: Belles Lettres.

778 Angiuoli SV, Salzberg SL. 2011. Mugsy: fast multiple alignment of closely related whole
779      genomes. Bioinformatics 27:334–342.

780 Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO, Huelsenbeck JP, Ronquist F,
781      Swofford DL, Cummings MP, et al. 2012. BEAGLE: An Application Programming
782      Interface and High-Performance Computing Library for Statistical Phylogenetics. Syst.
783      Biol. 61:170–173.

784 Ball W. 2016. Rome in the East: The Transformation of an Empire. 2nd ed. London & New
785      York: Routledge.

786 Barry C, Waring J, Stapledon R, Konstantinos A. 2012. Tuberculosis notifications in Australia,
787      2008 and 2009. Communicable Diseases Intelligence 36:86-94.

788 Begley V, De Puma RD. 1991. Rome and India: The Ancient Sea Trade. Madison: University of
789      Wisconsin Press.

790 Behr MA, Wilson MA, Gill WP, Salamon H, Schoolnik GK, Rane S, Small PM. 1999.
791      Comparative genomics of BCG vaccines by whole-genome DNA microarray. Science
792      284:1520–1523.

793 Benedictow OJ. 2004. The Black Death, 1346-1353: the complete history. Woodbridge, Suffolk,
794      UK ; Rochester, N.Y., USA: Boydell Press.

795 Bielejec F, Baele G, Rodrigo AG, Suchard MA, Lemey P. 2016. Identifying predictors of time-
796      inhomogeneous viral evolutionary processes. Virus Evol. 2:vew023.

797 Blench R. 1996. The Ethnographic Evidence for Long-distance Contacts between Oceania and
798      East Africa. In: Reade J, editor. The Indian Ocean in antiquity.

799 Blouin Y, Hauck Y, Soler C, Fabre M, Vong R, Dehan C, Cazajous G, Massoure P-L, Kraemer
800      P, Jenkins A, et al. 2012. Significance of the Identification in the Horn of Africa of an
801      Exceptionally Deep Branching Mycobacterium tuberculosis Clade. PLOS ONE
802      7:e52841.

803 Bos KI, Harkins KM, Herbig A, Coscolla M, Weber N, Comas I, Forrest SA, Bryant JM, Harris
804      SR, Schuenemann VJ, et al. 2014. Pre-Columbian mycobacterial genomes reveal seals as
805      a source of New World human tuberculosis. Nature 514:494–497.

806 Boussac M-F, Salles J-F, France eds. 1995. Athens, Aden, Arikamedu: essays on the
807     interrelations between India, Arabia, and the eastern Mediterranean. New Delhi:
808     Manohar : Distributed in South Asia by Foundation Books.

809 Brosch R, Gordon SV, Marmiesse M, Brodin P, Buchrieser C, Eiglmeier K, Garnier T, Gutierrez
810     C, Hewinson G, Kremer K, et al. 2002. A new evolutionary scenario for the
811     Mycobacterium tuberculosis complex. Proc. Natl. Acad. Sci. 99:3684–3689.

812 Butcher K. 2003. Roman Syria and the Near East. Los Angeles: J. Paul Getty Museum : Getty
813     Publications.

814 Centers for Disease Control and Prevention (CDC). 2016. Reported Tuberculosis in the United
815     States, 2015. Atlanta, GA: US Department of Health and Human Services.

816 Centre for Communicable Diseases and Infection Control (CCDIC). 2014. Tuberculosis
817     prevention and control in Canada a federal framework for action. Public Health Agency
818     of Canada.

819 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012.
820     A program for annotating and predicting the effects of single nucleotide polymorphisms,
821     SnpEff. Fly (Austin) 6:80–92.

822 Coffman AJ, Hsieh PH, Gravel S, Gutenkunst RN. 2016. Computationally Efficient Composite
823     Likelihood Statistics for Demographic Inference. Mol. Biol. Evol. 33:591–593.

824 Cole ST, Brosch R, Parkhill J, Garnier T, Churcher C, Harris D, Gordon SV, Eiglmeier K, Gas S,
825     III CEBI, et al. 1998. Erratum: Deciphering the biology of Mycobacterium tuberculosis
826     from the complete genome sequence. Nat. Lond. 396:190.

