

1 Multiple Introductions of the *Mycobacterium tuberculosis* Lineage 2– 2 Beijing into Africa over centuries

3 Liliana K. Rutaihwa^{1, 2, 3}, Fabrizio Menardo^{1, 2}, David Stucki^{1, 2}, Sebastian M. Gygli^{1, 2}, Serej D
4 Ley^{1, 2, 4, 5}, Bijaya Malla^{1, 2, 6}, Julia Feldmann^{1, 2}, Sonia Borrell^{1, 2}, Christian Beisel⁷, Kerren
5 Middelkoop^{8, 9}, E Jane Carter^{10, 11}, Lameck Diero¹¹, Marie Ballif¹², Levan Jugheli^{1, 2}, Klaus
6 Reither^{1, 2}, Lukas Fenner^{1, 2, 12}, Daniela Brites^{1, 2##} and Sebastien Gagneux^{1, 2##}

7 ¹ Department of Medical Parasitology and Infection Biology, Swiss Tropical and Public Health
8 Institute, Basel, Switzerland

9 ² University of Basel, Basel, Switzerland

10 ³ Ifakara Health Institute, Bagamoyo, Tanzania

11 ⁴ Papua New Guinea Institute of Medical Research, Goroka, Papua New Guinea

12
13 ⁵ DST-NRF Centre of Excellence for Biomedical Tuberculosis Research; South African Medical
14 Research Council Centre for Tuberculosis Research; Division of Molecular Biology and Human
15 Genetics, Faculty of Medicine and Health Sciences, Stellenbosch University, Cape Town.

16
17 ⁶ Department of Radiation Oncology, Inselspital, Bern University Hospital,
18 University of Bern, Bern, Switzerland.

19 ⁷ Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

20 ⁸ Desmond Tutu HIV Centre, Institute of Infectious Disease and Molecular Medicine, Cape Town
21 South-Africa

22 ⁹ Department of Medicine, University of Cape Town, Cape Town South-Africa

23
24 ¹⁰ Division of Pulmonary and Critical Care Medicine, Warren Alpert School of Medicine at Brown
25 University, Providence, Rhode Island, USA

26
27 ¹¹ Moi University School of Medicine, Eldoret, Kenya

28
29 ¹² Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland

30
31 # Equal contribution

32 * **Correspondence:** Sebastien Gagneux, sebastien.gagneux@swisstph.ch; Daniela Brites, d.brites@swisstph.ch
33

34 **Keywords:** tuberculosis, introduction, genetic diversity, migration, whole genome sequencing, drug
35 resistance

36 **Abstract**

37 The Lineage 2–Beijing (L2–Beijing) sub-lineage of *Mycobacterium tuberculosis* has received much
38 attention due to its high virulence, fast disease progression, and association with antibiotic resistance.
39 Despite several reports of the recent emergence of L2–Beijing in Africa, no study has investigated
40 the evolutionary history of this sub-lineage on the continent. In this study, we used whole genome
41 sequences of 817 L2 clinical strains from 14 geographical regions globally distributed to investigate
42 the origins and onward spread of this lineage in Africa. Our results reveal multiple introductions of
43 L2–Beijing into Africa linked to independent bacterial populations from East- and Southeast Asia.
44 Bayesian analyses further indicate that these introductions occurred during the past 300 years, with
45 most of these events pre-dating the antibiotic era. Hence, the success of L2–Beijing in Africa is most
46 likely due to its hypervirulence and high transmissibility rather than drug resistance.

47 1 Introduction

48 Tuberculosis (TB) is mainly caused by a group of closely related bacteria referred to as the
49 *Mycobacterium tuberculosis* Complex (MTBC). The MTBC comprises seven phylogenetic lineages
50 adapted to humans and several lineages adapted to different wild and domestic animal species
51 (Gagneux, 2018). The human-adapted lineages of the MTBC show a distinct geographic distribution,
52 with some “generalist” lineages such as Lineage (L)2 and L4 occurring all around the world and
53 others being geographically restricted “specialist” that include L5, L6 and L7 (Coscolla and
54 Gagneux, 2014; Stucki *et al.*, 2016). Africa is the only continent which is home to all seven human-
55 adapted lineages, including the three “specialist” lineages exclusively found on the continent. Current
56 evidence suggests that the MTBC overall originated in Africa (Gagneux, 2018) and subsequently
57 spread around the globe following human migratory events (Wirth *et al.*, 2008; Comas *et al.*, 2013).
58 The broad distribution of some of the “generalist” lineages and their presence in Africa has been
59 attributed to past exploration, trade and conquest. For instance, an important part of the TB epidemics
60 in sub-Saharan Africa is driven by the generalist Latin–American–Mediterranean (LAM) sublineage
61 of L4, which is postulated to have been introduced to the continent post-European contact (Stucki *et*
62 *al.*, 2016).

63 Among the different human-adapted MTBC lineages, the L2–Beijing sublineage has been of
64 particular interest (Merker *et al.*, 2015). L2–Beijing has expanded and emerged worldwide from East
65 Asia; its most likely geographical origin (Luo *et al.*, 2015; Merker *et al.*, 2015). In some parts of the
66 world, the recent emergence of L2–Beijing has been linked to increased transmission (Yang *et al.*,
67 2012; Holt *et al.*, 2018), high prevalence of multidrug–resistant TB (MDR–TB) (Borrell and
68 Gagneux, 2009), and to social and political instability, resulting into displacement of people and poor
69 health systems (Eldholm *et al.*, 2016). Increasingly, L2–Beijing is also being reported in Africa
70 (Bifani *et al.*, 2002; Affolabi *et al.*, 2009; Gehre *et al.*, 2016; Mbugi *et al.*, 2016), and evidence
71 suggests that L2–Beijing in African regions is becoming more prevalent over time (Cowley *et al.*,
72 2008; van der Spuy *et al.*, 2009; Glynn *et al.*, 2010). Some authors have hypothesized that the
73 introduction of L2–Beijing into South Africa resulted from the importation of slaves from Southeast
74 Asia during the 17th and 18th centuries and/or the Chinese labor forces arriving in the 1900s (van
75 Helden *et al.*, 2002). Alternatively, in West Africa, the presence of L2–Beijing was proposed to
76 reflect more recent immigration from Asia (Affolabi *et al.*, 2009; Gehre *et al.*, 2016). To a certain
77 extent, the recent expansion of L2–Beijing in parts of Africa has been associated with drug resistance
78 (Githui *et al.*, 2004; Klopper *et al.*, 2013) and higher transmissibility (Guerra-Assunção *et al.*, 2015).
79 In addition, a study in the Gambia showed a faster progression from latent infection to active TB
80 disease in patient house–hold contacts exposed to L2–Beijing (de Jong *et al.*, 2008).

81 Whilst L2–Beijing seems to be expanding in several regions of Africa, no study has formally
82 investigated the evolutionary history of L2–Beijing on the continent. In this study, we used whole
83 genome sequencing data from a global collection of L2 clinical strains to determine the most likely
84 geographical origin of L2–Beijing in Africa and its spread across the continent.

85 2 Results

86 2.1 Phylogenetic inference of L2 strains

87 We analyzed a total of 817 L2 genomes originating from 14 geographical regions including Eastern
88 and Southern Africa (Figure S1 and Table S1). We focused on seven geographical regions that had
89 more than 20 genomes each, and assigned the remainder to “Other”, including two genomes from
90 Western Africa (Figure 1A). The resulting phylogeny of L2 was divided into two main sublineages:
91 the L2–proto-Beijing and L2–Beijing, supporting previous results (Luo *et al.*, 2015; Shitikov *et al.*,
92 2017). The L2–proto-Beijing was the most basal L2 sublineage and was restricted to East- and
93 Southeast Asia. L2–Beijing, particularly the “modern” (also known as “typical”) sublineage, was
94 geographically widely distributed and included strains from Africa. We further characterized L2–
95 Beijing using the recently described unified classification scheme for L2 (Shitikov *et al.*, 2017).

96 2.2 The population structure of L2–Beijing in Eastern and Southern Africa

97 Our findings showed the population of African L2–Beijing to be heterogeneous (Figure 1B, Figure 2
98 and Table S2). Most of the African L2–Beijing strains were classified into several groups within the
99 “modern” sublineage, which included primarily the “Asian-African” sublineages (L2.2.4, L2.2.5 and
100 L2.2.7), consistent with previous findings (Merker *et al.*, 2015). We also identified the “ancient”
101 (atypical) strains among the African L2–Beijing. Given that “ancient” L2–Beijing strains (L2.2.1 –
102 L2.2.3) are generally uncommon (Luo *et al.*, 2015), it is interesting to observe such strains in both
103 African regions. In several instances, African L2–Beijing strains did not fall into any of the
104 previously defined groups (Figure 2). Of the two African regions studied here, East Africa had higher
105 proportion of previously uncharacterized L2–Beijing strains (50/101, 50%).

