

1 **Phylogenomics of *Mycobacterium africanum* reveals a new lineage and a complex**  
2 **evolutionary history**

3  
4 Mireia Coscolla<sup>1\*</sup>, Daniela Brites<sup>2,3</sup>, Fabrizio Menardo<sup>2,3</sup>, Chloe Loiseau<sup>2,3</sup>, Sonia  
5 Borrell<sup>2,3</sup>, Isaac Darko Otchere<sup>4</sup>, Adwoa Asante-Poku<sup>4</sup>, Prince Asare<sup>4</sup>, Leonor Sánchez-  
6 Busó<sup>5,6</sup>, Florian Gehre<sup>7,8</sup>, C. N'Dira Sanoussi<sup>9,10</sup>, Martin Antonio<sup>11</sup>, Affolabi Dissou<sup>12</sup>,  
7 Paula Ruiz-Rodriguez<sup>1</sup>, Janet Fyfe<sup>13</sup>, Erik C. Böttger<sup>14</sup>, Patrick Becket<sup>15-16</sup>, Stefan  
8 Niemann<sup>16</sup>, Abraham S. Alabi<sup>17</sup>, Martin P. Grobusch<sup>17,18,19</sup>, Robin Kobbe<sup>20</sup>, Julian  
9 Parkhill<sup>21</sup>, Christian Beisel<sup>22</sup>, Lukas Fenner<sup>23</sup>, Conor J. Meehan<sup>24</sup>, Simon R Harris<sup>25</sup>,  
10 Bouke C. De Jong<sup>8</sup>, Dorothy Yeboah-Manu<sup>4</sup>, Sebastien Gagneux<sup>2,3</sup>

11 \*corresponding author: mireia.coscolla@uv.es

12 1. I<sup>2</sup>SysBio, University of Valencia-FISABIO joint Unit, Valencia, Spain.  
13 2. Swiss Tropical and Public Health Institute, Basel, Switzerland  
14 3. University of Basel, Basel, Switzerland  
15 4. Noguchi Memorial Institute for Medical Research, University of Ghana, Legon, Accra, Ghana  
16 5. Centre for Genomic Pathogen Surveillance, Big Data Institute, Nuffield Department of  
17 Medicine, University of Oxford, Oxford, United Kingdom.  
18 6. Wellcome Sanger Institute, Wellcome Genome Campus, Cambridge, United Kingdom.  
19 7. Infectious Disease Epidemiology, Bernhard-Nocht-Institute for Tropical Medicine, Hamburg,  
20 Germany.  
21 8. Health Department, East African Community (EAC), Arusha, Tanzania  
22 9. Laboratoire de Référence des Mycobactéries, Cotonou, Bénin  
23 10. Unit of Mycobacteriology, Institute of Tropical Medicine, Antwerp, Belgium  
24 11. London School of Tropical medicine, London, UK  
25 12. Laboratoire de Référence des Mycobactéries, Ministry of Health, Cotonou, Benin  
26 13. Mycobacterium Reference Laboratory, Victoria Infectious Diseases Reference Laboratory,  
27 Peter Doherty Institute, Melbourne, Victoria, Australia  
28 14. Institute of Medical Microbiology, University of Zürich, Zürich, Switzerland  
29 15. Research Center Borstel, Molecular and Experimental Mycobacteriology, Borstel, Germany  
30 16. German Center for Infection Research, Partner Site Hamburg- Lübeck-Borstel-Riems,  
31 Borstel, Germany  
32 17. Centre de Recherches Médicales en Lambaréné (Cermel), Lambaréné, Gabon  
33 18. Institut für Tropenmedizin, Deutsches Zentrum fuer Infektionsforschung, University of  
34 Tübingen, Tübingen, Germany  
35 19. Center of Tropical Medicine and Travel Medicine, Department of Infectious Diseases,  
36 Amsterdam University Medical Centers, Amsterdam Infection & Immunity, Amsterdam Public  
37 Health, University of Amsterdam, Amsterdam, The Netherlands  
38 20. First Department of Medicine, Division of Infectious Diseases, University Medical Center  
39 Hamburg-Eppendorf; Germany  
40 21. Department of Veterinary Medicine, University of Cambridge, Madingley Road Cambridge,  
41 UK  
42 22. Department of Biosystems Science and Engineering, ETH Zürich, Basel, Switzerland  
43 23. Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland  
44 24. School of Chemistry and Biosciences, University of Bradford, Bradford, UK  
45 25. Microbotica Ltd, Biodata Innovation Centre, Wellcome Genome Campus, Hinxton,  
46 Cambridgeshire, UK

47 **Abstract**

48 Human tuberculosis is caused by members of the *Mycobacterium tuberculosis* Complex  
49 (MTBC). The MTBC comprises several human-adapted lineages known as *M.*  
50 *tuberculosis* sensu stricto as well as two lineages (L5 and L6) traditionally referred to as  
51 *M. africanum*. Strains of L5 and L6 are largely limited to West Africa for reasons  
52 unknown, and little is known on their genomic diversity, phylogeography and evolution.  
53 Here, we analyzed the genomes of 365 L5 and 326 L6 strains, plus five related genomes  
54 that had not been classified into any of the known MTBC lineages, isolated from patients  
55 from 21 African countries.  
56 Our population genomic and phylogeographical analyses show that the unclassified  
57 genomes belonged to a new group that we propose to name MTBC Lineage 9 (L9). While  
58 the most likely ancestral distribution of L9 was predicted to be East Africa, the most likely  
59 ancestral distribution for both L5 and L6 was the Eastern part of West Africa. Moreover,  
60 we found important differences between L5 and L6 strains with respect to their  
61 phylogeographical substructure, genetic diversity and association with drug resistance.  
62 In conclusion, our study sheds new light onto the genomic diversity and evolutionary  
63 history of *M. africanum*, and highlights the need to consider the particularities of each  
64 MTBC lineage for understanding the ecology and epidemiology of tuberculosis in Africa  
65 and globally.

66 **MAIN TEXT**

67

68 **Introduction**

69 Tuberculosis (TB) causes more human deaths than any other infectious disease, and it is  
70 among the top ten causes of death worldwide (1). Among the 30 high TB burden  
71 countries, half are in Sub-Saharan Africa (1). Africa also comprises the highest number  
72 of countries with the highest TB mortality (1). TB in humans and animals is caused by  
73 the *Mycobacterium tuberculosis* Complex (MTBC) (2), which includes different lineages,  
74 some referred to as *Mycobacterium tuberculosis* sensu stricto (Lineage 1 to Lineage 4 and  
75 Lineage 7) and others as *Mycobacterium africanum* (Lineage 5 and Lineage 6), a recently  
76 discovered Lineage 8 (3), as well as different animal-associated ecotypes such as *M.*  
77 *bovis*, *M. pinnipedii*, or *M. microti* among others (4, 5). Among the human-associated  
78 MTBC lineages, some are geographically widespread and others more restricted (6). The  
79 latter is particularly the case for Lineage (L) 7 that is limited to the Horn of Africa (7, 8),  
80 and L5 and L6 that are mainly found in West Africa (9). L5 and L6 differ substantially  
81 from the other lineages of the MTBC with respect to metabolism and in vitro growth (10,  
82 11). Several mutations in different genes of the electron transport chain and central carbon  
83 metabolic pathway can explain metabolic differences between L5 and L6 and the other  
84 lineages (12). L5 and L6 are also less virulent than other lineages in animal models, and  
85 appear to transmit less efficiently in clinical settings (13, 14). Even though L5 and L6 are  
86 mostly restricted to West-Africa, they show a prevalence of up to 50% among smear-  
87 positive TB cases in some West African countries (15-18). Hence, L5 and L6 contribute  
88 significantly to the overall burden of TB across sub-Saharan Africa. Compared to the  
89 other MTBC lineages, relatively little is known with regard to the ecology and evolution  
90 of L5 and L6 (5, 19). Two studies have found L5 to be associated with Ewe ethnicity in  
91 Ghana (20, 21), supporting the notion that this lineage might be locally adapted to this

92 particular human population (22). Several epidemiological associations suggest that L6  
93 might be attenuated for developing disease as compared to other lineages (reviewed in  
94 (9)). For example, L6 has been associated with slower progression from infection to  
95 disease in The Gambia (19). Other studies have linked L6 with HIV co-infection in TB  
96 patients from The Gambia and Ghana (19, 21), although other studies in Ghana and Mali  
97 have not seen such an association (23, 24). Human TB caused by *M. bovis* compared to  
98 *M. tuberculosis* has also been associated with HIV (25) and higher levels of  
99 immunosuppression as CD4 T cell counts  $\leq 200$  cells/ $\mu$ L (26), leading to the suggestion  
100 that L6 might be an opportunistic pathogen, similar to *M. bovis* in humans (27). L5 and  
101 L6 also differ in various molecular features relevant for patient diagnosis, such as a non-  
102 synonymous mutation in the MPT64 antigen (28) and reduced T cell response to ESAT6  
103 (29), leading to reduced detection by interferon gamma release assays of L5 and L6 in  
104 clinical samples (28, 30). To shed more light on the phylogeography, evolutionary history  
105 and population genetic characteristics of *M. africanum*, we analysed the largest set of  
106 whole genome data for L5 and L6 generated to date.

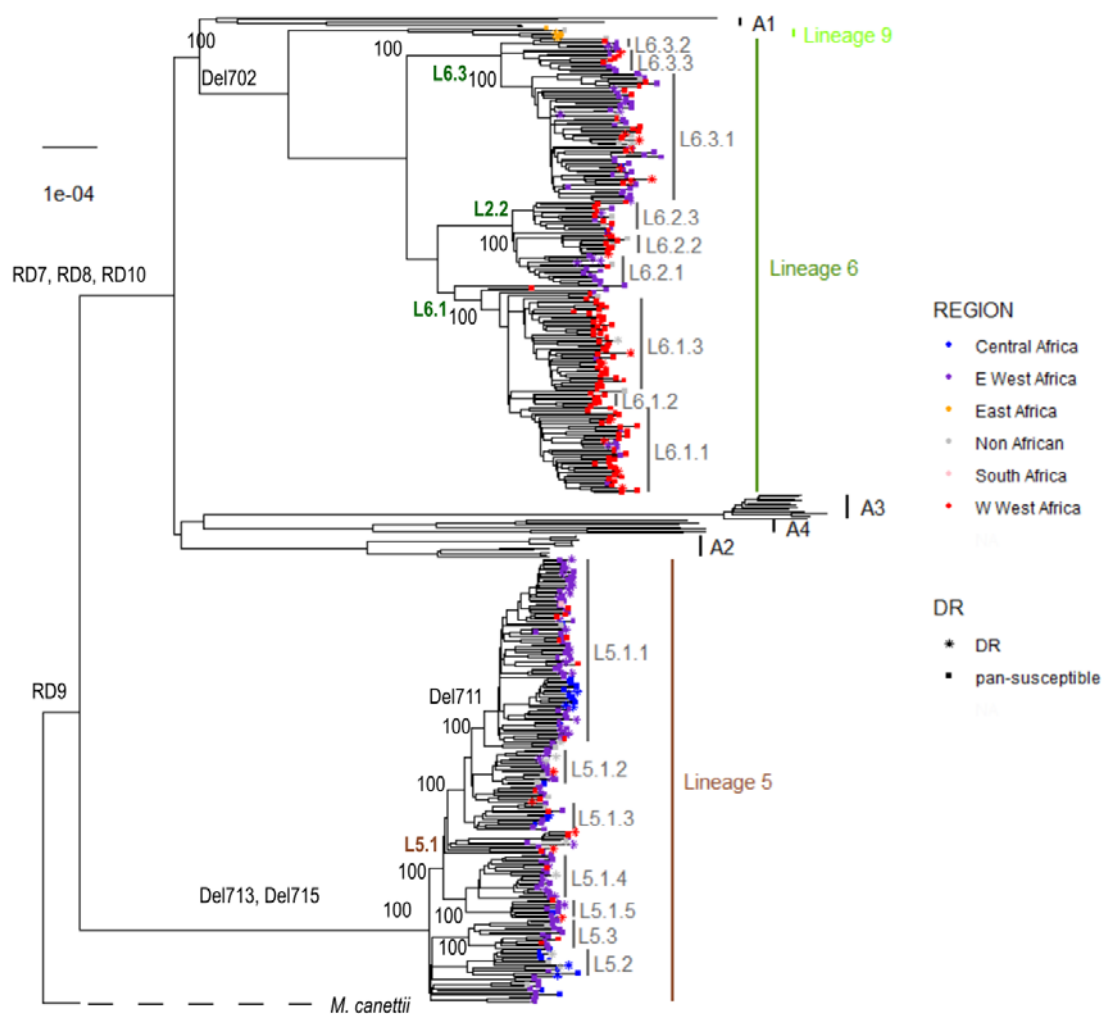
## 107 **Results**

### 108 **New MTBC Lineage: Lineage 9**

109 We analysed a total of 696 *M. africanum* genomes. These included 365 L5 and 326 L6  
110 genomes, as well as five related genomes that could not be classified into any of the  
111 known human- or animal-associated MTBC lineages (4, 31). Out of these 696 genomes,  
112 662 (95%) came from patient isolates originating in one of 21 countries of Sub-Saharan  
113 Africa. Another 34 (5%) strains were isolated outside Africa from patients with an origin  
114 other than Africa, or unknown (Table S1). To have a representative dataset and avoid  
115 overrepresentation of clustered strains, we removed 272 isolates that were redundant,  
116 while keeping the maximum phylogenetic diversity (>95% of the tree length) (32). The  
117 resulting non-redundant dataset comprised 424 genomes and showed a similar country  
118 distribution compared to the original dataset (Fig. S1).