827 Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, Parkhill J, Malla B, Berg S,
828     Thwaites G, et al. 2013. Out-of-Africa migration and Neolithic coexpansion of
829     Mycobacterium tuberculosis with modern humans. Nat. Genet. 45:1176–1182.

830 Comas I, Hailu E, Kiros T, Bekele S, Mekonnen W, Gumi B, Tschopp R, Ameni G, Hewinson
831     RG, Robertson BD, et al. 2015. Population Genomics of Mycobacterium tuberculosis in
832     Ethiopia Contradicts the Virgin Soil Hypothesis for Human Tuberculosis in Sub-Saharan
833     Africa. Curr. Biol. 25:3260–3266.

834 Coningham R, Young R. 2015. The Archaeology of South Asia: From the Indus to Asoka,
835     c.6500 BCE–200 CE. Cambridge University Press.

836 De Maio N, Wu C-H, O'Reilly KM, Wilson D. 2015. New Routes to Phylogeography: A
837     Bayesian Structured Coalescent Approximation. PLoS Genet 11:e1005421.

838 De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population
839     genetics and genomics. BMC Genet. 13:27.

840   Demay C, Liens B, Burguière T, Hill V, Couvin D, Millet J, Mokrousov I, Sola C, Zozio T,
841           Rastogi N. 2012. SITVITWEB – A publicly available international multimarker database
842           for studying Mycobacterium tuberculosis genetic diversity and molecular epidemiology.
843           Infect. Genet. Evol. 12:755–766.

844   DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del
845           Angel G, Rivas MA, Hanna M, et al. 2011. A framework for variation discovery and
846           genotyping using next-generation DNA sequencing data. Nat. Genet. 43:491–498.

847   Dilke O. 1985. Greek and Roman Maps. Baltimore: Johns Hopkins University Press

848   Dray S, Dufour A-B, others. 2007. The ade4 package: implementing the duality diagram for
849           ecologists. J. Stat. Softw. 22:1–20.

850   Drummond AJ, Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees.
851           BMC Evol. Biol. 7:214.

852   Duchêne S, Holt KE, Weill F-X, Le Hello S, Hawkey J, Edwards DJ, Fourment M, Holmes EC.
853           2016. Genome-scale rates of evolutionary change in bacteria. Microb. Genomics 2:
854           e000094.

855   Earl DA, vonHoldt BM. 2012. STRUCTURE HARVESTER: a website and program for
856           visualizing STRUCTURE output and implementing the Evanno method. Conserv. Genet.
857           Resour. 4:359–361.

858   Eldholm V, Balloux F. 2016. Antimicrobial Resistance in Mycobacterium tuberculosis: The Odd
859           One Out. Trends Microbiol. 24:637–648.

860   Eldholm V, Pettersson JH-O, Brynildsrud OB, Kitchen A, Rasmussen EM, Lillebaek T, Rønning
861           JO, Crudu V, Mengshoel AT, Debech N, et al. 2016. Armed conflict and population
862           displacement as drivers of the evolution and dispersal of Mycobacterium tuberculosis.
863           Proc. Natl. Acad. Sci. 113:13881–13886.

864   Erdkamp P. 2002. The Roman Army and the Economy. Amsterdam: Gieben.

865   Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of individuals using the
866           software structure: a simulation study. Mol. Ecol. 14:2611–2620.

867   Ewing GB, Jensen JD. 2015. The consequences of not accounting for background selection in
868           demographic inference. Mol. Ecol. 25:135–141.

869   Firdessa R, Berg S, Hailu E, Schelling E, Gumi B, Erenso G, Gadisa E, Kiros T, Habtamu M,
870           Hussein J, et al. 2013. Mycobacterial lineages causing pulmonary and extrapulmonary
871           tuberculosis, Ethiopia. Emerg. Infect. Dis. 19:460–463.