106 In summary, our findings show that African regions harbored distinct L2–Beijing populations. This is
107 unlike Eastern Europe and Central Asia, where L2–Beijing is dominated by a few highly similar
108 strains (Casali *et al.*, 2014; Eldholm *et al.*, 2016). Of note, L2–Beijing strains typical of Central Asia
109 and Eastern Europe were completely absent from the African populations (Figure 2).

110 2.3 Genetic diversity of L2–Beijing strains across geographic regions

111 The spatial distribution of L2–Beijing sublineages and the prevalence of “ancient” L2–Beijing strains
112 observed in this study and previously (Luo *et al.*, 2015; Merker *et al.*, 2015), suggest that L2–Beijing
113 has expanded worldwide from Asia. This view can further be supported by the measures of genetic
114 diversity of L2–Beijing in the different geographical regions (Figure 3). As expected, East- and
115 Southeast Asia contained the most genetically diverse L2 populations, which is consistent with
116 previous results (Luo *et al.*, 2015). Conversely, L2 populations in other geographies were less
117 genetically diverse, suggesting recent dissemination of L2 to these regions. Within Africa, Southern
118 Africa showed a higher diversity in L2–Beijing populations compared to Eastern Africa.

119 The genetic diversity within the African L2–Beijing populations not only reflects the number and
120 variety of source populations but also local patterns of diversification that occurred after their

121 introduction. Therefore, the higher genetic diversity of the L2–Beijing populations in Southern Africa
122 compared to Eastern Africa likely reflects both aspects.

123 **2.4 Multiple introductions of L2–Beijing from Asia into Africa**

124 Based on our reconstructed phylogeny, African L2–Beijing strains clustered into several unrelated
125 clades indicating multiple introductions into Africa (Figure 1B). We next investigated the most likely
126 geographical origins of those introductions. As anticipated, our ancestral reconstruction estimated
127 East Asia as the most likely origin of all L2 (posterior probability of 96.7%) and L2–Beijing
128 (posterior probability 92.5%). Our data further indicate that L2–Beijing was introduced into Africa
129 from East- and Southeast Asia on multiple occasions independently. Furthermore, we observed both
130 direct introductions from Asia into Africa as well as subsequent dispersal within the continent
131 (Figure S2 and S3).

132 Using stochastic mapping, we estimated a total of 13 introductions or migration events (M1 – M13)
133 into Africa (Figure 4). Eight of the African L2–Beijing introductions originated from East Asia and
134 five from Southeast Asia. Out of the 13 introductions, three (M3, M10 and M13) were present in both
135 African regions analyzed here, suggesting initial introductions from Asia followed by subsequent
136 spread within Africa. Overall, our analysis inferred more independent introductions into Southern
137 Africa than Eastern Africa, seven (M1, M4, M7-9, M11 and M12) and three (M2, M5 and M6),
138 respectively. Taken together, our data suggest that multiple migration events have shaped the
139 populations of L2–Beijing in Africa.

140 **2.5 Bayesian molecular dating**

141 Different hypotheses have been formulated on the possible timing of the introduction of L2–Beijing
142 into Africa (van Helden *et al.*, 2002). Here we used tip-calibration to date the phylogenetic tree of L2
143 and estimate the age of its introduction to Africa. For these analyses, we identified 308 strains among
144 the 817 for which the sampling year was known. These strains were sampled during a period of 19
145 years; 1995 - 2014 (Figure S4), were evenly distributed on the complete phylogenetic tree (Figure
146 S5) and included 40% members of the African L2–Beijing strains (Figure S6). Eleven of the 13
147 African introductions were represented in this dataset (M1-M3 and M6-M13).

148 We detected no overlap in the 95% credibility interval of the clock rate estimates of observed and
149 randomized datasets indicating that there was sufficient temporal signal in the dataset to perform
150 inference (see methods, Figure S7). Further, We found that the UCLD clock had the highest marginal
151 likelihood and a Bayes Factor of 27 with the second best fitting model, the strict clock (Table 1),
152 indicating strong evidence in favor of the UCLD clock (Kass and Raftery, 1995).

153 We performed a phylogenetic analysis with BEAST2 using the UCLD clock, and repeated the tip
154 randomization test (see methods). We found that the 95% credibility interval of the clock rate
155 estimates of observed and randomized datasets did not overlap, confirming that there was a temporal
156 signal in our sequence data (Figure S8). Additionally, under the UCLD model, the coefficient of
157 variation (COV), which is a summary of the branch rates distribution (standard deviation divided by

158 the mean), gives an indication on the clock-likeness of the data (Drummond *et al.*, 2006). A
159 coefficient of variation of zero indicates that the data fit a strict clock, whilst a greater COV indicates
160 a higher heterogeneity of rates through the phylogeny. We obtained a mean COV of 0.22 (95%
161 credibility interval= 0.1732, 0.2732), indicating a moderate level of rate variation across different
162 branches and thus supporting the results of the path sampling analysis that favored the UCLD model.

163 **2.6 Recent origins of the African L2–Beijing clades**

164 We used the UCLD clock model to infer the clock rate and divergent times of the 308 L2 strains with
165 known sampling dates and estimated a mean substitution rate of 1.34×10^{-7} [95% Highest Posterior
166 Density (HPD), 9.2867×10^{-8} - 1.7719×10^{-7}]. These estimates are in agreement with previously
167 reported rates from epidemiological studies (Walker *et al.*, 2013; Eldholm *et al.*, 2015). We estimated
168 the most recent common ancestor (MRCA) of the extant L2–Beijing of the 308 strains to the year
169 1225 [95% HPD, 900 - 1519] (Figure S9). For each African clade, we estimated the year of
170 introduction using the 0.975 quantile of the HPD of the age of the MRCA as the upper limit (most
171 recent possible year) and the 0.025 quantile of the HPD of the divergence time between the closest
172 non-African L2–Beijing strain (the closest outgroup) and the African clade of interest as lowest limit
173 (most ancient possible year). Our estimates placed the earliest introductions of the African L2–
174 Beijing (M1, M3, M7, and M12) in the 18th and 19th century (Figure 5 and Table S3). Four additional
175 migration events (M6, M9, M10, and M11) were estimated to have occurred between the beginning
176 of the 19th century and the first half of the 20th century. Finally, the three most recent introductions to
177 Africa happened in the second half of the 20th century (M2, M8, and M13). Diversity patterns of the
178 African clades exclusive to Eastern- and Southern Africa could further provide support for the recent
179 introductions of African L2–Beijing. We thus calculated the pairwise SNP distances within the
180 individual introductions to explore the local patterns of diversification associated with regional
181 epidemics after the introductions. Although strains within Southern African introductions were
182 relatively more distantly related, L2–Beijing strains from both African regions were on average 20 to
183 40 SNPs apart (Figure S10 and S11). The latter thresholds were proposed to correspond to strains
184 involved in transmission clusters of estimated 50 to 100 years (Meehan *et al.*, 2018), supporting the
185 relatively recent introductions of L2–Beijing into the African continent.

186 Overall, these results indicate that the different introductions of L2–Beijing to Africa occurred over a
187 period of 300 years. While the earliest introduction is unlikely to have happened after 1732 - 1874,
188 the most recent is unlikely to have occurred before 1946 - 1980.

189 **2.7 Introductions of L2–Beijing into Africa unrelated to drug resistance**

190 Because of the repeated association of L2–Beijing with antibiotic resistance (Borrell and Gagneux,
191 2009), the emergence and dissemination of L2–Beijing strains has been attributed to drug resistance.
192 However, our estimated timing of these introductions suggest that African L2–Beijing strains were
193 introduced prior the discovery of TB antibiotics, and thus must have involved drug-susceptible
194 strains (Figure 5). To explore this question further, we assessed the drug resistance profiles of L2–
195 Beijing strains linked to the various introduction events into the two African regions. We found that

196 all the Eastern African populations contained only drug-susceptible strains and that approximately
197 three-quarters of L2–Beijing strains in the Southern African populations were drug-susceptible, with
198 the remaining being either mono–, multi– or extensively drug-resistant (Figure 6 and Figure S12).
199 Taken together, these results suggest that the emergence of L2–Beijing in Africa, particularly in
200 Eastern Africa, was not driven by drug resistance. Moreover, our data indicate independent
201 acquisition of drug resistance for the resistant strains detected in the Southern African L2–Beijing
202 population (Figure 6), which might partly contribute to the subsequent spread of L2–Beijing in
203 Southern Africa but not in Eastern Africa.