119 We first focused our analysis on the five genomes that could not be classified into any of  
120 the known MTBC lineages. To explore the evolutionary relationship of these five  
121 genomes in the context of *M. africanum* diversity, we constructed the phylogeny of the  
122 424 *M. africanum* genomes plus a reference dataset of animal associated MTBC genomes  
123 we published previously (4). The resulting phylogeny (Fig. 1) corroborated the separation  
124 of L5 and L6, and the localization of L6 in a monophyletic clade together with the animal-  
125 associated lineages, as previously described (4).

126 **Fig. 1. Maximum likelihood phylogeny of 424 *M. africanum* genomes analysed**  
127 **together with reference animal associated genomes.** Support bootstrap values are  
128 indicated at the nodes. Nodes are coloured according to country or origin, and shape of  
129 the node indicates susceptible or drug resistance based on absence or presence at least  
130 one of the drug resistance mutations indicated in Table S8.



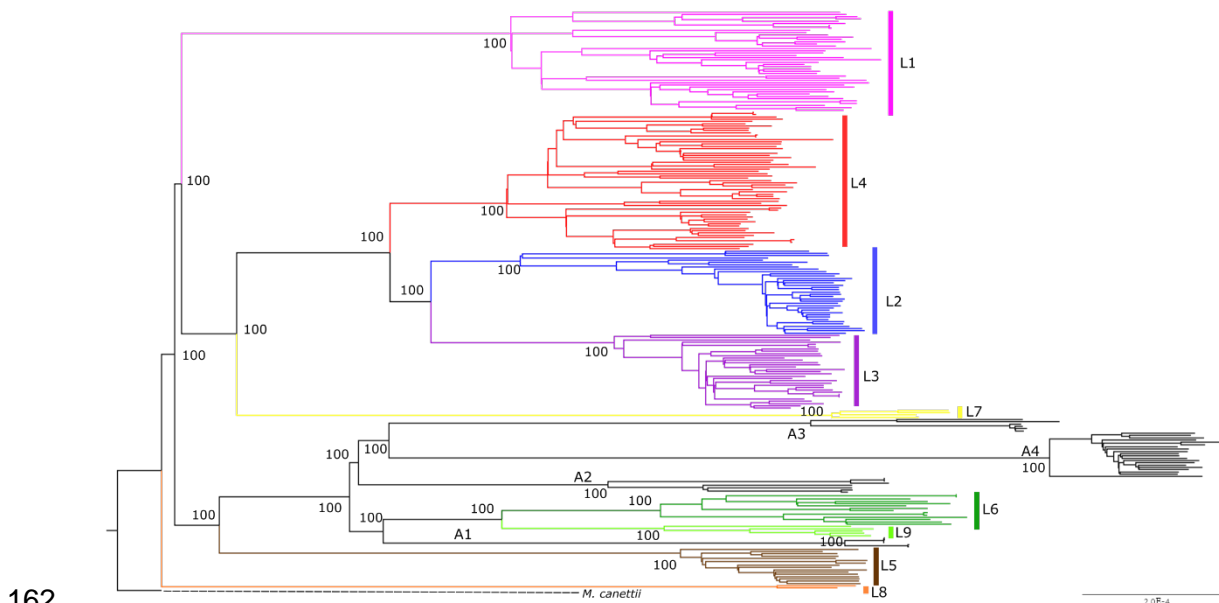
131

132 To further explore the phylogenetic position of these five genomes we constructed a  
133 genomic phylogeny with 248 reference genomes (3) including all eight human associated  
134 lineages and four animal associated lineages (Fig. 2). The five unclassified genomes  
135 appeared as a sister clade of L6, branching between L6 and the animal clade A1 (Fig. 2).  
136 The geographical origin of the five genomes differed from all other *M. africanum*  
137 genomes included in our analysis, as they were the only ones isolated from patients  
138 originating in East Africa (one from Djibouti, three from Somalia and one isolated in  
139 Europe but patient origin was unknown). By contrast, all L5 and L6 genomes came from  
140 patients originating in either West Africa (354 genomes) or Central Africa (37 genomes),

141 except for one isolated in South Africa (Fig. 1) and 28 isolated outside Africa and from  
142 unknown origin.

143 The five unclassified genomes showed the following in silico inferred spoligotype:  
144 772000007775671 (nnnnnnnonooooooooooooooooonnnnnnnnnnonnnnnnn) in the  
145 genome from Djibouti, 772700000003671  
146 (nnnnnnnonnnnonooooooooooooooooonnnnnnnnn) in all three Somalian genomes,  
147 and a very similar pattern 772600000003631  
148 (nnnnnnnonnnnonooooooooooooooooonnnnoonn) in the genome from Europe, for  
149 which the patient origin was unknown. We searched for these three spoligotypes in the  
150 international genotyping database SITVIT2, which includes 9,658 different spoligotypes  
151 from 103,856 strains isolated in 131 countries (33). Spoligotype 772600000003631 was  
152 not found among the 103,856 strains included in the database, and the other two  
153 spoligotypes can be considered extremely rare because they have been found only in three  
154 strains in the database: 772000007775671 in a strain isolated in France, and  
155 772700000003671 in two strains isolated in the Netherlands, although patient's origin is  
156 unknown.

157 **Fig. 2. Maximum likelihood phylogeny of five unclassified genomes analysed**  
158 **together with reference dataset of MTBC genomes.** The five unclassified genomes are  
159 coloured in light green and tagged as “L9”. Animal associated lineages A1 to A4 are  
160 indicated and coloured in black. Support bootstrap values are indicated at the deepest  
161 nodes.



162

163 The five unclassified genomes showed a mean distance of 1,191 SNPs to L6 genomes,  
164 1,632 SNPs to L5 genomes, and 1,491 SNPs to the animal-associated MTBC genomes.  
165 Those distances were higher than the corresponding intra-lineage differences: 342  
166 (standard deviation (SD) 3.65) within L5, 542 (SD 9.19) within L6, and 332.4 (SD 14.48)  
167 within the unclassified genomes. So, even when correcting for the diversity within each  
168 lineage, we still found that the five unclassified genomes were separated from the other  
169 lineages by 1,294, 582 and 654 SNPs of net distance to L5, L6 and the animal-associated  
170 lineages, respectively. Given the different geographical distribution and the substantial  
171 genetic separation, we classified these five genomes into a new MTBC lineage that we  
172 propose to call MTBC Lineage 9 (L9). The average intra-lineage diversity among these  
173 five L9 strains was 332 SNPs (SD=13). The maximum diversity within L9 was 514 SNPs  
174 between strain G00075 and strain G00074, with the smallest distance being 99 SNPs  
175 between strain G04304 and strain G00075.  
176 We looked for deleted regions in the L9 genomes that could be used as phylogenetic  
177 markers, as was done for other MTBC lineages in the past (34, 35) (6). We identified one  
178 region deleted in all L9 genomes that spanned from Rv1762c to Rv1765. However, this



179 region is not a robust phylogenetic marker because partially overlapping deletions can be  
180 found in other lineages. Specifically, Rv1762c is deleted in genomes from one of the  
181 animal associated lineages, Lineage A3, which includes the strain previously known as  
182 *M. orygis*, and the region between Rv1763c and Rv1765 is deleted in L6 genomes. Hence  
183 instead, we report a list of SNPs that can be used as phylogenetic markers for L9 (Table  
184 S2) given that they appear in all five L9 genomes and are absent from genomes from other  
185 lineages (32).

186 Given the low number of L9 genomes, we focused the remaining of our analysis on *M.*  
187 *africanum* L5 and L6.

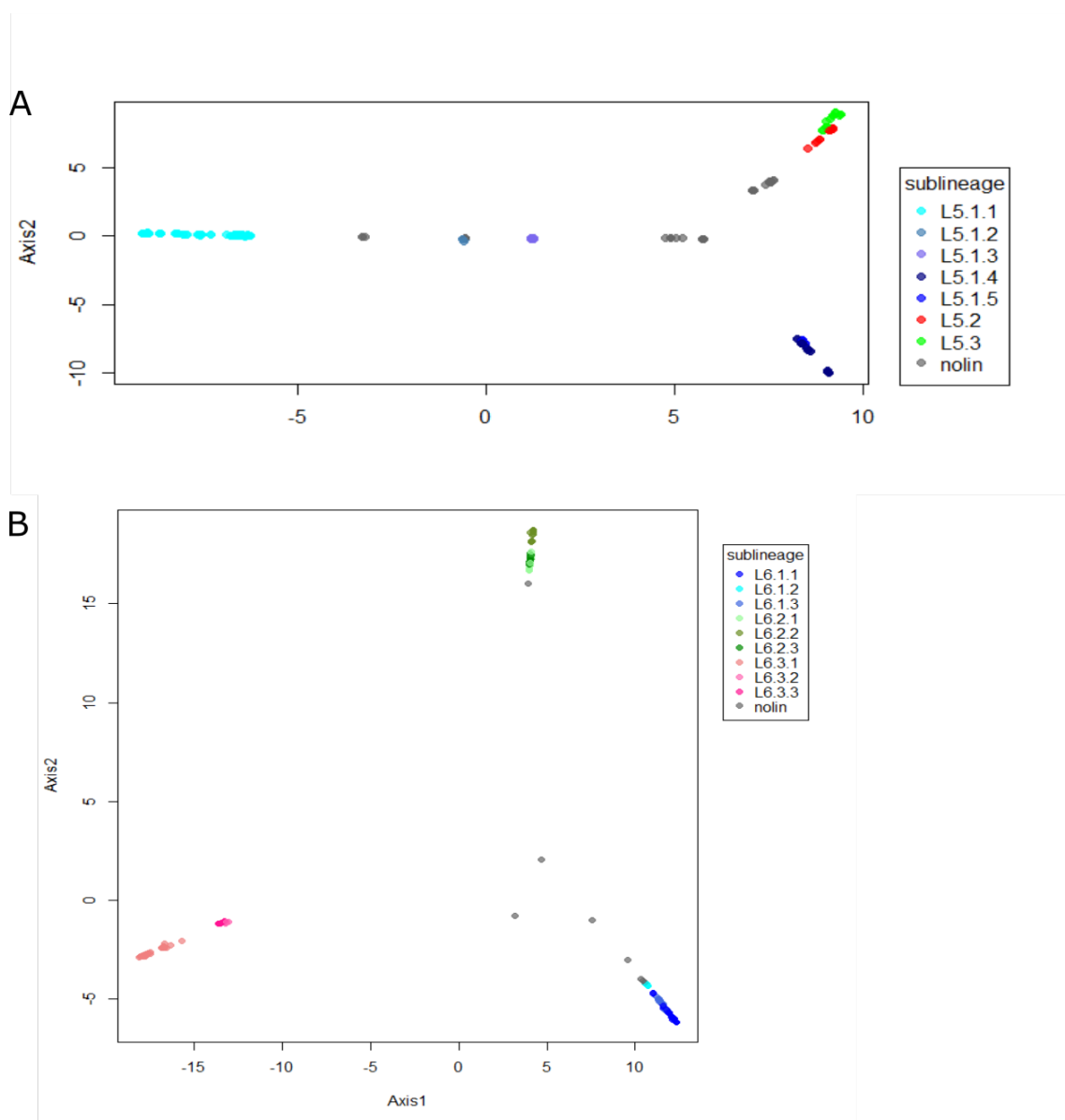
188

### 189 **Sublineages within L5 and L6**

190 Our extended genomic analysis of L5 and L6 confirmed the deletions of the previously  
191 described regions of difference (RDs), including RD7, RD8, RD9 and RD10 (34, 35), and  
192 RD713 and RD715 (6) as indicated in the phylogeny (Fig. 1). However, the deletion of  
193 RD711 could not be confirmed as a L5 marker as proposed previously (6), as it was only  
194 deleted in a subset of L5 genomes as reported recently (36). We found RD711-deleted  
195 genomes to form a monophyletic clade within L5; named L5.1.1 considering previous  
196 nomenclature as proposed in Ates et al. (36). In contrast, RD702 was found to be deleted  
197 in all L6 strains as shown previously (6), as well as in the newly defined L9 strains (Fig.  
198 1).