872   Gagneux S, DeRiemer K, Van T, Kato-Maeda M, de Jong BC, Narayanan S, Nicol M, Niemann
873           S, Kremer K, Gutierrez MC, et al. 2006. Variable host–pathogen compatibility in
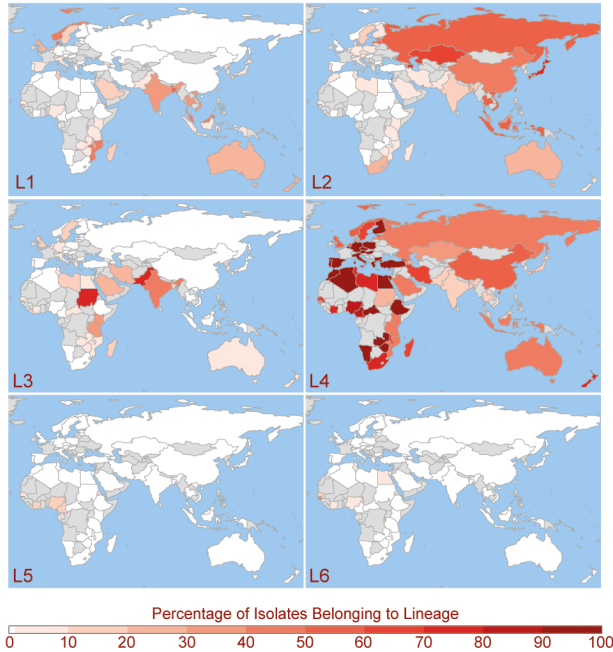874           Mycobacterium tuberculosis. Proc. Natl. Acad. Sci. U. S. A. 103:2869–2873.

875  García-Alcalde F, Okonechnikov K, Carbonell J, Cruz LM, Götz S, Tarazona S, Dopazo J,
876      Meyer TF, Conesa A. 2012. Qualimap: evaluating next-generation sequencing alignment
877      data. Bioinformatics 28:2678–2679.

878  Green MH. 2018. Climate and Disease in Medieval Eurasia. Oxf. Res. Encycl. Asian Hist.

879  Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the Joint
880      Demographic History of Multiple Populations from Multidimensional SNP Frequency
881      Data. PLoS Genet 5:e1000695.

882  Hansen V. 2012. The Silk Road: A New History. Oxford: Oxford University Press.

883  Hernandez RD. 2008. A flexible forward simulator for populations subject to selection and
884      demography. Bioinformatics 24:2786–2787.

885  Hershberg R, Lipatov M, Small PM, Sheffer H, Niemann S, Homolka S, Roach JC, Kremer K,
886      Petrov DA, Feldman MW, et al. 2008. High functional diversity in Mycobacterium
887      tuberculosis driven by genetic drift and human demography. PLoS Biol. 6:e311.

888  Hijmans RJ, Williams E, Vennes C. 2016. geosphere: Spherical Trigonometry. Available from:
889      https://cran.r-project.org/web/packages/geosphere/index.html

890  Hirsh AE, Tsolaki AG, DeRiemer K, Feldman MW, Small PM. 2004. Stable association between
891      strains of Mycobacterium tuberculosis and their human host populations. Proc Natl Acad
892      Sci U A 101:4871–4876.

893  Holt KE, McAdam P, Thai PVK, Thuong NTT, Ha DTM, Lan NN, Lan NH, Nhu NTQ, Hai HT,
894      Ha VTN, et al. 2018. Frequent transmission of the Mycobacterium tuberculosis Beijing
895      lineage and positive selection for the EsxW Beijing variant in Vietnam. Nat. Genet.
896      50:849–856.

897  Institute of Environmental Science and Research Ltd (ESR). 2015. Tuberculosis in New Zealand:
898      Annual Report 2014. Porirua: ESR.

899  Isaac BH. 2004. The limits of empire: the Roman army in the East. Oxford: Clarendon Press.

900  Jong D, C B, Hill PC, Aiken A, Awine T, Martin A, Adetifa IM, Jackson-Sillah DJ, Fox A,
901      Kathryn D, et al. 2008. Progression to Active Tuberculosis, but Not Transmission, Varies
902      by Mycobacterium tuberculosis Lineage in The Gambia. J. Infect. Dis. 198:1037–1043.