204 3 Discussion

205 This study investigated the most likely geographical origin of the L2–Beijing in Africa. In line with
206 previous findings (Luo *et al.*, 2015; Merker *et al.*, 2015), we identified East Asia as the most likely
207 place of origin of L2 and L2–Beijing. Our findings further revealed multiple independent
208 introductions of L2–Beijing into Africa linked to separate populations originating from both East-
209 and Southeast Asia. Some of these introductions were followed by further onward spread of L2–
210 Beijing within African regions. Finally, we demonstrate that most introductions of L2–Beijing on the
211 continent occurred before the age of antibiotics.

212 L2–Beijing has received much attention given its hypervirulence in infection models (Manca *et al.*,
213 2001; Ribeiro *et al.*, 2014), faster progression to disease and higher transmission potential in humans
214 (de Jong *et al.*, 2008; Holt *et al.*, 2018), frequent association with drug resistance, and recent
215 emergence in different regions of the world (Bifani *et al.*, 2002; Borrell and Gagneux, 2009; Fenner
216 *et al.*, 2013). Several studies indicate L2–Beijing originated in Asia and spread from there to the rest
217 of the world (Luo *et al.*, 2015; Merker *et al.*, 2015). Our results support this notion by identifying
218 “Asia” as the most likely geographical origin of both L2 and L2–Beijing based on our ancestral
219 reconstructions and the fact the L2–Beijing populations in Asia are much more diverse than in other
220 regions. In addition, our findings show that L2–Beijing was introduced into Africa multiple times
221 from both East- and Southeast Asia. The presence of L2–Beijing in South Africa has previously been
222 proposed to be due to the importation of slaves from Southeastern Asia by Europeans in the 17th and
223 18th centuries followed by the import of Chinese labor-forces in the early 1900s (van Helden *et al.*,
224 2002; Mokrousov *et al.*, 2005). Our Bayesian dating estimates predicted the earliest introductions of
225 L2–Beijing into Africa to have occurred in the 18th and 19th centuries, concurring with these
226 proposed time periods. However, our findings also point to later introductions of L2–Beijing into the
227 continent in the 19th and early 20th centuries. The timings of the latest three introductions in the
228 second half of the 20th century coincide with the decolonization and post-colonial period in Africa
229 when investments into infrastructure and other projects by Chinese enterprises substantially increased
230 (Yuan, 2006; Rice, 2011). These activities also brought many Chinese workers to Africa during a
231 time when TB was still very prevalent in China (Murray, 2018). Hence, many of these workers were
232 likely latently infected with L2–Beijing and might have later reactivated (Moreira Pescarini *et al.*,
233 2017). Overall, our findings suggest that L2–Beijing has emerged in Africa over the last 300 years
234 and that these introductions have occurred sporadically ever since.

235 The repeated association of L2–Beijing with drug resistance (Borrell and Gagneux, 2009) has led
236 some to propose that drug resistance is another reason why this sublineage might successfully
237 compete against and eventually replace other *M. tuberculosis* genotypes (Parwati, Van Crevel and
238 Van Soolingen, 2010). However, the underlying reasons for the association of L2–Beijing with drug
239 resistance remains unclear (Borrell and Trauner, 2017), and it is also far from universal, with several
240 reports from e.g. China and other regions finding no such association (Hanekom *et al.*, 2007; Yang *et al.*,
241 2012). Our results show that most introduction events of L2–Beijing into Africa pre-date the
242 antibiotic era, and because of that these introductions were most likely caused by drug-susceptible
243 strains. The notion that the initial emergence of L2–Beijing in Africa was not driven by drug

244 resistance is further supported by our findings that none of L2–Beijing strains from Eastern Africa
245 strains analyzed here were drug-resistant. Of note, our observations suggest that drug resistance in
246 South Africa was acquired via independent events post initial introductions from Asia. This is in
247 sharp contrast to the situation in Eastern Europe and Central Asia, where L2–Beijing is highly
248 prevalent but dominated by few recently expanded drug-resistant clones, which account for up to
249 60% of the L2–Beijing populations in some of these countries (Casali *et al.*, 2014; Eldholm *et al.*,
250 2016). The association of L2–Beijing with drug resistance in these regions were likely favored by the
251 economic and public health crises that followed the collapse of Soviet Union (Luo *et al.*, 2015;
252 Merker *et al.*, 2015).

253 Based on our finding that the original introductions of L2–Beijing into Africa involved drug-
254 susceptible strains and that the prevalence of drug-resistant L2–Beijing in Africa overall is
255 comparably low (WHO, 2017), we propose that some of the other characteristics of this sub-lineage,
256 in particular its high virulence, high transmissibility and rapid progression from infection to disease,
257 were responsible for the initial competitive success of L2–Beijing in Africa. Given that the MTBC
258 overall originated in Africa (Wirth *et al.*, 2008; Comas *et al.*, 2013), TB epidemics on the continent
259 were caused by many different “native” genotypes prior to foreign contacts (Comas *et al.*, 2015). The
260 emergence and expansion of “foreign” genotypes including L2–Beijing post-contact demonstrate
261 their ability to successfully compete against the existing genotypes on the continent, irrespective of
262 drug resistance. Following their initial establishment, poor TB treatment programs subsequently
263 selected for drug resistance in L2–Beijing but also in other MTBC lineages, which might have
264 facilitated their further spread in countries such as South Africa (Müller *et al.*, 2013).

265 This study is limited by the fact that we analyzed a globally diverse collection of L2 genomes
266 available in public repositories. Hence, these strains might not be fully representative of the
267 respective geographical regions. Moreover, our African L2–Beijing dataset came from convenient
268 sampling and comprised L2–Beijing mainly from Eastern and Southern Africa, as whole genome
269 data of L2–Beijing from the other African regions were unavailable at the time of the study.
270 However, the representation of African L2–Beijing in our sample reflects the overall prevalence of
271 this sub-lineage as recently reported for the continent (Mbugi *et al.*, 2016; Chihota *et al.*, 2018).
272 Moreover, although regions outside of Eastern- and Southern Africa were underrepresented, this is
273 unlikely to invalidate our findings regarding the multiple independents of L2–Beijing into Africa,
274 except by underestimating the number of true introductions.

275 In conclusion, this is the first study to address the geographical origins of L2–Beijing in Africa using
276 whole genome sequencing data. Our findings indicate multiple independent introductions of L2–
277 Beijing epidemics into Africa from East- and Southeast Asia during the last 300 years that were
278 unrelated to drug resistance. The TB epidemics in Africa have remained fairly stable over the last few
279 decades (WHO, 2017). However, Africa’s population growth and increasing urbanization (driven by
280 booming economies) are likely to have an impact on the future of TB in this continent, whether
281 directly by e.g. facilitating transmission or indirectly by promoting new risk factors such as diabetes
282 that increase TB susceptibility (Dye and Williams, 2010). It is therefore crucial to follow the TB
283 epidemics in the continent very closely, especially those related to hypervirulent strains such as L2–

284 Beijing, as these might take particular advantage of this expanding ecological niche (Cowley *et al.*,
285 2008).

286 4 Material and Methods

287 4.1 Identification of Lineage 2 strains and whole-genome sequencing

288 We obtained whole-genome sequencing data of L2 strains from the two previously largest studies
289 focusing on the evolutionary history and global spread of L2–Beijing strains (Luo *et al.*, 2015;
290 Merker *et al.*, 2015). We then identified additional published genomes as African representatives of
291 L2–Beijing strains from other studies (Guerra-Assunção *et al.*, 2015; Manson *et al.*, 2017).
292 Moreover, we newly sequenced 116 additional L2–Beijing strains using Illumina HiSeq 2000/2500
293 paired end technology ([PRJNA488343](#)). In total, we included 817 L2 genome sequences (Figure S1
294 and Table S1).

295 4.2 Whole genome sequence analysis and phylogenetic inference

296 We used a customized pipeline previously described to map short sequencing reads with BWA 0.6.2
297 to a reconstructed hypothetical MTBC ancestor used as reference (Comas *et al.*, 2013). SAMtools
298 0.1.19 was used to call single nucleotide polymorphisms (SNPs), and these SNPs were annotated
299 using ANNOVAR and customized scripts based on the *M. tuberculosis* H37Rv reference annotation
300 (AL123456.2). For downstream analyses, we excluded SNPs in repetitive regions, those annotated in
301 problematic regions such as ‘PE/PPE/PGRS’ and SNPs in drug-resistance associated genes. Small
302 insertions and deletions were also excluded from the analyses. Only SNPs with minimum coverage
303 of 20x and minimum mapping quality of 30 were kept. All SNPs classified by Samtools as having
304 frequencies of the major non-reference allele lower than 100% ($AF1 < 1$) within each genome were
305 considered to be heterogeneous and were treated as ambiguities, and otherwise considered fixed
306 ($AF1 = 1$). We concatenated fixed SNPs from the variable positions obtained, which yielded a 33,776
307 bp alignment. The alignment was then used to infer a maximum likelihood phylogeny using RAxML
308 8.3.2 with a general time reversible (GTR) model in RAxML and 1,000 rapid bootstrap inferences,
309 followed by a thorough maximum-likelihood search (Stamatakis, 2006).