199 **Fig. 3. Principal Component Analysis (PCA) based on genomic variable SNPs.** The  
200 PCA was conducted separately for L5 (A) and L6 (B). Colours indicate different  
201 sublineages and grey indicates genomes with no sublineage assigned “nolin”.

202



203

204 Our phylogeny revealed a different topology for L5 as compared to L6. Specifically, the  
205 L5 phylogeny showed little structure. Nevertheless, we managed to subdivide L5 into  
206 three main sublineages that were well differentiated and highly supported by bootstrap  
207 values >90, and named them consistent with previous nomenclature (36) as L5.1, L5.2  
208 and L5.3. Due to the high genomic diversity within L5.1, this group was further  
209 subdivided into five subgroups (Fig. 1), leading to a total of seven sublineages in this first  
210 and second level of subdivision. Sublineage classification was only partially corroborated  
211 by the results of the PCA performed on whole genome SNPs (Fig. 3A), where these  
212 sublineages were not clearly separated. By contrast, L6 showed a more differentiated tree  
10

213 structure with three clearly differentiated monophyletic sublineages (L6.1, L6.2 and L6.3)  
214 at the first level that could be further subdivided into a second level subdivision with at  
215 least three other subgroups each (Fig. 1), resulting in a total of nine sublineages. The first  
216 level of subdivision was strong for L6, where L6.1, L6.2 and L6.3 were clearly separated  
217 using PCA (Fig. 3B). However, sublineages at the second level of subdivision were not  
218 that clearly separated (Fig. 3B). To explore the robustness of the classification beyond  
219 PCA, we estimated genetic differentiation for each of these sublineages using the fixation  
220 index (FST) based on Wright's F-statistic (37) as measure of population differentiation  
221 due to genetic structure. We conducted a hierarchical analysis comparing the population  
222 structure at the two levels of subdivision: one level with the three main groups for both  
223 L5 and L6, and a second level with all seven and nine sublineages of L5 and L6,  
224 respectively. The L5 population structure showed the highest differentiation within all  
225 seven sublineages, where the highest population differentiation index  $F_{st}=0.48$  ( $p$ -  
226  $value<0.000001$ ), and the lowest population differentiation index was found between the  
227 three main sublineages at the first level of subdivision ( $F_{st}=0.14$ ,  $p$ - $value=0.04915$ ).  
228 Similarly,  $F_{st}$  between all seven L5 sublineages showed moderate differentiation with  
229 pairwise  $F_{st}$  values between 0.3 and 0.5 (Table S3) and net pairwise differences between  
230 76 and 206 SNPs (Table S4). Conversely, for L6, the higher differentiation was at the  
231 first level of subdivision, that is between the main sublineages (L6.1, L6.2, L6.3, with  
232 47% of the variation,  $F_{st}=0.47$ ,  $p=0.0035$ ), mirroring the PCA results. Even  
233 differentiation at the second level of subdivision, that is between all nine sublineages of  
234 L6, showed more structure than for L5, with  $F_{st}$  values ranging between 0.25 and 0.75  
235 (Table S5), and net pairwise differences of between 73 and 493 SNPs (Table S6). A list  
236 of SNPs found exclusively in each of the L5 and L6 sublineages is shown in Table S7.  
237

## 238 **Phylogeography**

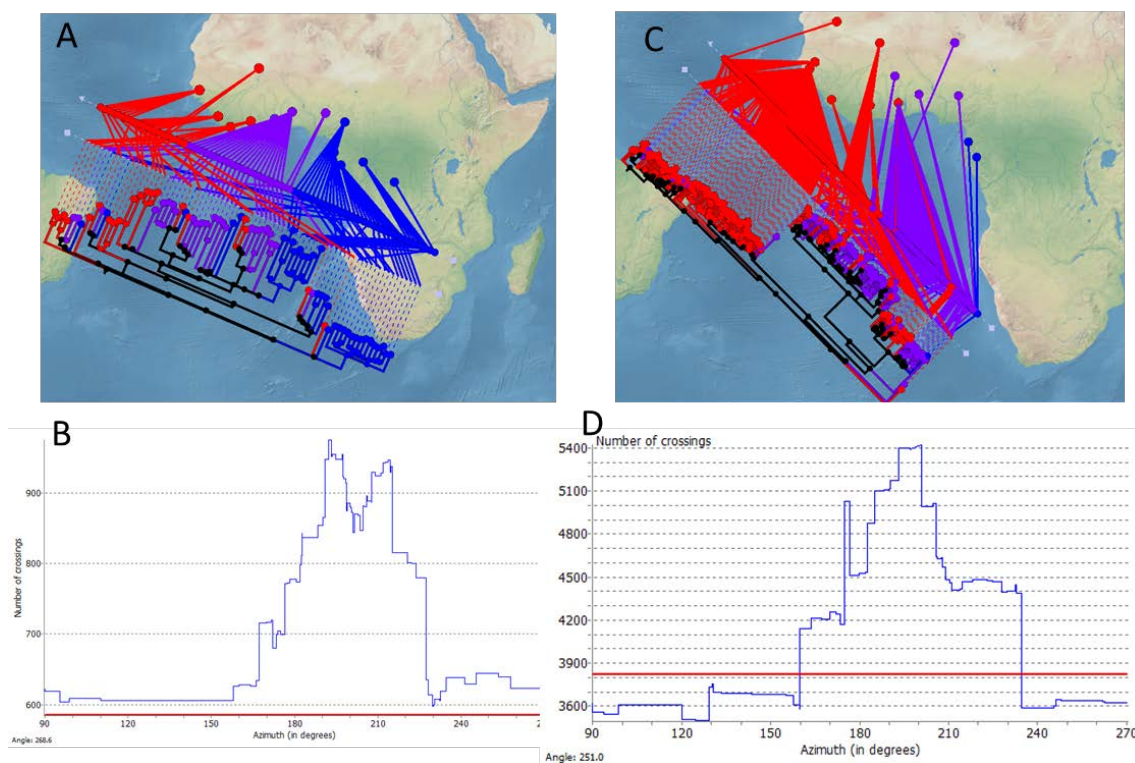
239 To explore the phylogeographical structure of L5, L6, and L9, we mapped the  
240 geographical origin of the genomes onto the phylogenetic tree as a coloured point at the  
241 end of each branch (Fig. 1). We grouped the different countries represented in the dataset  
242 into five regions in Africa: East, South, Central, and the Western part of West Africa  
243 (<sup>W</sup>West Africa) and the Eastern part of West Africa (<sup>E</sup>West Africa). We observed that  
244 most sublineages showed a characteristic geographical association at the regional level.  
245 At least five sublineages within L6 (all three L6.1 and two L6.2) showed a majority of  
246 genomes originating in <sup>W</sup>West Africa, mostly The Gambia. By contrast, a few scattered  
247 L6 genomes, one sublineages within L6.2 and all three L6.3 genomes came from <sup>E</sup>West  
248 Africa, mostly Ghana. Only a few L6 strains were found in Central Africa (N=2) or  
249 outside Africa (N=15). L5 showed a different phylogeographical structure with most  
250 genomes originating in <sup>E</sup>West Africa (mostly Ghana) and two groups (L5.2 and one  
251 sublineage within L5.1.1) in Central Africa. Only a few dispersed genomes originated  
252 from <sup>W</sup>West-Africa.

253 To verify the geographic separation within L5 and L6, we conducted an independent  
254 phylogeographic analysis using the GenGIS software, where each whole genome SNP  
255 phylogeny was superimposed onto the five main African regions defined previously (Fig.  
256 4A and C). We found several orientations of the tree's geographical axis resulting in less  
257 crossings than expected by chance in L6 ( $p < 0.001$ , 10,000 permutations; Fig. 4D). By  
258 contrast, for L5 we did not find any lineage axis with less crossing than expected by  
259 chance (Fig. 4B). These results indicate a marked geographical structure within L6 but  
260 not within L5. To further verify the different phylogeographical structures within L5 and  
261 L6, we calculated population differentiation indices considering each African region as a  
262 different population for each lineage. This analysis revealed some phylogeographical

263 substructure within L6, where the percentage of variation attributed to different regions  
264 within Africa was 15% ( $F_{st}=0.15$ ,  $p<0.00001$ ). By contrast, L5 did not show any well-  
265 marked population differentiation, as the percentage of the variance attributed to  
266 population differences was only 6.6 %, with the rest of the variation attributed to intra-  
267 population differences ( $F_{st}=0.036$ ,  $p<0.00001$ ). This result further supports the  
268 observation of higher geographical structure within L6 than L5.

269

270 **Fig. 4. Phylogeographical structure in L5 and L6.** Linear axis plot between the  
271 genomic phylogeny and the geographical origin of the genomes for L5 (A) and L6 (C).  
272 Histograms show the number of crossing for each inclination of the axis, and the red line  
273 indicates the number of crossing expected by chance for L5 (B) and L6 (D).



274

275

276 Finally, we explored possible differences in geographic range. Our dataset was  
277 geographically biased because it was designed to assemble as many L5 and L6 genomes  
13

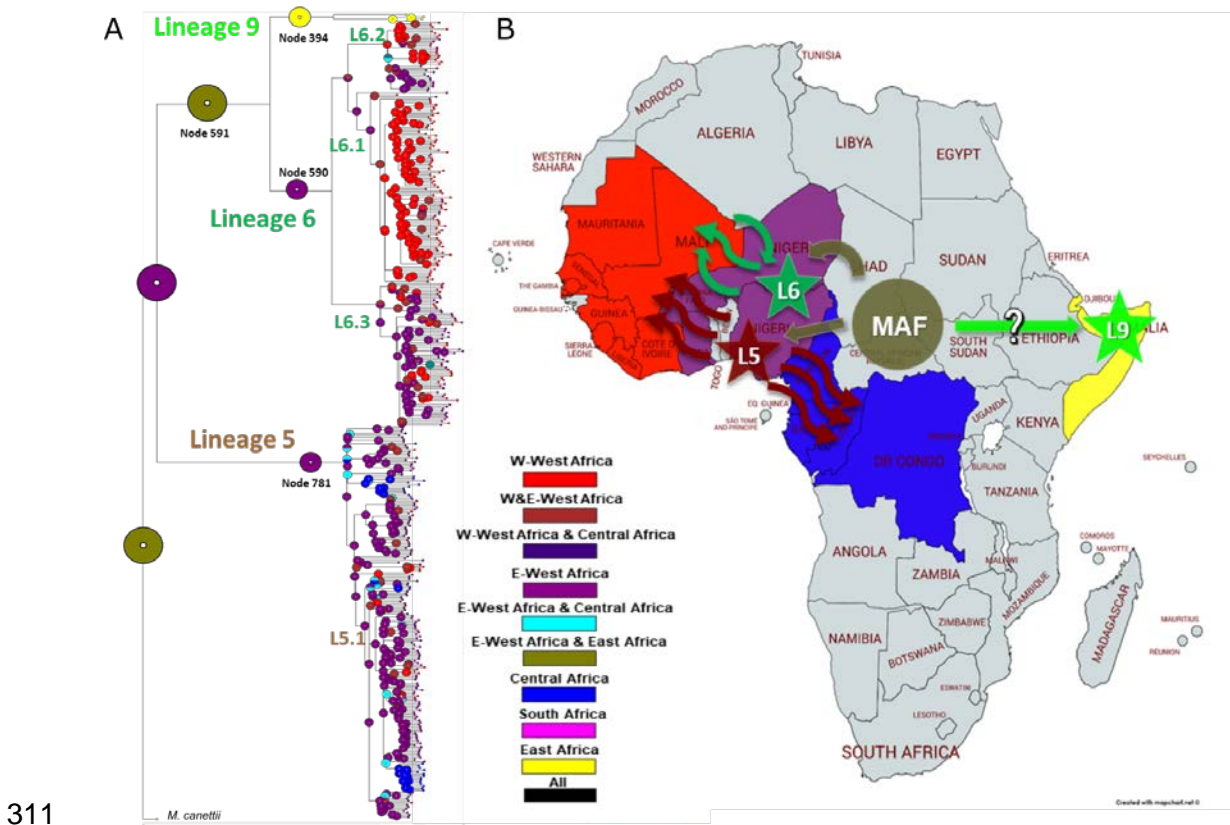
278 from as many countries as possible. We therefore analysed our genome dataset together  
279 with two other large datasets where samples were not genome sequenced but genotyped  
280 using spoligotyping to compare the geographic distributions of L5 and L6 (33, 38). This  
281 combined dataset included N=733 L5 from 27 African countries and N=1,031 L6 from  
282 18 African countries. We expected that a broader geographical distribution of a specific  
283 lineage associated with a lower probability that two individuals selected randomly will  
284 belong to the same country. We used the Simpson's Index (D) to measure the probability  
285 that two individuals randomly selected from a sample will belong to the same country.  
286 We found a larger diversity of countries of origin in L5 than in L6 (D=2.29 vs D=1.78, t-  
287 test  $\alpha < 0.05$ ) indicating a broader geographic distribution of L5.