903  Kamvar ZN, Tabima JF, Grünwald NJ. 2014. Poppr: an R package for genetic analysis of
904      populations with clonal, partially clonal, and/or sexual reproduction. PeerJ 2:e281.

905  Kay GL, Sergeant MJ, Zhou Z, Chan JZ-M, Millard A, Quick J, Szikossy I, Pap I, Spigelman M,
906      Loman NJ, et al. 2015. Eighteenth-century genomes show that mixed infections were
907      common at time of peak tuberculosis in Europe. Nat. Commun. 6:6717.

908    Kent RK. 1979. The Possibilities of Indonesian Colonies in Africa with Reference to
909        Madagascar. In: Mouvements de Populations dans L'Ocean Indie. Paris: H. Champion. p.
910        93–105.

911    Lapierre M, Blin C, Lambert A, Achaz G, Rocha EPC. 2016. The Impact of Selection, Gene
912        Conversion, and Biased Sampling on the Assessment of Microbial Demography. Mol.
913        Biol. Evol. 33:1711–1725.

914    Leinonen R, Sugawara H, Shumway M. 2011. The Sequence Read Archive. Nucleic Acids Res.
915        39:D19–D21.

916    Lemey P, Rambaut A, Drummond AJ, Suchard MA. 2009. Bayesian Phylogeography Finds Its
917        Roots. PLOS Comput. Biol. 5:e1000520.

918    Lemey P, Rambaut A, Welch JJ, Suchard MA. 2010. Phylogeography Takes a Relaxed Random
919        Walk in Continuous Space and Time. Mol. Biol. Evol. 27:1877–1885.

920    Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
921        ArXiv13033997 Q-Bio. Available from: http://arxiv.org/abs/1303.3997

922    Luo T, Comas I, Luo D, Lu B, Wu J, Wei L, Yang C, Liu Q, Gan M, Sun G, et al. 2015.
923        Southern East Asian origin and coexpansion of Mycobacterium tuberculosis Beijing
924        family with Han Chinese. Proc. Natl. Acad. Sci. 112:8136–8141.

925    Luttwak EN. 1976. The grand strategy of the Roman Empire: from the first century A.D. to the
926        third. London: Weidenweld & Nicholson.

927    Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
928        EMBnet.journal 17:10–12.

929    Millar F. 1993. The Roman Near East, 31 B.C.-A.D. 337. Cambridge, Mass: Harvard University
930        Press.

931    Mortimer TD, Weber AM, Pepperell CS. 2018. Signatures of Selection at Drug Resistance Loci
932        in Mycobacterium tuberculosis. mSystems 3:e00108-17.

933    Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, Minchin PR, O'Hara
934        RB, Simpson GL, Solymos P, et al. 2017. vegan: Community Ecology Package.
935        Available from: https://cran.r-project.org/web/packages/vegan/index.html

936    O'Neill MB, Mortimer TD, Pepperell CS. 2015. Diversity of Mycobacterium tuberculosis across
937        Evolutionary Scales. PLOS Pathog. 11:e1005257.

938    Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-sites:
939        rapid efficient extraction of SNPs from multi-FASTA alignments. Microbial Genetics  2 .

940    Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R
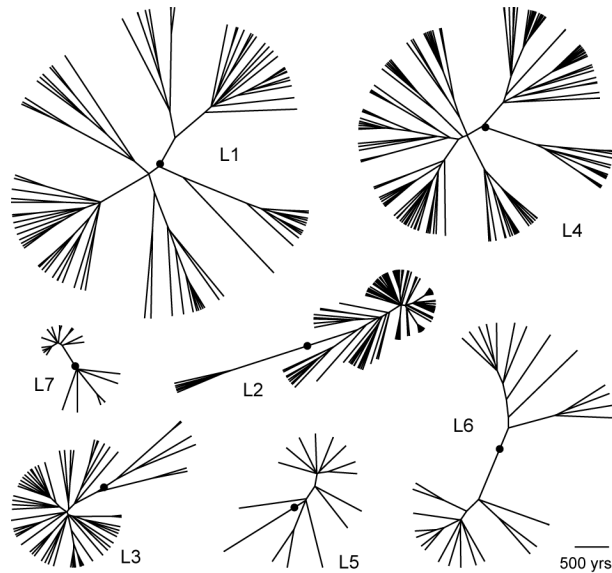941        language. Bioinformatics 20:289–290.