310 4.3 Phylogeographic analyses

311 4.3.1 Reconstruction of ancestral state

312 To investigate the likely geographic origin of L2–Beijing strains in Africa, we inferred the historical
313 biogeography of L2 using the RASP software (Yu *et al.*, 2015) on a representative subset of 430
314 genomes due to software’s sample limitation. We achieved this subset by performing hierarchical
315 clustering implemented in *pvclust* package in R (Suzuki and Shimodaira, 2006) on a distance matrix
316 and randomly removing clustered genomes. We applied a Bayesian based method in RASP to
317 reconstruct ancestral geographical states on the phylogeny of 430 L2 genomes. We used geographical
318 regions as proxy for origins of the L2 strains and loaded them as distributions. We then ran Bayesian
319 analysis with 5 chains and 500 generations.

320 4.3.2 Stochastic character mapping

321 To determine the number of introduction events of L2–Beijing into African regions, we applied
322 stochastic mapping (SIMMAP) on the 817 L2 phylogeny using phytools package 0.6.44 in R (Revell,

323 2012). Geographical origin of the L2 strains was treated as a discrete trait and modeled onto the
324 phylogeny using ARD model with 100 replicates. This model permits independent region-to-region
325 transfer. We referred to the resulting introductions as migration events “M”, considering only those
326 introductions with more than 5 genomes.

327 **4.3.3 Population genetic analyses**

328 **4.3.3.1 Nucleotide diversity (π)**

329 We calculated the mean pair-wise nucleotide diversity per site (π) measured by geographic region.
330 We excluded geographic regions represented by less than 20 genomes. Confidence intervals were
331 obtained by bootstrapping through resampling using the sample function in R with replacement and
332 the respective lower and upper confidence levels by calculating 2.5th and 97.5th quartiles.

333 **4.3.3.2 Pairwise SNP distances**

334 We used *dist.dna* function of ape package implemented in R (Paradis, Claude and Strimmer, 2004) to
335 calculate pairwise SNP distances with raw mutation counts and pairwise deletions for gaps. Mean
336 pairwise SNP distance to all strains of the same geographic population was calculated per strain and
337 the distribution of the mean SNP pairwise distance plotted. The mean pairwise SNP distances were
338 assumed not to be normally distributed and we therefore used Wilcoxon rank-sum test to test the
339 differences among geographic regions. Additionally, we calculated pairwise SNP distances within
340 African L2–Beijing populations for migration events with more than 10 genomes each.

341 **4.3.4 Drug Resistance**

342 To distinguish between drug-susceptible and drug-resistant strains, we used genotypic drug resistance
343 molecular markers previously described (Steiner *et al.*, 2014). We categorized strains into:
344 susceptible as having no drug resistance specific mutations; monoresistant as having mutations
345 conferring resistance to a single drug; MDR as having mutations conferring resistance to isoniazid
346 and rifampicin; and extensively drug-resistant (XDR) as having mutations conferring resistance to
347 fluoroquinolones and aminoglycosides in addition to being MDR (Table S4).

348 **4.3.5 Bayesian molecular dating**

349 **4.3.5.1 Data preparation and preliminary analysis**

350 To estimate the historical period in which L2–Beijing was introduced to Africa, we performed a set
351 of Bayesian phylogenetic analyses using tip-calibration (Rieux and Khatchikian, 2017). Among the
352 817 studied L2 strains, we had information on the year of sampling for 308. We performed all further
353 analysis on this subset of 308 strains. We excluded all genomic positions that were invariable in this
354 subset and all positions that were undetermined (missing data or deletions) in more than 25% of the
355 strains, and obtained an alignment of 10,769 polymorphic positions.

356 In tip dating analysis it is important to test whether the dataset contains strong enough temporal
357 signal (Rieux and Balloux, 2016). To do this, we performed a tip randomization test (Ramsden *et al.*,
358 2008) as follows. We used BEAST2 v. 2.4.8 (Bouckaert *et al.*, 2014) to run a phylogenetic analysis

359 with a HKY + GAMMA model (Hasegawa, Kishino and Yano, 1985), a constant population size
360 prior on the tree and a strict molecular clock. Additionally, we used the years in which the strains
361 were sampled to time-calibrate the tree, and we modified the extensible markup language (xml) file
362 to specify the number of invariant sites as indicated by the developers of BEAST2 here:
363 <https://groups.google.com/forum/#!topic/beast-users/QfBHMOqImFE> (*strict_preliminary.xml*). We
364 ran three independent runs (245 million generations in total), and we used Tracer 1.7 (Rambaut *et al.*,
365 2018) to identify the burn-in (8 million generations), to assess that the different runs converged, and
366 to estimate the effective sample size (ESS) for all parameters, the posterior and the likelihood (ESS >
367 110 for all parameters). We then used TipDatingBeast (Rieux and Khatchikian, 2017) to generate 20
368 replicates of the xml file in which the sampling dates were randomly reassigned to different strains.
369 For each replication, we ran the same BEAST2 analysis as for the original (observed) dataset (one
370 run per replicate, 50 million generations, 10% burn-in). We used TipDatingBeast to parse the log
371 files output of BEAST2 and compare the clock rate estimates for the observed data and the
372 randomized replications. The estimates of the molecular clock rate did not overlap between the
373 observed and the randomized dataset, indicating that there is a clear temporal signal and that we
374 could proceed with further analysis (Figure S5).

375 4.3.5.2 Model selection

376 To identify the model that best fits the data, we estimated the marginal likelihood of three different
377 clock models: UCED and UCLD (Drummond *et al.*, 2006). We used the Model selection package of
378 BEAST2 to run a path sampling analysis (Lartillot and Philippe, 2006) following the
379 recommendations of the BEAST2 developers (<http://www.beast2.org/path-sampling/>). We used the
380 following settings: 100 steps, 4 million generations per step, alpha = 0.3, pre-burn-in = 1 million
381 generations, burn-in for each step= 40% (**PS.xml*). For these analyses, we used proper priors as
382 suggested by (Baele *et al.*, 2012).

383 4.3.5.3 UCLD analysis

384 Since the model selection analysis indicated that the UCLD clock was the best fitting model, we
385 repeated the analysis using the UCLD and the same settings used in the path sampling analysis,
386 sampling every 10,000 generations. We ran three independent runs (800 million generations in total),
387 we used Tracer 1.7 (Rambaut *et al.*, 2018) to identify the burn-in (10 million generations), to assess
388 that the different runs converged and to estimate the effective sample size (ESS) for all parameters,
389 the posterior and the likelihood (ESS > 260 for all parameters) (*UCLD_final.xml* and Table S5)

390 We checked the sensitivity to the priors by running one analysis of 250 million generation sampling
391 from the prior, and compared the parameter estimates with the analysis using the data. We observed
392 the posterior distribution and the prior distribution of all parameters are very distinct (Table S6),
393 indicating that the parameter estimates are influenced by the data and not by the priors (Bromham *et*
394 *al.*, 2018).

395 We repeated the tip randomization test with the UCLD model as described above (20 replicates, one
396 run per replicate, 105 million generations per replicate or more, burn-in 10%), and again we found a
397 temporal signal (Figure S8).

398 To summarize the results, we sampled the trees from the three runs (5% burn-in corresponding to 10
399 million generations or more, sampling every 25,000 generation). We then summarized the 31,758
400 sampled trees, created a maximum clade credibility tree using the software TreeAnnotator from the
401 BEAST2 package and used FigTree version 1.4.2 for visualization (Figure S9).

402 **5 Conflict of Interest**

403 We declare no conflict of interest.

404 **6 Author Contributions**

405 LKR, DB, FM, DS, LF and SG planned the study, SDL, BM and JF performed the experiments,
406 LKR, DB, FM, SMG, SDL, BM, CB, SB, KM, MB, LJ, KR and LF contributed strains and prepared
407 the data, LKR, DB, FM and SG analyzed the data, LKR, DB, FM and SG drafted the manuscript. All
408 authors critically reviewed the manuscript.