288

### 289 **The ancestral geographical distribution of L5, L6 and L9**

290 Next, we explored the most likely geographical origin of L5 and L6 using four methods  
291 based on a Bayesian approach (39). The probabilities of ancestral distribution areas for  
292 the principal nodes were always congruent with at least two methods, but the results of  
293 the two other methods were either inconclusive or showed minor discrepancies (Fig. 5A  
294 and Fig. S2). For L5, two of the four methods inferred <sup>E</sup>West Africa as the most likely  
295 origin (marginal probability was 1.0 using both Bayesian binary and S-DIVA), while the  
296 other two were inconclusive (marginal probabilities were <sup>E</sup>West - Central: 0.94 and 0.58  
297 with Bay Area and DEC, respectively; node 783 in Fig. 5A and Fig. S2). For L6, two  
298 methods also pointed to <sup>E</sup>West Africa as the most likely origin (0.77, 1.0, of marginal  
299 probability using Bayesian binary and S-DIVA, respectively) and two methods supported  
300 both regions of West Africa as equally likely (0.94 and 0.58 using Bay area and DEC,  
301 respectively; node 592 in Fig. 5A and Fig. S2). The ancestral distribution of L9 was  
302 predicted to be East Africa based on all four methods (node 396 in Fig. 5A and Fig. S2).

303 The ancestral distribution of the common ancestor between L6 and L9 was not that clearly  
304 predicted because of marginal probabilities of the methods supporting <sup>E</sup>West Africa (0.65,  
305 0.57 using BMBM, and DEC; node 591 Fig. 5A and Fig. S2A) and two methods  
306 supporting both regions in West Africa (0.5 using S-DIVA and Bay Area).  
307 By contrast, the ancestral distribution for L5, L6 and L9 showed more consistency, where  
308 <sup>E</sup>West Africa was supported by three methods (0.74, 1.0 and 0.57 using S-DIVA, BMBM  
309 and DEC, respectively) and one method predicting both <sup>E</sup>West Africa and East Africa  
310 with a marginal probability of 0.99 (Bay Area: node 784 Fig. 5A and Fig. S4).



312 **Fig. 5. Geographical ancestral distributions of L5, L6 and L9.** A. Ancestral area  
313 reconstruction by the Bayesian binary model onto the maximum likelihood phylogeny.  
314 Circles represent the probabilities of ancestral ranges, and the most likely ancestral areas  
315 are indicated by their corresponding colour code. B. Four geographical areas considered  
316 in this analysis are coloured in the map, the most likely areas ancestral areas for each

317 lineage shown as stars, and movements of strains inferred from phylogeny indicated as  
318 arrows.

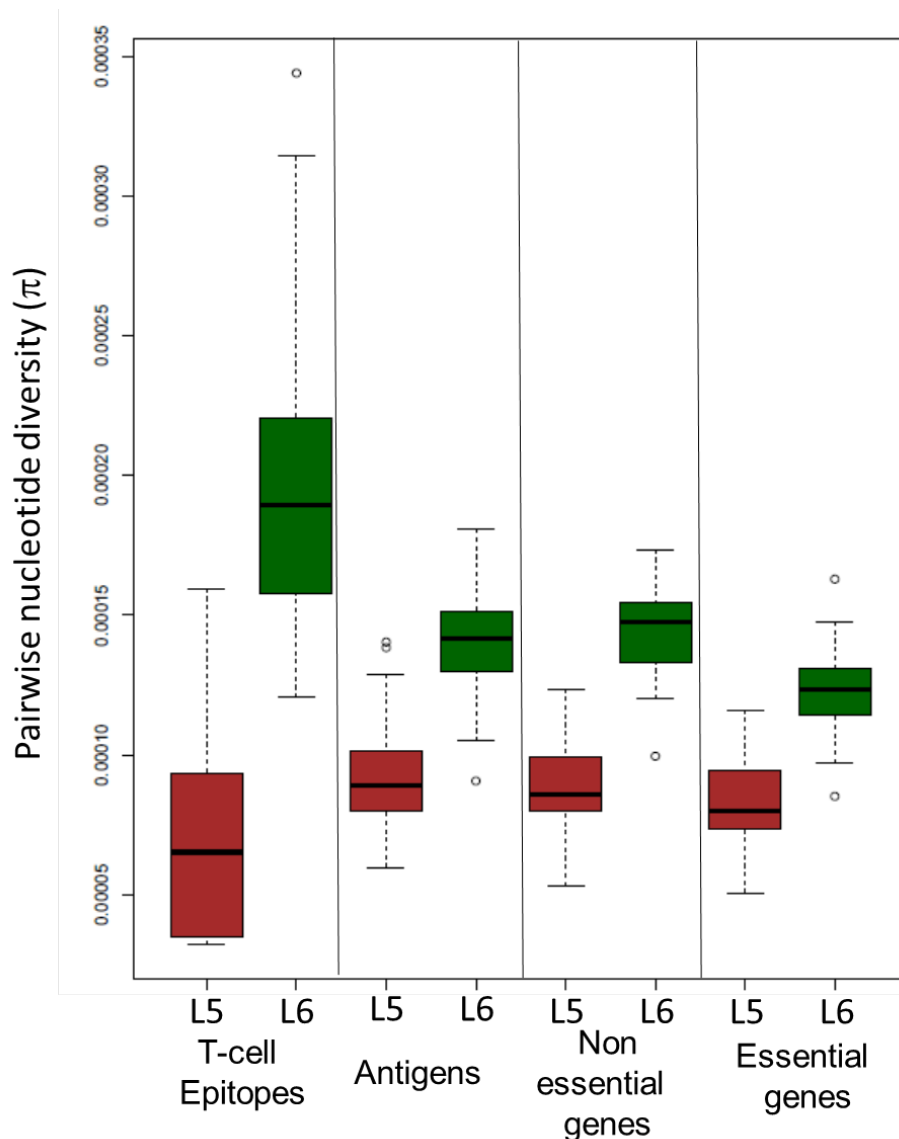
319

### 320 **Differences in genetic diversity between lineages**

321 In support of our previous findings based on a more limited dataset (40), we found that  
322 L6 was significantly more genetically diverse than L5 with significantly higher number  
323 of SNPs between pairs of sequences (median values 553 vs 321; p-value < 2.2e-15), and  
324 significantly higher average nucleotide diversity ( $1.4 \times 10^{-4}$  vs  $8.7 \times 10^{-5}$ ; p-value < 2.2e-  
325 15). To explore if this trend was consistent across the whole genome, we studied the  
326 nucleotide diversity in different regions that might be under different selection pressures:  
327 essential genes, non-essential genes, antigens, and T cell epitopes (Fig. 6). Although the  
328 genetic diversity was higher in all these different gene categories for L6 (Fig. 6), epitopes  
329 showed an inverted pattern in diversity between lineages (Fig. 6). Specifically, epitopes  
330 in L6 showed significantly higher genetic diversity than non-essential genes (Wilcoxon  
331 signed rank test p-value < 2.2e-15), while the opposite was found for L5, with epitopes  
332 showing significantly lower genetic diversity than non-essential genes (Wilcoxon signed  
333 rank test p-value < 2.2e-15).

334 **Fig. 6. Nucleotide diversity ( $\pi$ ).** Comparison of pairwise nucleotide diversity ( $\pi$ )  
335 between L5 and L6 across gene categories.





336

### 337 **Drug resistance mutations**

338 Antibiotic pressure is a strong selective force in bacteria including MTBC. Hence, we  
339 explored the difference in drug resistance determinants between L5 and L6. We found  
340 that among the 424 genomes analysed, 89 (21%) showed at least one genetic marker of  
341 antimycobacterial drug resistance, with 24 (6%) being multi-drug resistant (which is  
342 resistance to at least isoniazid and rifampicin, Table S8). The most common resistance  
343 found was for streptomycin, with 60 genomes showing 13 different resistance-conferring  
344 mutations. The next most common was resistance to rifampicin and isoniazid, with 32

345 and 29 genomes, respectively. Additional resistance was found to ethambutol,  
346 fluoroquinolones, ethionamide, pyrazinamide and aminoglycosides (Table S8). L5  
347 genomes were more likely than L6 genomes to carry mutations associated with any  
348 resistance (OR 2.05 [95% confidence interval (CI) 1.26-3.31],  $p$ -value= $2.29 \times 10^{-3}$  using  
349 Fisher's Exact Test). However, this was not due to a single antibiotic resistance profile  
350 because both lineages did not differ significantly when comparing the number of drug  
351 resistance mutations to fluoroquinolones ( $p$ -value=0.32), ethambutol ( $p$ -value=0.32),  
352 isoniazid ( $p$ -value=0.32), rifampicin ( $p$ -value=0.2), or streptomycin ( $p$ -value=0.34).  
353 Contrary to a previous report by Ates et al. (36), we found no evidence of differences in  
354 drug resistance genotype between L5.2 and other L5 genomes (OR 1.21 [95% CI 0.36-  
355 4.11],  $p$ -value=0.49, Fisher's Exact Test).

356

## 357 **Discussion**

358 *M. africanum* has traditionally been considered a single entity and a separate species from  
359 what classically has been referred to as *M. tuberculosis* sensu stricto. The results  
360 presented here provide novel insights into the genomic particularities of the different  
361 lineages within *M. africanum*: L5, L6 and a new group described in this study, L9.  
362 Differences between these three lineages further emphasize the need to consider these  
363 lineages as separate phylogenetic and ecologic variants within the MTBC.

364 Unexpectedly, our study of the global diversity of *M. africanum* revealed the presence of  
365 another MTBC lineage in Africa: L9, which is genetically close to L6. But unlike L5 and  
366 L6 that predominately occur in West Africa, L9 seems to be restricted to the East of  
367 Africa. Given that only five L9 isolates were included in our study, future studies are  
368 needed to confirm this observation (7, 8). In this respect, L9 is similar to L7 and the

369 recently described L8 (3), which are also mainly restricted to East Africa, but genetically  
370 more distant. We found clinical strains of L9 to be rare compared to L5 and L6, and this  
371 observation also resembles the situation for L7 and L8. As mentioned, we cannot dismiss  
372 the notion that this might be due to limited sampling, but the observation that clinical  
373 strains from L7, L8 and L9 originate in East Africa and are generally rare, while L5 and  
374 L6 are more prevalent and distributed across West and Central, raises the question of  
375 whether the reduced prevalence of L7, L8 and L9 is due to biological reasons, or social-  
376 environmental causes that renders L7, L8 and L9 to be less successful. The lack of  
377 experimental and epidemiological data on L7, L8 and L9 impedes a profound discussion  
378 on the matter. However, the fact that L9 is genetically closer to L6 and L5 than to L7 and  
379 L8, speaks against a common intrinsic biological determinant shared by L7, L8 and L9.  
380 Instead, convergence in the biology of the strains and/or in the socio-demography of the  
381 host is a more likely driver of the evolutionary history of L7, L8 and L9.