942    Parkin D, Barnes R eds. 2002. Ships and the development of maritime technology in the Indian
943         Ocean. London: RoutledgeCurzon.

944    Pepperell CS, Casto AM, Kitchen A, Granka JM, Cornejo OE, Holmes EC, Birren B, Galagan J,
945         Feldman MW. 2013. The Role of Selection in Shaping Diversity of Natural M.
946         tuberculosis Populations. PLoS Pathog 9:e1003543.

947    Pepperell CS, Granka JM, Alexander DC, Behr MA, Chui L, Gordon J, Guthrie JL, Jamieson
948         FB, Langlois-Klassen D, Long R, et al. 2011. Dispersal of Mycobacterium tuberculosis
949         via the Canadian fur trade. Proc. Natl. Acad. Sci. 108:6526–6531.

950    Pfister R, Bellinger L. 1945. The Excavations at Dura-Europos.  Final Report IV: The Textiles.
951         New Haven: Yale University Press.

952    Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees
953         for Large Alignments. PLOS ONE 5:e9490.

954    Pritchard JK, Stephens M, Donnelly P. 2000. Inference of Population Structure Using Multilocus
955         Genotype Data. Genetics 155:945–959.

956    Public Health England (PHE). 2016. Tuberculosis in England: 2016. London, UK: PHE.

957    R Development Core Team. R: A Language and Environment for Statistical Computing. Vienna,
958         Austria: R Foundation for Statistical Computing Available from: http: //www.R-
959         project.org/

960    Ramachandran S, Deshpande O, Roseman CC, Rosenberg NA, Feldman MW, Cavalli-Sforza
961         LL. 2005. Support from the relationship of genetic and geographic distance in human
962         populations for a serial founder effect originating in Africa. Proc. Natl. Acad. Sci. U. S.
963         A. 102:15942–15947.

964    Rauh NK. 2003. Merchants, sailors and pirates in the Roman world. Stroud: Tempus.

965    Ray HP. 2003. The archaeology of seafaring in ancient South Asia. Cambridge ; New York:
966         Cambridge University Press.

967    Ray HP, Salles J-F, Institute of Southeast Asian Studies, Maison de l'Orient méditerranéen
968         ancien (Lyon F, National Institute of Science, Technology and Development Studies,
969         France, Ambassade (India), Centre for Human Sciences. 1996. Tradition and
970         archaeology: early maritime contacts in the Indian Ocean. New Delhi: Manohar
971         Publishers.

972    Rocha EPC, Smith JM, Hurst LD, Holden MTG, Cooper JE, Smith NH, Feil EJ. 2006.
973         Comparisons of dN/dS are time dependent for closely related bacterial genomes. J.
974         Theor. Biol. 239:226–235.

975    Salles J-F. 1996. Achaemenid and Hellenistic Trade in the Indian Ocean. In: The Indian Ocean in
976         antiquity. p. 251–267.

977    Sartre M. 1991. L'Orient romain: provinces et sociétés provinciales en Méditerranée orientale
978        d'Auguste aux Sévères (31 avant J.-C-235 après J.-C.). Paris: Seuil.

979    Shabbeer A, Cowan LS, Ozcaglar C, Rastogi N, Vandenberg SL, Yener B, Bennett KP. 2012.
980        TB-Lineage: An online tool for classification and analysis of strains of Mycobacterium
981        tuberculosis complex. Infect. Genet. Evol. 12:789–797.

982    Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,
983        Ideker T. 2003. Cytoscape: A Software Environment for Integrated Models of
984        Biomolecular Interaction Networks. Genome Res. 13:2498–2504.

985    South A. 2016. rworldmap: Mapping Global Data. Available from: https://cran.r-
986        project.org/web/packages/rworldmap/index.html

987    Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
988        large phylogenies. Bioinformatics 30:1312–1313.