409 **7 Acknowledgments and funding information**

410 We would like to thank Sebastián Duchêne and Yan Yu for their technical support and Linda-Gail
411 Bekker for contributing strains. All bioinformatics analyses were performed at the scientific
412 computing core facility of the University of Basel, sciCORE (<http://scicore.unibas.ch/>). This work
413 was supported by the Swiss National Science Foundation (grants 310030_166687 to SG), the
414 European Research Council (309540-EVODRTB to SG) and SystemsX.ch This research was also
415 partially supported (strain collection) by a funding supplement from the National Institutes of Allergy
416 and Infectious Diseases (NIAID) under award numbers U01 AI069924 (IeDEA Southern Africa) and
417 U01 AI069911 (IeDEA East Africa).

418 **8 References**

- 419 Affolabi, D. *et al.* (2009) ‘Possible Outbreak of Streptomycin-Resistant *Mycobacterium tuberculosis*
420 Beijing in Benin’, *Emerging Infectious Diseases*, 15(7), pp. 1123–1125. doi:
421 10.3201/eid1507.080697.
- 422 Baele, G. *et al.* (2012) ‘Accurate Model Selection of Relaxed Molecular Clocks in Bayesian
423 Phylogenetics’, *Molecular Biology and Evolution*. Oxford University Press, 30(2), pp. 239–243. doi:
424 10.1093/molbev/mss243.
- 425 Bifani, P. J. *et al.* (2002) ‘Global dissemination of the *Mycobacterium tuberculosis* W-Beijing family
426 strains’, *Trends in Microbiology*, 10(1), pp. 45–52. doi: 10.1016/S0966-842X(01)02277-6.
- 427 Borrell, S. and Gagneux, S. (2009) ‘Infectiousness, reproductive fitness and evolution of drug-
428 resistant *Mycobacterium tuberculosis*’, *Int J Tuberc Lung Dis*, 13(12), pp. 1456–1466.
- 429 Borrell, S. and Trauner, A. (2017) ‘Strain Diversity and the Evolution of Antibiotic Resistance’, in
430 Gagneux, S. (ed.) *Strain Variation in the Mycobacterium tuberculosis Complex: Its Role in Biology,*
431 *Epidemiology and Control*. Cham: Springer International Publishing, pp. 263–279. doi: 10.1007/978-
432 3-319-64371-7_14.
- 433 Bouckaert, R. *et al.* (2014) ‘BEAST 2: A Software Platform for Bayesian Evolutionary Analysis’,
434 *PLoS Comput Biol*, 10(4), p. 1003537. doi: 10.1371/journal.pcbi.1003537.
- 435 Bromham, L. *et al.* (2018) ‘Bayesian molecular dating: opening up the black box’, *Biological*
436 *Reviews*. Wiley/Blackwell (10.1111), 93(2), pp. 1165–1191. doi: 10.1111/brv.12390.
- 437 Casali, N. *et al.* (2014) ‘Evolution and transmission of drug-resistant tuberculosis in a Russian
438 population’, *Nature Publishing Group*, 46(3), pp. 279–86. doi: 10.1038/ng.2878.
- 439 Chihota, V. N. *et al.* (2018) ‘Geospatial distribution of *Mycobacterium tuberculosis* genotypes in
440 Africa’, *PLOS ONE*, 13(8), p. e0200632. doi: 10.1371/journal.pone.0200632.
- 441 Comas, I. *et al.* (2013) ‘Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium*
442 *tuberculosis* with modern humans’, *Nature Publishing Group*, 45(10). doi: 10.1038/ng.2744.
- 443 Comas, I. *et al.* (2015) ‘Population Genomics of *Mycobacterium tuberculosis* in Ethiopia Contradicts
444 the Virgin Soil Hypothesis for Human Tuberculosis in Sub-Saharan Africa’, *Current Biology*, 25(24),
445 pp. 3260–3266. doi: 10.1016/j.cub.2015.10.061.
- 446 Coscolla, M. and Gagneux, S. (2014) ‘Consequences of genomic diversity in *Mycobacterium*
447 *tuberculosis*’, *Seminars in Immunology*. (Immunity to *Mycobacterium tuberculosis*), 26(6), pp. 431–
448 444. doi: 10.1016/j.smim.2014.09.012.
- 449 Cowley, D. *et al.* (2008) ‘Recent and Rapid Emergence of W-Beijing Strains of *Mycobacterium*
450 *tuberculosis* in Cape Town, South Africa’, *Clinical Infectious Diseases*. Oxford University Press,
451 47(10), pp. 1252–1259. doi: 10.1086/592575.
- 452 Drummond, A. J. *et al.* (2006) ‘Relaxed phylogenetics and dating with confidence’, *PLoS Biol*, 4(5),
453 p. 88. doi: 10.1371/journal.pbio.0040088.

- 454 Dye, C. and Williams, B. G. (2010) ‘The Population Dynamics and Control of Tuberculosis’,
455 *Science*, 328(5980), pp. 856–861. doi: 10.1126/science.1185449.
- 456 Eldholm, V. *et al.* (2015) ‘Four decades of transmission of a multidrug-resistant Mycobacterium
457 tuberculosis outbreak strain’, *Nature Communications*, 6(1), p. 7119. doi: 10.1038/ncomms8119.
- 458 Eldholm, V. *et al.* (2016) ‘Armed conflict and population displacement as drivers of the evolution
459 and dispersal of Mycobacterium tuberculosis’, *Proceedings of the National Academy of Sciences*.
460 National Academy of Sciences, 113(48), pp. 13881–13886. doi: 10.1073/PNAS.1611283113.
- 461 Fenner, L. *et al.* (2013) ‘HIV Infection Disrupts the Sympatric Host–Pathogen Relationship in
462 Human Tuberculosis’, *PLoS Genet*, 9(3), p. e1003318. doi: 10.1371/journal.pgen.1003318.
- 463 Gagneux, S. (2018) ‘Ecology and evolution of Mycobacterium tuberculosis’, *Nature Publishing
464 Group*, 16. doi: 10.1038/nrmicro.2018.8.
- 465 Gehre, F. *et al.* (2016) ‘A Mycobacterial Perspective on Tuberculosis in West Africa: Significant
466 Geographical Variation of *M. africanum* and Other *M. tuberculosis* Complex Lineages’, *PLoS Negl
467 Trop Dis*, 10(3), p. e0004408. doi: 10.1371/journal.pntd.0004408.
- 468 Githui, W. A. *et al.* (2004) ‘Identification of MDR-TB Beijing/W and other Mycobacterium
469 tuberculosis genotypes in Nairobi, Kenya.’, *The international journal of tuberculosis and lung
470 disease : the official journal of the International Union against Tuberculosis and Lung Disease*, 8(3),
471 pp. 352–60.
- 472 Glynn, J. R. *et al.* (2010) ‘Changes in Mycobacterium tuberculosis Genotype Families Over 20 Years
473 in a Population Based Study in Northern Malawi’, *PLoS ONE*, 5(8), p. e12259. doi:
474 10.1371/journal.pone.0012259.
- 475 Guerra-Assunção, J. *et al.* (2015) ‘Large-scale whole genome sequencing of *M. tuberculosis* provides
476 insights into transmission in a high prevalence area’, *Elife*, 4, pp. 285–292. doi: 10.7554/eLife.05166.
- 477 Hanekom, M. *et al.* (2007) ‘A recently evolved sublineage of the Mycobacterium tuberculosis
478 Beijing strain family is associated with an increased ability to spread and cause disease.’, *Journal of
479 clinical microbiology*. American Society for Microbiology, 45(5), pp. 1483–90. doi:
480 10.1128/JCM.02191-06.
- 481 Hasegawa, M., Kishino, H. and Yano, T. (1985) ‘Dating of the human-ape splitting by a molecular
482 clock of mitochondrial DNA.’, *Journal of molecular evolution*, 22(2), pp. 160–74.
- 483 van Helden, P. D. *et al.* (2002) ‘Strain families of Mycobacterium tuberculosis’, *Trends in
484 Microbiology*, 10(4), pp. 167–168. doi: 10.1016/S0966-842X(02)02317-X.
- 485 Holt, K. E. *et al.* (2018) ‘Frequent transmission of the Mycobacterium tuberculosis Beijing lineage
486 and positive selection for the EsxW Beijing variant in Vietnam’, *Nature Genetics*, 50(6), pp. 849–
487 856. doi: 10.1038/s41588-018-0117-9.
- 488 de Jong, B. C. *et al.* (2008) ‘Progression to active tuberculosis, but not transmission, varies by
489 Mycobacterium tuberculosis lineage in The Gambia.’, *The Journal of infectious diseases*, 198(7), pp.
490 1037–43. doi: 10.1086/591504.