382 Our phylogeographic analyses localized the common ancestors of L5 and L6 to <sup>E</sup>West  
383 Africa. We could detect that several subgroups of L5 moved from <sup>E</sup>West Africa to Central  
384 Africa, while L6 subgroups moved mostly within West Africa. One of these events  
385 resulted in half of the L6 genomes in our dataset moving from <sup>E</sup>West Africa to <sup>W</sup>West  
386 Africa and with few dispersals back to <sup>E</sup>West Africa (Fig. 5B). The ancestral  
387 reconstruction of L6 and L9 did not provide any clear signal, with <sup>E</sup>West Africa and East  
388 Africa equally supported. For the ancestral distribution of all *M. africanum*, there was no  
389 consensus, but three out of four methods agreed on <sup>E</sup>West Africa being the most likely  
390 place of origin. That would imply that L5 and L6 diversified there, and L9 migrated to  
391 East Africa. Remarkably, the fact that L9 is only present in East Africa, similar to L7 and  
392 L8 (3, 7, 8), suggests either one or two migration events to East Africa, depending on the  
393 ancestral distribution of the MTBC as a whole (Fig. 5B). Because *M. canettii*, the most

394 closely related species of *M. tuberculosis* is highly restricted to East Africa, we and others  
395 have proposed that East Africa is the likely origin of the MTBC (41-43). If confirmed,  
396 the current geographical distribution of L5, L6 and L9 could be explained by a migration  
397 of their common ancestor from East Africa to West Africa, with the ancestor of L9 then  
398 moving back to East Africa. Alternatively, if the origin of the MTBC was in Central or  
399 West Africa, the current distribution would reflect at least three migration events to East  
400 Africa: one for the ancestor of L8, one for the ancestor of L7 and one for the ancestor of  
401 L9.

402 The work presented here also demonstrates differences in the population structure of L5  
403 compared to L6. While L6 showed a marked phylogenetic structure comprising distinct  
404 sublineages associated with different geographical regions, the classification of L5 into  
405 sublineages was not so clearly supported despite the fact that L5 showed broader  
406 geographical range compared to L6. Additionally, our work confirms previous  
407 observations of differences in the genomic diversity, where L6 shows a higher diversity  
408 compared to L5 (40). In particular, human T cell epitopes in L6 were more diverse than  
409 non-essential genes, while the opposite was true for L5. Several studies have shown that  
410 human T cell epitopes in the human-adapted MTBC are overall more conserved than non-  
411 essential genes (44-46). This observation gave rise to the hypothesis that the MTBC might  
412 benefit from T cell recognition that drives lung pathology, leading to enhanced bacterial  
413 transmission (47). The fact that L6 differs in this respect from L5 and the other human-  
414 adapted MTBC lineages, indicates a potential different ecological niche, including  
415 possible animal reservoirs (17), which would also be supported by the phylogenetic  
416 proximity of L6 to the animal-adapted lineages of the MTBC (Figure 1).

417 We found L5 genomes more likely to carry any drug resistance-conferring mutations than  
418 L6. This result was consistent with previous findings from Ghana where L5 was  
20

419 compared to L4 (48). Due to the dominance of L5 genomes from Ghana in our dataset,  
420 we cannot rule out that our observation might have been partially driven by the Ghanaean  
421 genomes. However, unlike the previous report from Ghana, our study found L5 to be  
422 associated with any resistance, as opposed to specifically with a single antibiotic. In  
423 addition, contrary to the previous study from Ates et al. (36) based on a smaller dataset,  
424 our larger sampling indicated no association between drug resistance and a specific  
425 sublineage of L5 (36).

426 Our main study limitation is sampling bias, leading to an overrepresentation of isolates  
427 from the Gambia and Ghana. The overrepresentation of genomes from these two countries  
428 could contort our observation regarding the genomic diversity and population structure.  
429 Moreover, including more genomes from other countries will likely reveal additional sub-  
430 lineages within L5 and L6. Importantly, as stated in the previous paragraph, the  
431 association between L5 and drug resistance can be partially driven by a similar situation  
432 reported in Ghana previously. However, the differences we found between L5 and L6 is  
433 unlikely driven by this overrepresentation, because each country was enriched with  
434 strains from one of the two lineages.

435 In summary, we describe a large-scale whole-genome sequencing and a comprehensive  
436 phylogenomic analysis of clinical isolates classically referred to as “*M. africanum*” from  
437 21 countries across Africa. Our findings have unravelled hidden diversity, a complex  
438 evolutionary history, and differential patterns of variation between lineages. Our results  
439 contribute to a better understanding of the MTBC lineages restricted to parts of Africa.  
440 These findings might assist in unraveling the molecular signatures of adaptations, and  
441 inform the development of targeted interventions for controlling TB in that part of the  
442 world.

## 443 **Methods**

### 444 ***M. africanum* dataset**

445 We analysed 697 L5 and L6 genomes to determine the genetic diversity, phylogeography  
446 and population structure of *M. africanum* (Table S1). This dataset included 495 newly  
447 sequenced genomes and 88 genomes from a previous study (4). Geographical origin was  
448 determined as the country of origin of the patient, and when not available the country of  
449 isolation. Because the number of different countries was too high to be shown properly  
450 in the figures, and some of them only included very few genomes, we grouped countries  
451 together into five African regions following definitions in (38): three big regions such as  
452 South, East and Central Africa, and two regions within West Africa, where most of the  
453 isolates come from. Western part of West Africa includes Gambia, Senegal, Mauritania,  
454 Sierra Leone, Liberia, Guinea, Ivory Coast, and Mali while the Eastern part of West  
455 Africa includes Ghana, Nigeria, Benin Niger, Burkina Faso). African maps were built  
456 using Mapchart® (<https://mapchart.net/africa.html>)

### 457 **Bacterial Culture, DNA extraction and Whole-Genome sequencing**

458 Archived MTBC isolates were revived by sub-culturing on Lowenstein Jensen media  
459 slants supplemented with 0.4% sodium pyruvate or with 0.3% glycerol to enhance the  
460 growth of the different lineages and incubated at 37 °C. Five loops full of colonies were  
461 harvested at the late exponential phase into 2 mL cryo-vials containing 1 mL of sterile  
462 nuclease-free water, inactivated at 98 °C for 60 minutes for DNA extraction using the  
463 previously described hybrid DNA extraction method (48). The MTBC lineages were then  
464 confirmed by spoligotyping and long sequence polymorphisms and sent for whole  
465 genome sequencing.

466 The MTBC isolates were grown in 7H9-Tween 0.05% medium (BD) +/- 40mM sodium  
467 pyruvate. We extracted genomic DNA after harvesting the bacterial cultures in the late  
468 exponential phase of growth using the CTAB method (49).

469 Sequencing libraries were prepared using NEXTERA XT DNA Preparation Kit  
470 (Illumina, San Diego, USA). Multiplexed libraries were paired-end sequenced on  
471 Illumina HiSeq2500 (Illumina, San Diego, USA) with 151 or 101 cycles when sequenced  
472 at the Genomics Facility Basel, HiSeq 2500 (100 bp, paired end) when sequenced at the  
473 Wellcome Sanger Institute, or on Illumina MiSeq (250 and 300 bp, paired end) or  
474 NextSeq (150 bp, paired end) according to the manufacturer's instruction (Illumina, San  
475 Diego, USA) when sequenced at the genomics facilities in Research Center Borstel.

#### 476 **Bioinformatics analysis:**

##### 477 **Mapping and variant calling of Illumina reads.**

478 The obtained FASTQ files were processed with Trimmomatic v 0.33  
479 (SLIDINGWINDOW: 5:20) (50) to clip Illumina adaptors and trim low quality reads.  
480 Any reads shorter than 20 bp were excluded for the downstream analysis. Overlapping  
481 paired-end reads were merged with SeqPrep v 1.2 (overlap size = 15)  
482 (<https://github.com/jstjohn/SeqPrep>). We used BWA v 0.7.13 (mem algorithm) (51) to  
483 align the resultant reads to the reconstructed ancestral sequence of *M. tuberculosis*  
484 obtained in (44). Duplicated reads were marked by the Mark Duplicates module of Picard  
485 v 2.9.1 (<https://github.com/broadinstitute/picard>) and excluded. To avoid false positive  
486 calls, Pysam v 0.9.0 (<https://github.com/pysam-developers/pysam>) was used to exclude  
487 reads with an alignment score lower than  $(0.93 * \text{read\_length}) - (\text{read\_length} * 4 * 0.07)$ ,  
488 corresponding to more than 7 miss-matches per 100 bp. SNPs were called with Samtools  
489 v 1.2 mpileup (52) and VarScan v 2.4.1 (53) using the following thresholds: minimum

490 mapping quality of 20, minimum base quality at a position of 20, minimum read depth at  
491 a position of 7X and without strand bias. Only SNPs considered to have reached fixation  
492 within an isolate were considered (at a within-host frequency of  $\geq 90\%$ ). Conversely,  
493 when the within-isolate SNP frequency was  $\leq 10\%$  the ancestor state was called. Mixed  
494 infections or contaminations were discarded by excluding genomes with more than 1000  
495 variable positions with within-host frequencies between 90% and 10% and genomes for  
496 which the number of within-host SNPs was higher than the number of fixed SNPs.  
497 Additionally, we excluded genomes with average read depth  $< 15 X$  (after all the referred  
498 filtering steps). All SNPs were annotated using snpEff v4.11 (54), in accordance with the  
499 *M. tuberculosis* H37Rv reference annotation (NC\_000962.3). SNPs falling in regions  
500 such as PPE and PE-PGRS, phages, insertion sequences and in regions with at least 50  
501 bp identities to other regions in the genome were excluded from the analysis as in (55).  
502 Customized scripts were used to calculate mean coverages per gene corrected by the size  
503 of the gene. Gene deletions were determined as regions with no coverage to the reference  
504 genome.

### 505 **Phylogenetic reconstruction and ancestry estimation**

506 All 695 genomes were used to produce an alignment containing only polymorphic sites.  
507 The alignment was used to infer a maximum likelihood phylogenetic tree using the MPI  
508 parallel version of RAxML (56). We used the General Time Reversible model of  
509 nucleotide substitution under the Gamma model of rate heterogeneity and performed  
510 1000 alternative runs on distinct starting trees combined with rapid bootstrap inference.  
511 To correct the likelihood for ascertainment bias introduced by only using polymorphic  
512 site, we used Lewis correction (57). The software Tremmer (32) was used to remove  
513 redundancy in the collection of 695 whole genome SNP alignment with the stop option -  
514 *RTL* 0.95, i.e. keeping 95% of the original tree length. The resulting reduced dataset of



515 424 genomes was kept for subsequent analysis. First we used the reduced dataset plus a  
516 collection of 35 representative animal genomes to produce an alignment containing only  
517 polymorphic sites and inferred a maximum likelihood phylogenetic tree as described  
518 above. The best-scoring Maximum Likelihood topology is shown. The phylogeny was  
519 rooted using *Mycobacterium canettii*. The topology was annotated and coloured using the  
520 package *ggtree* (58) from R (59) and InkScape.

521 We inferred the biogeographic histories of L5 and L6 using Statistical-Dispersal Analysis  
522 (S-DIVA) and Bayesian Binary MCMC (BBM) Method For Ancestral State, Dispersal-  
523 Extinction-Cladogenesis (DEC), and Bayesian inference for discrete Areas (BayArea)  
524 implemented in RASP v4.0 (39). Because we did not have the geographical origin of 18  
525 samples, we used a phylogeny containing only samples from Africa where the isolation  
526 or place of birth of the patient was known. The possible ancestral ranges at each node on  
527 a selected tree were obtained. For S-DIVA the number of maximum areas was kept as 2.  
528 For BBM analysis, chains were run simultaneously for 500000 generations. The state was  
529 sampled every 100 generations. Estimated Felsenstein 1981 + Gamma was used with null  
530 root distribution.