989    Thorpe HA, Bayliss SC, Hurst LD, Feil EJ. 2017. Comparative Analyses of Selection Operating
990        on Nontranslated Intergenic Regions of Diverse Bacterial Species. Genetics 206:363–
991        376.

992    Vogt B. 1996. Bronze Age Maritime Trade in the Indian Ocean: Harappan Traits on the Oman
993        Peninsula. In: Reade J, editor. The Indian Ocean in antiquity. p. 107–132.

994    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
995        Wortman J, Young SK, et al. 2014. Pilon: An Integrated Tool for Comprehensive
996        Microbial Variant Detection and Genome Assembly Improvement. PLOS ONE
997        9:e112963.

998    White Z, Painter J, Douglas P, Abubakar I, Njoo H, Archibald C, Halverson J, Robson J, Posey
999        DL. 2017. Immigrant Arrival and Tuberculosis among Large Immigrant- and Refugee-
1000        Receiving Countries, 2005-2009. Tuberc. Res. Treat 2017.

1001    Wink A. 2002. From the Mediterranean to the Indian Ocean: Medieval History in Geographic
1002        Perspective. Comp. Stud. Soc. Hist. 44:416–445.

1003    World Health Organization (WHO). 2017. Global tuberculosis report 2017. Geneva: WHO.

1004    Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. 2017. ggtree: an r package for visualization and
1005        annotation of phylogenetic trees with their covariates and other associated data. Methods
1006        Ecol. Evol. 8:28–36.

1007    Zarins J. 1996. Obsidian in the Larger Context of Predynastic/Archaic Egyptian Red Sea Trade.
1008        In: Reade J, editor. The Indian Ocean in antiquity. p. 107–132.

1009    Figures and Tables
1010



1011
1012    **Fig. 1.** Geographic distributions of *Mycobacterium tuberculosis* lineages 1-6. Spoligotypes from

1013    the SITVIT WEB database (*n* = 42,358) were assigned to lineages 1-6. Countries are colored

1014    from white to dark-red based on the percentage of isolates from the country belonging to each

1015    lineage. Unsampled countries and those with less than 10 isolates in the database are shown in

1016    gray. Lineage 7 (not pictured) is found exclusively in Ethiopia.