- 491 Kass, R. E. and Raftery, A. E. (1995) ‘Bayes Factors’, *Journal of the American Statistical*
492 *Association*, 90(430), pp. 773–795. doi: 10.1080/01621459.1995.10476572.
- 493 Klopper, M. *et al.* (2013) ‘Emergence and Spread of Extensively and Totally Drug-Resistant
494 Tuberculosis, South Africa’, *Emerging Infectious Diseases*, 19(3), pp. 449–455. doi:
495 10.3201/EID1903.120246.
- 496 Lartillot, N. and Philippe, H. (2006) ‘Computing Bayes Factors Using Thermodynamic Integration’,
497 *Systematic Biology*. Edited by P. Lewis. Oxford University Press, 55(2), pp. 195–207. doi:
498 10.1080/10635150500433722.
- 499 Luo, T. *et al.* (2015) ‘Southern East Asian origin and coexpansion of Mycobacterium tuberculosis
500 Beijing family with Han Chinese’, *PNAS*, p. 201424063. doi: 10.1073/pnas.1424063112.
- 501 Manca, C. *et al.* (2001) ‘Virulence of a Mycobacterium tuberculosis clinical isolate in mice is
502 determined by failure to induce Th1 type immunity and is associated with induction of IFN- α/β ’,
503 *Proceedings of the National Academy of Sciences*, 98(10), pp. 5752–5757.
- 504 Manson, A. L. *et al.* (2017) ‘Genomic analysis of globally diverse Mycobacterium tuberculosis
505 strains provides insights into the emergence and spread of multidrug resistance’, *Nature GeNetics*,
506 49(3). doi: 10.1038/ng.3767.
- 507 Mbugi, E. V. *et al.* (2016) ‘Mapping of Mycobacterium tuberculosis Complex Genetic Diversity
508 Profiles in Tanzania and Other African Countries’, *PLOS ONE*. Edited by S. Sreevatsan. World
509 Bank, 11(5), p. e0154571. doi: 10.1371/journal.pone.0154571.
- 510 Meehan, C. J. *et al.* (2018) ‘The relationship between transmission time and clustering methods in 1
511 Mycobacterium 2’, *bioRxiv*. doi: 10.1101/302232.
- 512 Merker, M. *et al.* (2015) ‘Evolutionary history and global spread of the Mycobacterium tuberculosis
513 Beijing lineage’, *Nature Publishing Group*, 47. doi: 10.1038/ng.3195.
- 514 Mokrousov, I. *et al.* (2005) ‘Origin and primary dispersal of the Mycobacterium tuberculosis Beijing
515 genotype: clues from human phylogeography.’, *Genome research*. Cold Spring Harbor Laboratory
516 Press, 15(10), pp. 1357–64. doi: 10.1101/gr.3840605.
- 517 Moreira Pescarini, J. *et al.* (2017) ‘Migration to middle-income countries and tuberculosis-global
518 policies for global economies’, *Globalization and Health*, 13(1), p. 15. doi: 10.1186/s12992-017-
519 0236-6.
- 520 Müller, B. *et al.* (2013) ‘The heterogeneous evolution of multidrug-resistant Mycobacterium
521 tuberculosis’, *Trends in Genetics*, 29(3), pp. 160–169. doi: 10.1016/j.tig.2012.11.005.
- 522 Murray, J. F. (2018) ‘Tuberculosis in China before, during, and after the Sino-Japanese War’, in
523 *Tuberculosis and War*. Karger Publishers, pp. 204–212. doi: 10.1159/000481489.
- 524 Paradis, E., Claude, J. and Strimmer, K. (2004) ‘APE: Analyses of Phylogenetics and Evolution in R
525 language’, *Bioinformatics*. Oxford University Press, 20(2), pp. 289–290. doi:
526 10.1093/bioinformatics/btg412.

- 527 Parwati, I., Van Crevel, R. and Van Soolingen, D. (2010) ‘Possible underlying mechanisms for
528 successful emergence of the Mycobacterium tuberculosis Beijing genotype strains’. doi:
529 10.1016/S1473-3099(09)70330-5.
- 530 Rambaut, A. *et al.* (2018) ‘Software for Systematics and Evolution Posterior Summarization in
531 Bayesian Phylogenetics Using Tracer 1.7’, *Systematic Biology*, 0(0), pp. 1–5. doi:
532 10.1093/sysbio/syy032.
- 533 Ramsden, C. *et al.* (2008) ‘High Rates of Molecular Evolution in Hantaviruses’, *Molecular Biology
534 and Evolution*, 25(7), pp. 1488–1492. doi: 10.1093/molbev/msn093.
- 535 Revell, L. J. (2012) ‘phytools: an R package for phylogenetic comparative biology (and other
536 things)’, *Methods in Ecology and Evolution*. Wiley/Blackwell (10.1111), 3(2), pp. 217–223. doi:
537 10.1111/j.2041-210X.2011.00169.x.
- 538 Ribeiro, S. C. M. *et al.* (2014) ‘Mycobacterium tuberculosis strains of the modern sublineage of the
539 Beijing family are more likely to display increased virulence than strains of the ancient sublineage’,
540 *J. Clin. Microbiol.*, 52(7), pp. 2615–2624. doi: 10.1128/JCM.00498-14.
- 541 Rice, X. (2011) ‘China ’ s economic invasion of Africa’, *Africa*, pp. 1–8.
- 542 Rieux, A. and Balloux, F. (2016) ‘Inferences from tip-calibrated phylogenies: a review and a
543 practical guide’, *Molecular Ecology*, 25(9), pp. 1911–1924. doi: 10.1111/mec.13586.
- 544 Rieux, A. and Khatchikian, C. E. (2017) ‘tipdatingbeast : an r package to assist the implementation of
545 phylogenetic tip-dating tests using beast’, *Molecular Ecology Resources*. Wiley/Blackwell (10.1111),
546 17(4), pp. 608–613. doi: 10.1111/1755-0998.12603.
- 547 Shitikov, E. *et al.* (2017) ‘Evolutionary pathway analysis and unified classification of East Asian
548 lineage of Mycobacterium tuberculosis’, *Scientific Reports*, 7(1), p. 9227. doi: 10.1038/s41598-017-
549 10018-5.
- 550 van der Spuy, G. D. *et al.* (2009) ‘Changing Mycobacterium tuberculosis population highlights
551 clade-specific pathogenic characteristics’, *Tuberculosis*. Churchill Livingstone, 89(2), pp. 120–125.
552 doi: 10.1016/J.TUBE.2008.09.003.
- 553 Stamatakis, A. (2006) ‘RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
554 thousands of taxa and mixed models’, *Bioinformatics*, 22(21), pp. 2688–2690. doi:
555 10.1093/bioinformatics/btl446.
- 556 Steiner, A. *et al.* (2014) ‘KvarQ: targeted and direct variant calling from fastq reads of bacterial
557 genomes’, *BMC Genomics*, 15(1), p. 881. doi: 10.1186/1471-2164-15-881.
- 558 Stucki, D. *et al.* (2016) ‘Mycobacterium tuberculosis lineage 4 comprises globally distributed and
559 geographically restricted sublineages’, *Nature Genetics*, 48(12). doi: 10.1038/ng.3704.
- 560 Suzuki, R. and Shimodaira, H. (2006) ‘Pvclust: an R package for assessing the uncertainty in
561 hierarchical clustering’, *Bioinformatics*. Oxford University Press, 22(12), pp. 1540–1542. doi:
562 10.1093/bioinformatics/btl117.

- 563 Walker, T. M. *et al.* (2013) ‘Whole-genome sequencing to delineate *Mycobacterium tuberculosis*
564 outbreaks: a retrospective observational study.’, *The Lancet. Infectious diseases*, 13(2), pp. 137–46.
565 doi: 10.1016/S1473-3099(12)70277-3.
- 566 WHO (2017) *Global tuberculosis report. Geneva: World Health Organization.*
- 567 Wirth, T. *et al.* (2008) ‘Origin, Spread and Demography of the *Mycobacterium tuberculosis*
568 Complex’, *PLoS Pathog*, 4(9), p. e1000160. doi: 10.1371/journal.ppat.1000160.
- 569 Yang, C. *et al.* (2012) ‘*Mycobacterium tuberculosis* Beijing Strains Favor Transmission but Not
570 Drug Resistance in China’, *Clinical Infectious Diseases*, 55(9), pp. 1179–1187. doi:
571 10.1093/cid/cis670.
- 572 Yu, Y. *et al.* (2015) ‘RASP (Reconstruct Ancestral State in Phylogenies): A tool for historical
573 biogeography’. doi: 10.1016/j.ympev.2015.03.008.
- 574 Yuan, W. (2006) *La Chine et l’Afrique : [1956-2006]*. [S. l.]: China Intercontinental Press.
- 575

576 **Data Availability Statement**

577 The additional datasets i.e. additional xml files and scripts for this study can be found here
578 https://github.com/SwissTPH/TBRU_L2Africa.git

579 **Figure Legends**

580 **Figure 1.** Global phylogeny and geographical distribution of L2 strains. **(A)** Geographical origin
581 (according to United Nations geoscheme) for the 817 L2 strains. The geographical origins with less
582 than 20 strains are colored dark gray and those with missing data light gray. **(B)** Maximum likelihood
583 phylogeny inferred from 33,776 variable single nucleotide positions of the 817 strains. Taxa are
584 colored according to the geographical origin of the strains and the clades are highlighted according to
585 previously defined sublineages. L2 defining markers i.e. deletions (RD) are also mapped onto the
586 phylogeny.