531

### 532 **Population structure and genetic diversity**

533 Genetic structure indices and corrected pairwise SNP differences between the five  
534 African regions where genomes are grouped (Western West Africa, Easter West Africa,  
535 Central Africa, South Africa, and East Africa) were calculated using Analysis of  
536 MOlecular VAriance (AMOVA) using information on the allelic content of haplotypes,  
537 as well as their frequencies implemented in Arlequin 3.5.2.2 (60). The significance of  
538 the covariance components was tested using 20000 permutations by non-parametric  
539 permutation procedures.

540 Pairwise SNP differences and mean nucleotide diversity per site ( $\pi$ ) was calculated using  
541 the R package *ape* (61).  $\pi$  was calculated as the mean number of pairwise mismatches  
542 among L5 and L6 divided by the total length of queried genome base pairs, which  
543 comprise the total length of the genome after excluding repetitive regions (see above)  
544 (62). Confidence intervals for  $\pi$  were obtained by bootstrapping (1000 replicates) by re-  
545 sampling with replacement the nucleotide sites of the original alignments of polymorphic  
546 positions using the function *sample* in R (59). Lower and upper levels of confidence were  
547 obtained by calculating the 2.5th and the 97.5th quantiles of the  $\pi$  distribution obtained  
548 by bootstrapping. Population structure was evaluated using Principal Component  
549 Analysis (PCA) on SNP differences using *adegenet* (63) and plotted using the *plot*  
550 function in R.

551 To further explore geographical structure we evaluated the relation between the genomic  
552 phylogeny and the geographical origin of the genomes for each lineage separately using  
553 linear axis analysis in GenGISvs2.2.2 (64). The default GenGIS Africa map was used and  
554 a maximum likelihood phylogenetic tree was constructed from whole genome SNPs as  
555 described above for each lineage separately. A linear axis plot (10000 permutations) was  
556 run at significance level p-value = 0.001. If there is geographical separation, we expect  
557 the geographical distribution of the genomes to fit the phylogenetic tree structure. Fitting  
558 the tree is determined by finding a linear axis where the ordering of leaf nodes matches  
559 the ordering of sample sites according to the geographical distribution of each leaf node.  
560 If we draw a line between each leaf nodes in the phylogeny and its geographical  
561 distribution, a perfect match will result in minimum crossing of lines. Consequently,  
562 marked phylogeographical structure will result in significantly less crossing than the  
563 number of crossings expected by chance.

564 Simpson's Index (D) for geographical diversity were calculated using three different  
565 datasets: 1) the current dataset (N=424), 2) 489 L5 and L6 strains obtained from the  
566 SITVIT2 database (33), a publicly available database that contains available genotyping  
567 (spoligotyping and MIRU-VNTRs), demographic and epidemiologic information on  
568 111,635 clinical isolates, and 3) 837 genomes genotyped as L5 and L6 from 3580 strains  
569 from West Africa (38).

### 570 **Antimycobacterial resistance determining mutations and genes**

571 We have used a compiled list of resistant mutations for 11 antibiotics compiled from two  
572 independent curated datasets (65).

### 573 **Acknowledgments**

#### 574 **General**

575 Library preparation and sequencing was done in the Genomics Facility at ETH Zürich,  
576 Basel, Switzerland. Calculations were performed at sciCORE (<http://scicore.unibas.ch/>)  
577 scientific computing core facility at University of Basel.

#### 578 **Funding**

579 This work was supported by the Swiss National Science Foundation (grants  
580 310030\_188888, IZRJZ3\_164171, IZLSZ3\_170834 and CRSII5\_177163), the  
581 European Research Council (883582-ECOEVDRTB) and Wellcome (grant number  
582 098051). M.C. is supported by Ramón y Cajal program from Ministerio de Ciencia and  
583 grants from ESCMID, Ministerio de Ciencia (RTI2018-094399-A-I00) and Generalitat  
584 Valenciana (SEJI/2019/011). Authors declare no conflict of interest.

#### 585 **Author contributions:**

586 Mireia Coscolla: Conceptualization, Data curation, Formal analysis, Investigation,  
587 Methodology, Visualization, Writing – original draft

- 588 Chloe Marie Loiseau: Data curation, Methodology, Writing –review & editing
- 589 Daniela Brites: Data curation, Conceptualization, Investigation, Formal Analysis,  
590 Writing –review & editing
- 591 Fabrizio Menardo Formal analysis, Investigation, Writing –review & editing
- 592 Sonia Borrell: Resources, Writing –review & editing
- 593 C. D’Nira Sanoussi: Investigation, Writing – review & editing
- 594 Conor Meehan: Conceptualization, Investigation, Writing – review & editing
- 595 Isaac Darko Otchere: Data curation, Formal analysis, Writing – review & editing
- 596 Leonor Sanchez-Busó: Data curation, Writing – review & editing
- 597 Julian Parkhill: Conceptualization, Supervision, Writing –review & editing
- 598 Patrick Beckert: Data curation, Writing – review & editing.
- 599 Stefan Niemann: Resources, Supervision, Writing –review & editing
- 600 Dissou Affolabi: Resources, Writing – review & editing
- 601 Prince Assare: Data curation, Formal analysis, Writing – review & editing
- 602 Florian Gehre: Data curation, review & editing, Writing – review & editing
- 603 Martin Antonio: Resources, Writing – review & editing
- 604 Adwoa Asante-Poku: Data curation, Writing – review & editing
- 605 Paula Ruíz-Rodriguez: Visualization, Methodology, Writing – review & editing
- 606 Janet Fyfe: Resources, Writing – review & editing

- 607 Robin Kobbe: Resources, Writing – review & editing
- 608 Martin P. Grobusch: Resources, Writing – review & editing
- 609 Abraham S. Alabi: Resources, Writing – review & editing
- 610 Lukas Fenner: Resources, Writing – review & editing
- 611 Erik C. Boettger: Resources, Writing – review & editing
- 612 Beisel Christian: Methodology, Writing – review & editing
- 613 Simon Harris: Conceptualization, Funding acquisition, Project administration,  
614 Supervision, Writing – review & editing
- 615 Dorothy Yeboah-Manu: Conceptualization, Funding acquisition, Project administration,  
616 Supervision, Writing – review & editing
- 617 Bouke de Jong: Conceptualization, Funding acquisition, Resources, Supervision,  
618 Project administration, Writing – review & editing,
- 619 Sebastien Gagneux: Conceptualization, Funding acquisition, Resources, Supervision,  
620 Project administration, Writing – original draft.

621 **Competing interests:**

622 Authors declare no competing interest

623 **Data and materials availability**

624 All raw data generated for this study have been submitted to the European Genome-  
625 phenome Archive (EGA; <https://www.ebi.ac.uk/ega/>) under the study accession numbers  
626 PRJEB38317 and PRJEB38656. Individual runs accession number for new and published  
627 sequences are indicated in Table S1.

628

## 629 **Supplementary figures**

630 **Fig. S1. Lineage and country distribution.** Genomes analysed for the initial dataset (A)  
631 and the non-redundant dataset (B). L5 genomes are indicated in brown bars, L6 genomes  
632 in green bars and L9 genomes in light green bars.

633 **Fig. S2. Ancestral area reconstruction onto the maximum likelihood phylogeny.**  
634 Circles represent the probabilities of ancestral ranges, and the most likely ancestral areas  
635 are indicated by their corresponding color code. The inset map represents the four  
636 geographical areas considered in this analysis. Results for all four methods are shown:  
637 Bayesian binary (A), DIVA (B) DEC (C) and BayArea (C).

638

## 639 **Supplementary tables**

640 **Table S1. Genomes analysed.** Genome identifier, sequence accession numbers and  
641 sequencing statistics.

642 **Table S2. L9 specific mutations.** Synonymous and non-synonymous mutations in all  
643 lineage 9 genomes and absent in other strains from the dataset.

644 **Table S3. Pairwise  $F_{ST}$  values for L5 sublineages.**

645 **Table S4. Population average pairwise differences between L5 sublineages.**

646 **Table S5. Pairwise  $F_{ST}$  values for L6 sublineages.**

647 **Table S6. Population average pairwise differences between L6 sublineages.**

648 **Table S7. Sublineages SNPs.** SNPs defining L5 and L6 sublineages. SNPs in previously  
649 reported drug resistant genes were excluded.

650 **Table S8. Drug resistance mutations and genomes harbouring those mutations.**

651

## 652   **References**

- 653   1.     W. H. Organization, "Global tuberculosis report 2019" (ISBN 978-92-4-156571-  
654         4, 2019).
- 655   2.     M. A. Riojas, K. J. McGough, C. J. Rider-Riojas, N. Rastogi, M. H. Hazbon,  
656         Phylogenomic analysis of the species of the *Mycobacterium tuberculosis*  
657         complex demonstrates that *Mycobacterium africanum*, *Mycobacterium bovis*,  
658         *Mycobacterium caprae*, *Mycobacterium microti* and *Mycobacterium pinnipedii*  
659         are later heterotypic synonyms of *Mycobacterium tuberculosis*. *Int J Syst Evol*  
660         *Microbiol* **68**, 324-332 (2018).
- 661   3.     J. C. Semuto Ngabonziza, C. Loiseau, M. Marceau, A. Jouet, F. Menardo, O.  
662         Tzfadia, R. Antoine, E. B. Niyigena, W. Mulders, K. Fissette, M. Diels, C. Gaudin,  
663         S. Duthoy, W. Ssenooba, E. André, M. K. Kaswa, Y. M. Habimana, D. Brites, D.  
664         Affolabi, J. B. Mazarati, B. C. de Jong, L. Rigouts, S. Gagneux, C. J. Meehan, P.  
665         Supply, **A sister lineage of the *Mycobacterium tuberculosis* complex**  
666         **discovered in the AfricanGreat Lakes region.** *Nature communications*  
667         **Accepted**, (2020).
- 668   4.     D. Brites, C. Loiseau, F. Menardo, S. Borrell, M. B. Boniotti, R. Warren, A.  
669         Dippenaar, S. D. C. Parsons, C. Beisel, M. A. Behr, J. A. Fyfe, M. Coscolla, S.  
670         Gagneux, A New Phylogenetic Framework for the Animal-Adapted  
671         *Mycobacterium tuberculosis* Complex. *Front Microbiol* **9**, 2820 (2018).
- 672   5.     S. Gagneux, Ecology and evolution of *Mycobacterium tuberculosis*. *Nat Rev*  
673         *Microbiol* **16**, 202-213 (2018).
- 674   6.     S. Gagneux, K. DeRiemer, T. Van, M. Kato-Maeda, B. C. de Jong, S. Narayanan,  
675         M. Nicol, S. Niemann, K. Kremer, M. C. Gutierrez, M. Hilty, P. C. Hopewell, P. M.  
676         Small, Variable host-pathogen compatibility in *Mycobacterium tuberculosis*.  
677         *Proceedings of the National Academy of Sciences* **103**, 2869-2873 (2006).

- 678 7. R. Firdessa, S. Berg, E. Hailu, E. Schelling, B. Gumi, G. Erenso, E. Gadisa, T.  
679 Kiros, M. Habtamu, J. Hussein, J. Zinsstag, B. D. Robertson, G. Ameni, A. J.  
680 Lohan, B. Loftus, I. Comas, S. Gagneux, R. Tschopp, L. Yamuah, G. Hewinson,  
681 S. V. Gordon, D. B. Young, A. Aseffa, Mycobacterial Lineages Causing  
682 Pulmonary and Extrapulmonary Tuberculosis, Ethiopia. *Emerging infectious*  
683 *diseases* **19**, 460-463 (2013).
- 684 8. Y. Blouin, Y. Hauck, C. Soler, M. Fabre, R. Vong, C. Dehan, G. Cazajous, P. L.  
685 Massoure, P. Kraemer, A. Jenkins, E. Garnotel, C. Pourcel, G. Vergnaud,  
686 Significance of the Identification in the Horn of Africa of an Exceptionally Deep  
687 Branching *Mycobacterium tuberculosis* Clade. *PLoS ONE* **7**, (2012).
- 688 9. B. C. de Jong, M. Antonio, S. Gagneux, *Mycobacterium africanum*-Review of an  
689 Important Cause of Human Tuberculosis in West Africa. *Plos Neglected Tropical*  
690 *Diseases* **4**, (2010).
- 691 10. W. H. Haas, G. Bretzel, B. Amthor, K. Schilke, G. Krommes, S. Rusch-Gerdes,  
692 V. Sticht-Groh, H. J. Bremer, Comparison of DNA fingerprint patterns of isolates  
693 of *Mycobacterium africanum* from east and west Africa. *J Clin Microbiol* **35**, 663-  
694 666 (1997).
- 695 11. M. Kato-Maeda, P. J. Bifani, B. N. Kreiswirth, P. M. Small, The nature and  
696 consequence of genetic variability within *Mycobacterium tuberculosis*. *The*  
697 *Journal of Clinical Investigation* **107**, 533-537 (2001).
- 698 12. B. Ofori-Anyinam, A. J. Riley, T. Jobarteh, E. Gitteh, B. Sarr, T. I. Faal-Jawara,  
699 L. Rigouts, M. Senghore, A. Kehinde, N. Onyejebu, M. Antonio, B. C. de Jong, F.  
700 Gehre, C. J. Meehan, Comparative genomics shows differences in the electron  
701 transport and carbon metabolic pathways of *Mycobacterium africanum* relative to  
702 *Mycobacterium tuberculosis* and suggests an adaptation to low oxygen tension.  
703 *Tuberculosis (Edinb)* **120**, 101899 (2020).
- 704 13. M. Coscolla, S. Gagneux, Consequences of genomic diversity in *Mycobacterium*  
705 *tuberculosis*. *Seminars in Immunology* **26**, 431-444 (2014).