1017

**Fig. 2.** Maximum clade credibility phylogenies of *Mycobacterium tuberculosis* lineages 1-6. Bayesian analyses were performed on each lineage alignment with the general time reversible model of nucleotide substitution with a gamma distribution to account for rate heterogeneity between sites, a strict molecular clock, and Bayesian skyline plot demographic models. The most recent common ancestor 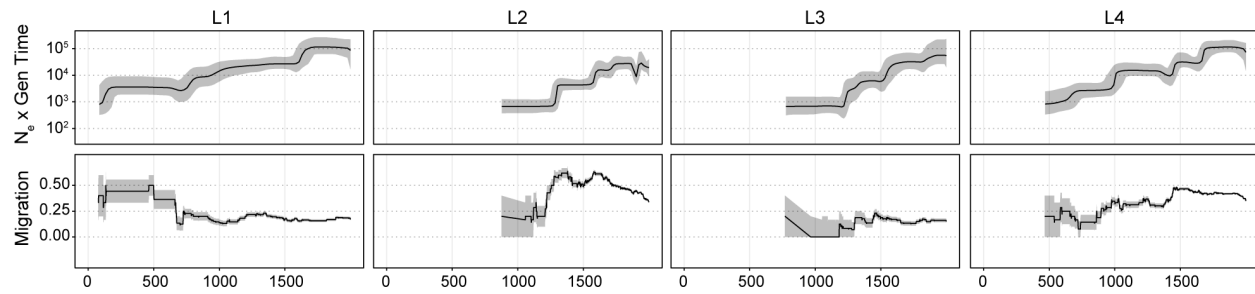(MRCA) of each lineage is indicated with a black circle; the MRCA of individual lineage phylogenies were informed by the phylogeny of the entire Old World collection, which was dated using a substitution rate of $5 \times 10^{-8}$ substitutions/site/year (Kay et al. 2015).
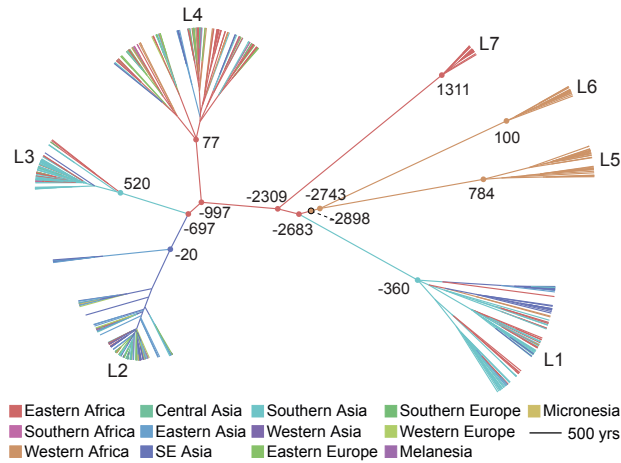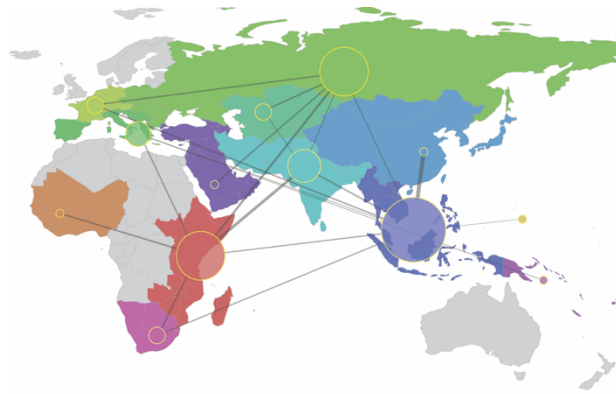
**Fig. 3.** Patterns of effective population size and migration through time of *Mycobacterium tuberculosis* lineages 1-4. Bayesian skyline plots (top panels) show inferred changes in effective population size ($N_e$) through time deduced from lineage specific analyses. Black lines denote median $N_e$ and gray shading the 95% highest posterior density. Estimated migration through time (see *Methods*) for each lineage is shown in the bottom panels. Gray shading depicts the rates inferred after the addition or subtraction of a single migration event, and demonstrate the uncertainty of rate estimates, particularly from the early history of each lineage. Dates are shown in calendar years and are based on scaling the phylogeny of the Old World collection with a substitution rate of 5 x $10^{-8}$ substitutions/site/year (Kay et al. 2015).

1039

**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.
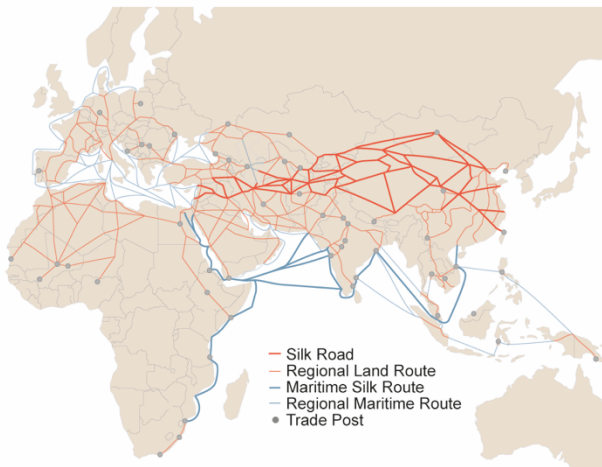
1040
1041
1042
1043
1044
1045
1046

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

1055

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

1060 **Table 1.** Genetic diversity of Old World *M.tb* across lineages 1-7. TMRCA estimates reflect
1061 scaling of results to evolutionary rates calibrated from ancient DNA [median $5.00 \times 10^{-8}$
1062 substitutions/ site/ year (Kay et al. 2015) and are written as calendar years. To account for
1063 uncertainty in this rate estimate, our lower and upper TMRCA estimates reflect scaling of our
1064 results with the low and high bounds of the 95% highest posterior density estimates of the rate
1065 reported from ancient DNA analysis (i.e. $4.06 \times 10^{-8}$ and $5.87 \times 10^{-8}$, respectively).

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

1066