587 **Figure 2.** Population structure of L2 strains. Frequency in proportions of L2 sub-lineages across
588 seven geographical regions.

589 **Figure 3.** Genetic diversity of L2 strains within geographical regions. **(A)** Nucleotide diversity (π)
590 per site of L2 strains by geography. Error bars are the 95% confidence intervals. **(B)** Pairwise genetic
591 SNP distance of L2 by geography (p values were obtained from Wilcoxon rank-sum tests). Each box
592 represents the 25% and 75% quartiles and the line denotes the median.

593 **Figure 4.** Stochastic mapping of the geographic origin on L2 phylogeny. **(A)** Maximum likelihood
594 phylogeny of the 817 MTB Lineage 2 strains. Branches are colored according to their geographical
595 region inferred from stochastic mapping of geographic origin of L2 strains onto the phylogeny. The
596 13 migration events to Africa (M1–M13) are indicated. **(B)** Proposed scenario for the multiple
597 introduction events of L2–Beijing into Africa. **(C)** Plot summarizing the number introduction events
598 to Eastern and Southern Africa from East- and Southeast Asia.

599 **Figure 5.** Estimated time in median ages for the introductions of African L2–Beijing (M1-M3 and
600 M6-M13). Introductions to Eastern Africa are colored in red and those to Southern Africa in blue.
601 Migration M10 contained L2–Beijing from both Southern and Eastern Africa. Dotted line marks the
602 year of first anti-TB drug discovery (1943). The error bars correspond to the 95% HPD.

603 **Figure 6.** Proportions of drug resistance profiles for L2 strains in seven geographical regions.

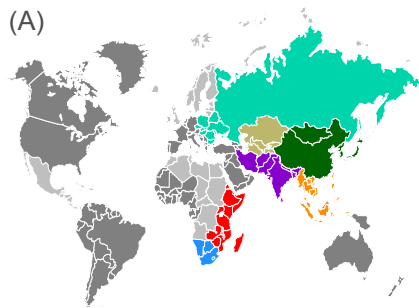
604 **Tables**

605 Table 1: Model selection based on path sampling Log-Marginal Likelihood

Clock model	Log(e) Marginal Likelihood	Bayes factor (UCLD vs model)
UCLD	-5374827	
Strict	-5374854	27
UCED	-5374897	70

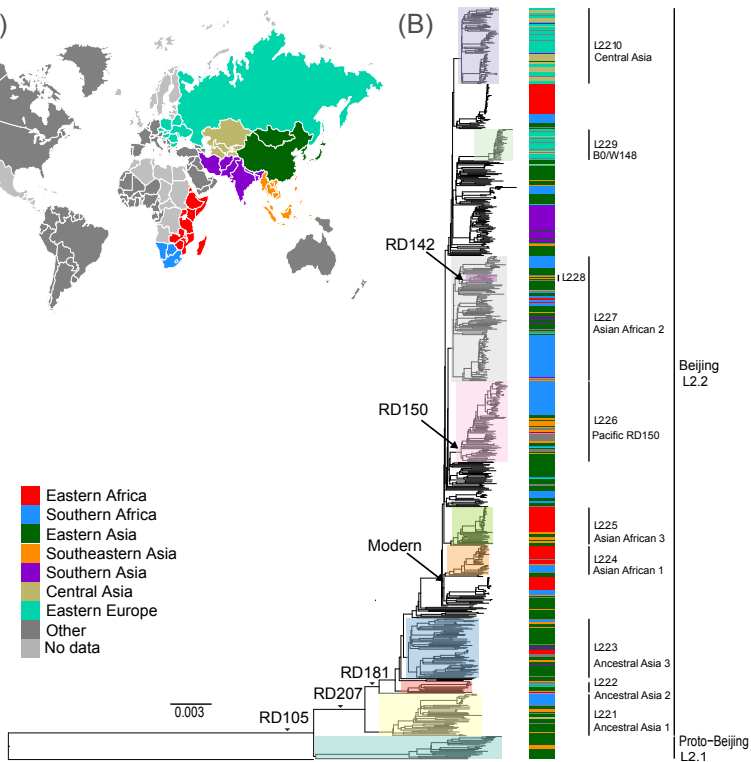
606

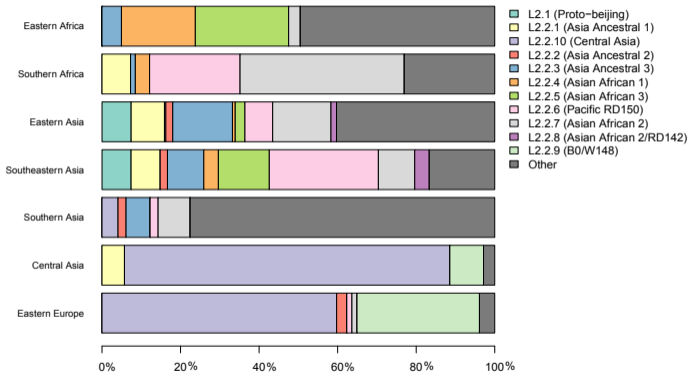
(A)