- 706 14. P. Asare, A. Asante-Poku, D. A. Prah, S. Borrell, S. Osei-Wusu, I. D. Otchere, A.  
707 Forson, G. Adjapong, K. A. Koram, S. Gagneux, D. Yeboah-Manu, Reduced  
708 transmission of *Mycobacterium africanum* compared to *Mycobacterium*  
709 *tuberculosis* in urban West Africa. *International journal of infectious diseases* :  
710 *IJID* : official publication of the International Society for Infectious Diseases **73**,  
711 30-42 (2018).
- 712 15. M. Huet, N. Rist, G. Boube, D. Potier, [Bacteriological study of tuberculosis in  
713 Cameroon]. *Rev Tuberc Pneumol (Paris)* **35**, 413-426 (1971).
- 714 16. G. Kallenius, T. Koivula, S. Ghebremichael, S. E. Hoffner, R. Norberg, E.  
715 Svensson, F. Dias, B. I. Marklund, S. B. Svenson, Evolution and clonal traits of  
716 *Mycobacterium tuberculosis* complex in Guinea-Bissau. *J Clin Microbiol* **37**,  
717 3872-3878 (1999).
- 718 17. B. C. de Jong, M. Antonio, S. Gagneux, *Mycobacterium africanum*—Review of  
719 an Important Cause of Human Tuberculosis in West Africa. *PLoS Negl Trop Dis*  
720 **4**, (2010).
- 721 18. S. Homolka, E. Post, B. Oberhauser, A. G. George, L. Westman, F. Dafaee, S.  
722 Rüscher-Gerdes, S. Niemann, High genetic diversity among *Mycobacterium*  
723 *tuberculosis* complex strains from Sierra Leone. *BMC Microbiol* **8**, 103 (2008).
- 724 19. B. C. de Jong, I. Adetifa, B. Walther, P. C. Hill, M. Antonio, M. Ota, R. A.  
725 Adegbola, Differences between tuberculosis cases infected with *Mycobacterium*  
726 *africanum*, West African type 2, relative to Euro-American *Mycobacterium*  
727 *tuberculosis*: an update. *FEMS Immunology & Medical Microbiology* **58**, 102-105  
728 (2010).
- 729 20. A. Asante-Poku, D. Yeboah-Manu, I. D. Otchere, S. Y. Aboagye, D. Stucki, J.  
730 Hattendorf, S. Borrell, J. Feldmann, E. Danso, S. Gagneux, *Mycobacterium*  
731 *africanum* Is Associated with Patient Ethnicity in Ghana. *PLoS Negl Trop Dis* **9**,  
732 e3370 (2015).

- 733 21. A. Asante-Poku, I. D. Otchere, S. Osei-Wusu, E. Sarpong, A. Baddoo, A. Forson,  
734 C. Laryea, S. Borrell, F. Bonsu, J. Hattendorf, C. Ahorlu, K. A. Koram, S.  
735 Gagneux, D. Yeboah-Manu, Molecular epidemiology of *Mycobacterium*  
736 *africanum* in Ghana. *BMC Infect Dis* **16**, 385 (2016).
- 737 22. D. Brites, S. Gagneux, Co-evolution of *Mycobacterium tuberculosis* and Homo  
738 sapiens. *Immunological Reviews* **264**, 6-24 (2015).
- 739 23. C. G. Meyer, G. Scarisbrick, S. Niemann, E. N. Browne, M. A. Chinbuah, J.  
740 Gyapong, I. Osei, E. Owusu-Dabo, T. Kubica, S. Rusch-Gerdes, T. Thye, R. D.  
741 Horstmann, Pulmonary tuberculosis: virulence of *Mycobacterium africanum* and  
742 relevance in HIV co-infection. *Tuberculosis (Edinb)* **88**, 482-489 (2008).
- 743 24. B. Diarra, M. Kone, A. C. G. Togo, Y. D. S. Sarro, A. B. Cisse, A. Somboro, B.  
744 Degoga, M. Tolofoudie, B. Kone, M. Sanogo, B. Baya, O. Kodio, M. Maiga, M.  
745 Belson, S. Orsega, M. Krit, S. Dao, Maiga, II, R. L. Murphy, L. Rigouts, S.  
746 Doumbia, S. Diallo, B. C. de Jong, *Mycobacterium africanum* (Lineage 6) shows  
747 slower sputum smear conversion on tuberculosis treatment than *Mycobacterium*  
748 *tuberculosis* (Lineage 4) in Bamako, Mali. *PLoS One* **13**, e0208603 (2018).
- 749 25. M. C. Hlavsa, P. K. Moonan, L. S. Cowan, T. R. Navin, J. S. Kammerer, G. P.  
750 Morlock, J. T. Crawford, P. A. LoBue, Human Tuberculosis due to *Mycobacterium*  
751 *bovis* in the United States, 1995-2005. *Clinical Infectious Diseases* **47**, 168-175  
752 (2008).
- 753 26. D. Park, H. Qin, S. Jain, M. Preziosi, J. J. Minuto, W. C. Mathews, K. S. Moser,  
754 C. A. Benson, Tuberculosis due to *Mycobacterium bovis* in Patients Coinfected  
755 with Human Immunodeficiency Virus. *Clinical Infectious Diseases* **51**, 1343-1346  
756 (2010).
- 757 27. B. C. de Jong, P. C. Hill, R. H. Brookes, J. K. Otu, K. L. Peterson, P. M. Small,  
758 R. A. Adegbola, "*Mycobacterium africanum*: a new opportunistic pathogen in HIV  
759 infection?" in *Aids* (England, 2005), vol. 19, pp. 1714-1715.

- 760 28. N. D. C. Sanoussi, B. C. deJong, M. Odoun, K. Arekpa, M. Ali Ligali, O. Bodi, S.  
761 Harris, B. Ofori-Anyinam, D. Yeboah-Manu, I. Otchere, Darko, A. Asante-Poku,  
762 S. Gagneux, M. Coscolla, L. Rigouts, D. Affolbi, Low sensitivity of the MPT64  
763 identification test to detect lineage 5 of the *Mycobacterium tuberculosis* complex.  
764 *J Med Microbiol* **67**, 1718-1727 (2018).
- 765 29. B. C. de Jong, P. C. Hill, R. H. Brookes, S. Gagneux, D. J. Jeffries, J. K. Otu, S.  
766 A. Donkor, A. Fox, K. P. McAdam, P. M. Small, R. A. Adegbola, *Mycobacterium*  
767 *africanum* elicits an attenuated T cell response to early secreted antigenic target,  
768 6 kDa, in patients with tuberculosis and their household contacts. *J Infect Dis*  
769 **193**, 1279-1286 (2006).
- 770 30. B. Ofori-Anyinam, F. Kanuteh, S. C. Agbla, I. Adetifa, C. Okoi, G. Dolganov, G.  
771 Schoolnik, O. Secka, M. Antonio, B. C. de Jong, F. Gehre, Impact of the  
772 *Mycobacterium africanum* West Africa 2 Lineage on TB Diagnostics in West Africa:  
773 Decreased Sensitivity of Rapid Identification Tests in The Gambia. *PLOS*  
774 *Neglected Tropical Diseases* **10**, e0004801 (2016).
- 775 31. S. Lipworth, R. Jajou, A. de Neeling, P. Bradley, W. van der Hoek, G. Maphalala,  
776 M. Bonnet, E. Sanchez-Padilla, R. Diel, S. Niemann, Z. Iqbal, G. Smith, T. Peto,  
777 D. Crook, T. Walker, D. van Soolingen, SNP-IT Tool for Identifying Subspecies  
778 and Associated Lineages of *Mycobacterium tuberculosis* Complex. *Emerg Infect*  
779 *Dis* **25**, 482-488 (2019).
- 780 32. F. Menardo, C. Loiseau, D. Brites, M. Coscolla, S. M. Gygli, L. K. Rutaiwa, A.  
781 Trauner, C. Beisel, S. Borrell, S. Gagneux, Treemmer: a tool to reduce large  
782 phylogenetic datasets with minimal loss of diversity. *BMC Bioinformatics* **19**, 164  
783 (2018).
- 784 33. D. Couvin, A. David, T. Zozio, N. Rastogi, Macro-geographical specificities of the  
785 prevailing tuberculosis epidemic as seen through SITVIT2, an updated version of  
786 the *Mycobacterium tuberculosis* genotyping database. *Infection, Genetics and*  
787 *Evolution*, (2018).

- 788 34. R. Brosch, S. V. Gordon, M. Marmiesse, P. Brodin, C. Buchrieser, K. Eiglmeier,  
789 T. Garnier, C. Gutierrez, G. Hewinson, K. Kremer, L. M. Parsons, A. S. Pym, S.  
790 Samper, D. van Soolingen, S. T. Cole, A new evolutionary scenario for the  
791 *Mycobacterium tuberculosis* complex. *Proc Natl Acad Sci U S A* **99**, 3684-3689  
792 (2002).
- 793 35. S. Mostowy, D. Cousins, J. Brinkman, A. Aranaz, M. A. Behr, Genomic deletions  
794 suggest a phylogeny for the *Mycobacterium tuberculosis* complex. *J Infect Dis*  
795 **186**, 74-80 (2002).
- 796 36. L. S. Ates, A. Dippenaar, F. Sayes, A. Pawlik, C. Bouchier, L. Ma, R. M. Warren,  
797 W. Sougakoff, L. Majlessi, J. W. J. van Heijst, F. Brossier, R. Brosch, Unexpected  
798 Genomic and Phenotypic Diversity of *Mycobacterium africanum* Lineage 5  
799 Affects Drug Resistance, Protein Secretion, and Immunogenicity. *Genome*  
800 *biology and evolution* **10**, 1858-1874 (2018).
- 801 37. S. Wright, Genetical Structure of Populations. *Nature* **166**, 247-249 (1950).
- 802 38. F. Gehre, S. Kumar, L. Kendall, M. Ejo, O. Secka, B. Ofori-Anyinam, E. Abatih,  
803 M. Antonio, D. Berkvens, B. C. de Jong, A Mycobacterial Perspective on  
804 Tuberculosis in West Africa: Significant Geographical Variation of *M. africanum*  
805 and Other *M. tuberculosis* Complex Lineages. *PLoS Negl Trop Dis* **10**, e0004408  
806 (2016).
- 807 39. Y. Yu, A. J. Harris, C. Blair, X. He, RASP (Reconstruct Ancestral State in  
808 Phylogenies): a tool for historical biogeography. *Mol Phylogenet Evol* **87**, 46-49  
809 (2015).
- 810 40. I. D. Otchere, M. Coscolla, L. Sanchez-Buso, A. Asante-Poku, D. Brites, C.  
811 Loiseau, C. Meehan, S. Osei-Wusu, A. Forson, C. Laryea, A. I. Yahayah, A.  
812 Baddoo, G. A. Ansa, S. Y. Aboagye, P. Asare, S. Borrell, F. Gehre, P. Beckert,  
813 T. A. Kohl, S. N'Dira, C. Beisel, M. Antonio, S. Niemann, B. C. de Jong, J. Parkhill,  
814 S. R. Harris, S. Gagneux, D. Yeboah-Manu, Comparative genomics of