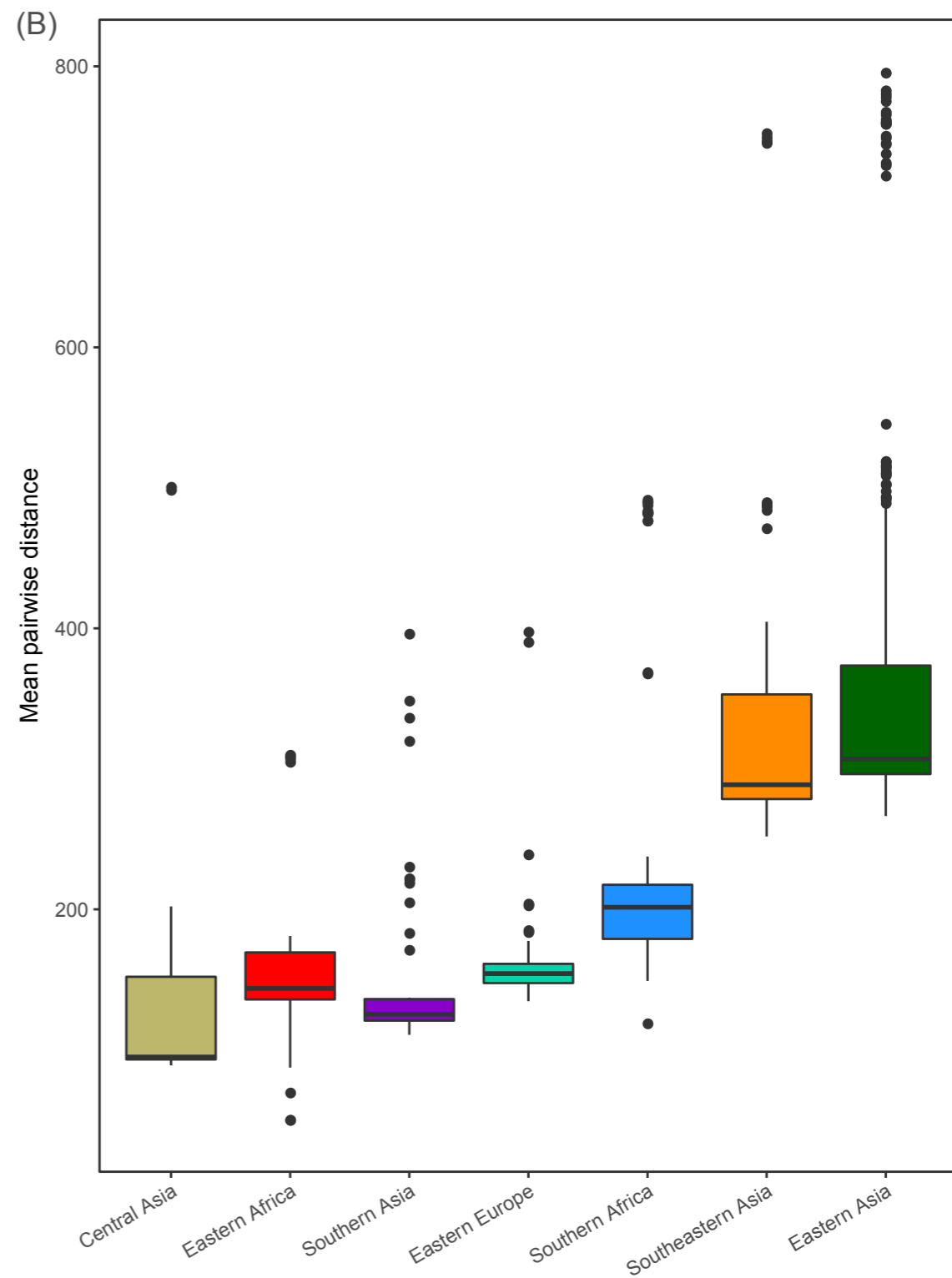
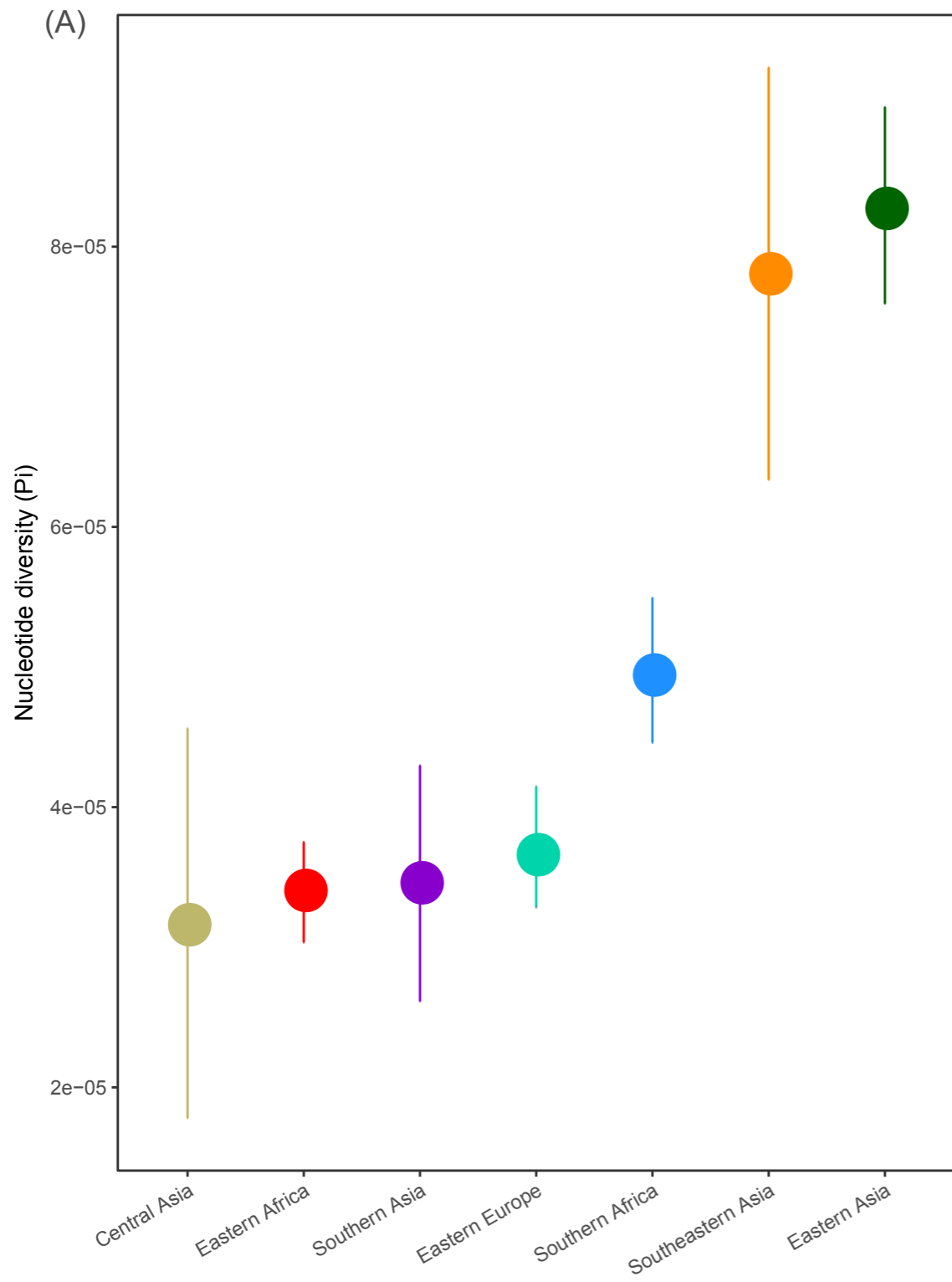


- Eastern Africa
- Southern Africa
- Eastern Asia
- Southeastern Asia
- Southern Asia
- Central Asia
- Eastern Europe
- Other
- No data

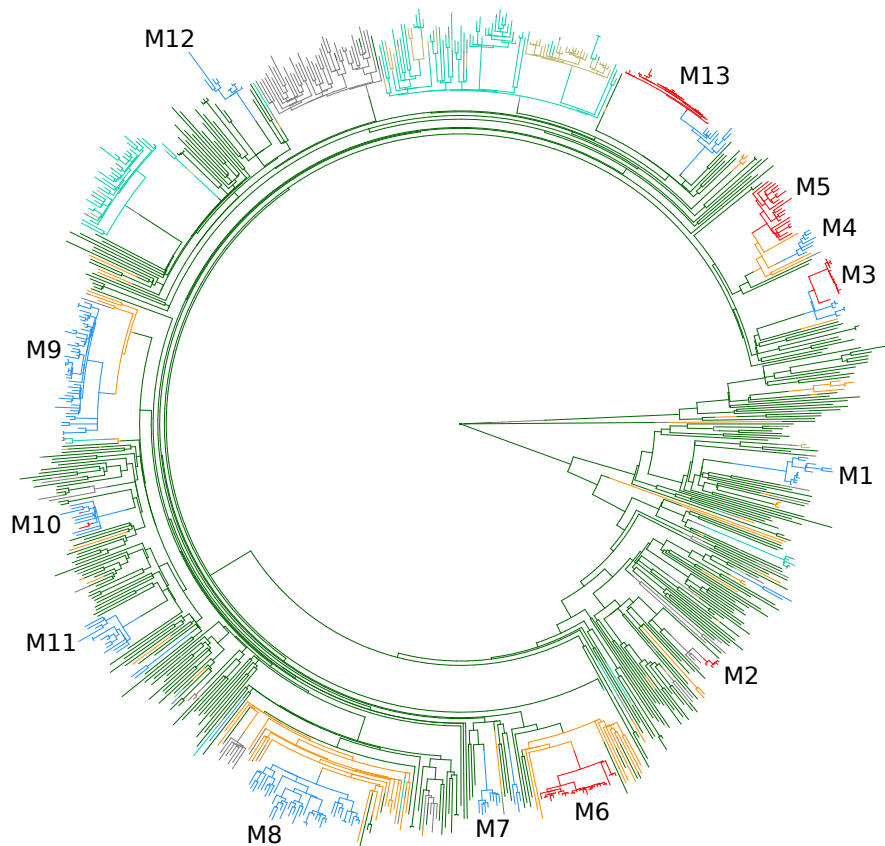
(B)



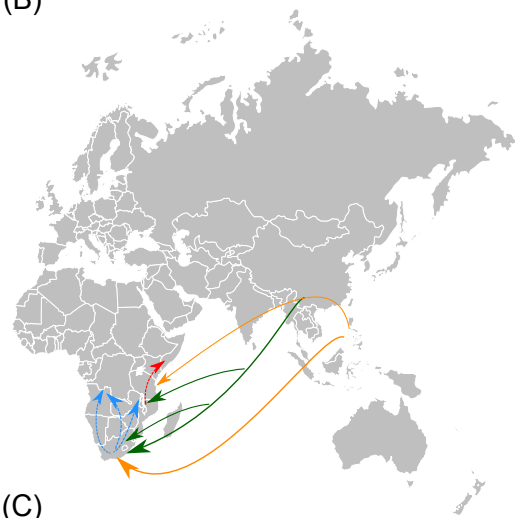




(A)



(B)



(C)