- 815 *Mycobacterium africanum* Lineage 5 and Lineage 6 from Ghana suggests distinct  
816 ecological niches. *Sci Rep* **8**, 11269 (2018).
- 817 41. R. Hershberg, M. Lipatov, P. M. Small, H. Sheffer, S. Niemann, S. Homolka, J.  
818 C. Roach, K. Kremer, D. A. Petrov, M. W. Feldman, S. Gagneux, High Functional  
819 Diversity in *Mycobacterium tuberculosis* Driven by Genetic Drift and Human  
820 Demography. *PLoS Biology* **6**, e311 (2008).
- 821 42. P. Supply, M. Marceau, S. Mangenot, D. Roche, C. Rouanet, V. Khanna, L.  
822 Majlessi, A. Criscuolo, J. Tap, A. Pawlik, L. Fiette, M. Orgeur, M. Fabre, C.  
823 Parmentier, W. Frigui, R. Simeone, E. C. Boritsch, A. S. Debie, E. Willery, D.  
824 Walker, M. A. Quail, L. Ma, C. Bouchier, G. Salvignol, F. Sayes, A. Cascioferro,  
825 T. Seemann, V. Barbe, C. Loch, M. C. Gutierrez, C. Leclerc, S. D. Bentley, T. P.  
826 Stinear, S. Brisse, C. Medigue, J. Parkhill, S. Cruveiller, R. Brosch, Genomic  
827 analysis of smooth tubercle bacilli provides insights into ancestry and  
828 pathoadaptation of *Mycobacterium tuberculosis*. *Nature Genetics* **45**, 172-179  
829 (2013).
- 830 43. I. Comas, M. Coscolla, T. Luo, S. Borrell, K. E. Holt, M. Kato-Maeda, J. Parkhill,  
831 B. Malla, S. Berg, G. Thwaites, D. Yeboah-Manu, G. Bothamley, J. Mei, L. Wei,  
832 S. Bentley, S. R. Harris, S. Niemann, R. Diel, A. Aseffa, Q. Gao, D. Young, S.  
833 Gagneux, Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium*  
834 *tuberculosis* with modern humans. *Nat Genet* **45**, 1176-1182 (2013).
- 835 44. I. Comas, J. Chakravarti, P. M. Small, J. Galagan, S. Niemann, K. Kremer, J. D.  
836 Ernst, S. Gagneux, Human T cell epitopes of *Mycobacterium tuberculosis* are  
837 evolutionarily hyperconserved. *Nat Genet* **42**, 498-503 (2010).
- 838 45. M. Coscolla, R. Copin, J. Sutherland, F. Gehre, B. de Jong, O. Owolabi, G.  
839 Mbayo, F. Giardina, J. D. Ernst, S. Gagneux, *M. tuberculosis* T Cell Epitope  
840 Analysis Reveals Paucity of Antigenic Variation and Identifies Rare Variable TB  
841 Antigens. *Cell Host Microbe* **18**, 538-548 (2015).

- 842 46. C. S. Lindestam Arlehamn, S. Paul, F. Mele, C. Huang, J. A. Greenbaum, R. Vita,  
843 J. Sidney, B. Peters, F. Sallusto, A. Sette, Immunological consequences of  
844 intragenus conservation of *Mycobacterium tuberculosis* T-cell epitopes.  
845 *Proceedings of the National Academy of Sciences of the United States of*  
846 *America* **112**, E147-E155 (2015).
- 847 47. J. D. Ernst, The immunological life cycle of tuberculosis. *Nat Rev Immunol* **12**,  
848 581-591 (2012).
- 849 48. I. D. Otchere, A. Asante-Poku, S. Osei-Wusu, A. Baddoo, E. Sarpong, A. H.  
850 Ganiyu, S. Y. Aboagye, A. Forson, F. Bonsu, A. I. Yahayah, K. Koram, S.  
851 Gagneux, D. Yeboah-Manu, Detection and characterization of drug-resistant  
852 conferring genes in *Mycobacterium tuberculosis* complex strains: A prospective  
853 study in two distant regions of Ghana. *Tuberculosis (Edinb)* **99**, 147-154 (2016).
- 854 49. J. T. Belisle, M. G. Sonnenberg, Isolation of genomic DNA from mycobacteria.  
855 *Methods Mol Biol* **101**, 31-44 (1998).
- 856 50. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: a flexible trimmer for Illumina  
857 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 858 51. H. Li, R. Durbin, Fast and accurate long-read alignment with Burrows-Wheeler  
859 transform. *Bioinformatics* **26**, 589-595 (2010).
- 860 52. H. Li, A statistical framework for SNP calling, mutation discovery, association  
861 mapping and population genetical parameter estimation from sequencing data.  
862 *Bioinformatics* **27**, 2987-2993 (2011).
- 863 53. D. C. Koboldt, Q. Zhang, D. E. Larson, D. Shen, M. D. McLellan, L. Lin, C. A.  
864 Miller, E. R. Mardis, L. Ding, R. K. Wilson, VarScan 2: somatic mutation and copy  
865 number alteration discovery in cancer by exome sequencing. *Genome Res* **22**,  
866 568-576 (2012).
- 867 54. P. Cingolani, A. Platts, L. Wang le, M. Coon, T. Nguyen, L. Wang, S. J. Land, X.  
868 Lu, D. M. Ruden, A program for annotating and predicting the effects of single

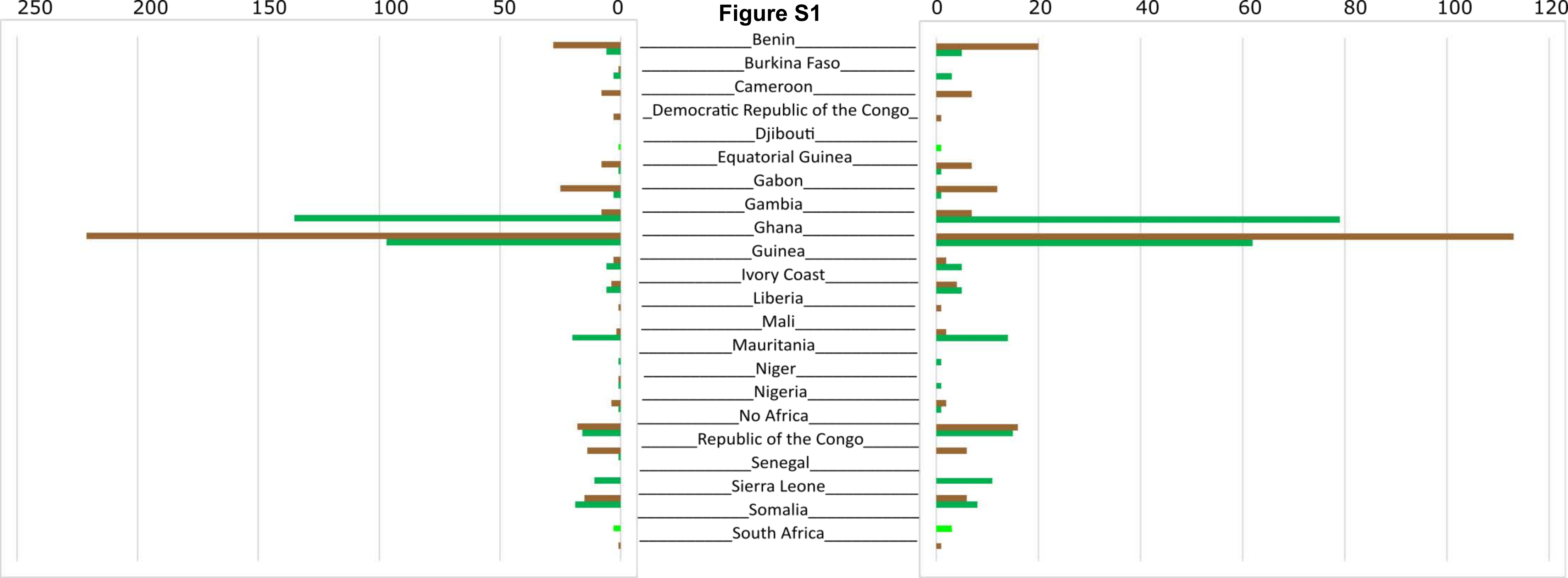
- 869 nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*  
870 *melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80-92 (2012).
- 871 55. D. Stucki, D. Brites, L. Jeljeli, M. Coscolla, Q. Liu, A. Trauner, L. Fenner, L.  
872 Rutaihwa, S. Borrell, T. Luo, Q. Gao, M. Kato-Maeda, M. Ballif, M. Egger, R.  
873 Macedo, H. Mardassi, M. Moreno, G. T. Vilanova, J. Fyfe, M. Globan, J. Thomas,  
874 F. Jamieson, J. L. Guthrie, A. Asante-Poku, D. Yeboah-Manu, E. Wampande, W.  
875 Ssengooba, M. Joloba, W. H. Boom, I. Basu, J. Bower, M. Saraiva, S. E.  
876 Vasconcellos, P. Suffys, A. Koch, R. Wilkinson, L. Gail-Bekker, B. Malla, S. D.  
877 Ley, H. P. Beck, B. C. de Jong, K. Toit, E. Sanchez-Padilla, M. Bonnet, A. Gil-  
878 Brusola, M. Frank, V. N. Penlap Beng, K. Eisenach, I. Alani, P. W. Ndung'u, G.  
879 Revathi, F. Gehre, S. Akter, F. Ntoumi, L. Stewart-Isherwood, N. E. Ntinginya, A.  
880 Rachow, M. Hoelscher, D. M. Cirillo, G. Skenders, S. Hoffner, D. Bakonyte, P.  
881 Stakenas, R. Diel, V. Crudu, O. Moldovan, S. Al-Hajoj, L. Otero, F. Barletta, E. J.  
882 Carter, L. Diero, P. Supply, I. Comas, S. Niemann, S. Gagneux, *Mycobacterium*  
883 *tuberculosis* lineage 4 comprises globally distributed and geographically  
884 restricted sublineages. *Nat Genet* **48**, 1535-1543 (2016).
- 885 56. A. Stamatakis, RAxML-VI-HPC: maximum likelihood-based phylogenetic  
886 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-  
887 2690 (2006).
- 888 57. P. O. Lewis, A likelihood approach to estimating phylogeny from discrete  
889 morphological character data. *Syst Biol* **50**, 913-925 (2001).
- 890 58. Y. Guangchuang, S. D. K., Z. Huachen, G. Yi, L. T. Tsan-Yuk, ggtree: an r  
891 package for visualization and annotation of phylogenetic trees with their  
892 covariates and other associated data. *Methods in Ecology and Evolution* **8**, 28-  
893 36 (2017).
- 894 59. R. C. Team. (Vienna, Austria, 2018).

- 895 60. L. Excoffier, G. Laval, S. Scheneider, Arlequin (version 3.0): An integrated  
896 software package for population genetics data analysis. *Evolutionary*  
897 *Bioinformatics* **1**, 47-50 (2005).
- 898 61. E. Paradis, J. Claude, K. Strimmer, APE: Analyses of Phylogenetics and  
899 Evolution in R language. *Bioinformatics* **20**, 289-290 (2004).
- 900 62. D. L. Hartl, A. G. Clarck, *Principles of population genetics* (Sinauer Associates,  
901 Inc, Sunderland, MA, 2006).
- 902 63. T. Jombart, adegenet: a R package for the multivariate analysis of genetic  
903 markers. *Bioinformatics* **24**, 1403-1405 (2008).
- 904 64. D. H. Parks, T. Mankowski, S. Zangooui, M. S. Porter, D. G. Armanini, D. J. Baird,  
905 M. G. Langille, R. G. Beiko, GenGIS 2: geospatial analysis of traditional and  
906 genetic biodiversity, with new gradient algorithms and an extensible plugin  
907 framework. *PLoS One* **8**, e69885 (2013).
- 908 65. J. L. Payne, F. Menardo, A. Trauner, S. Borrell, S. M. Gygli, C. Loiseau, S.  
909 Gagneux, A. R. Hall, Transition bias influences the evolution of antibiotic  
910 resistance in *Mycobacterium tuberculosis*. *PLoS Biol* **17**, e3000265 (2019).

911



Figure S1



**Figure S2**

bioRxiv preprint doi: <https://doi.org/10.1101/2020.06.10.141783>; this version posted June 10, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

