1    **TITLE: An African origin for *Mycobacterium bovis***

2

3

4    Chloé Loiseau[1,2*], Fabrizio Menardo[1,2*], Abraham Aseffa[3], Elena Hailu[3], Balako Gumi[4],

5    Gobena Ameni[5], Stefan Berg[6], Leen Rigouts[7,8,9], Suelee Robbe-Austerman[10], Jakob

6    Zinsstag[1,2], Sebastien Gagneux[1,2*] and Daniela Brites[1,2*]

7

8    [1]Swiss Tropical and Public Health Institute, Basel, Switzerland

9    [2]University of Basel, Basel, Switzerland

10    [3]Armauer Hansen Research Centre, Addis Ababa, Ethiopia

11    [4]Bule Hora University, Department of Animal Science and Range Management, Bule Hora

12    Town, Ethiopia

13    [5]Addis Ababa University, Aklilu Lemma Institute of Pathobiology, Addis Ababa, Ethiopia

14    [6]Animal & Plant Health Agency (APHA), Bacteriology Department, Weybridge, Surrey,

15    United Kingdom

16    [7]Mycobacteriology Unit, Department of Biomedical Sciences, Institute of Tropical Medicine,

17    Antwerp, Belgium

18    [8]Collection of Mycobacterial Cultures (BCCM/ITM), Institute of Tropical Medicine,

19    Antwerp, Belgium

20    [9]Department of Biomedical Sciences, Antwerp University, Antwerp, Belgium

21    [10]National Veterinary Services Laboratories, United States Department of Agriculture, Ames,

22    Iowa, USA

23

24    Corresponding authors:d.brites@swisstph.ch; sebastien.gagneux@swisstph.ch

25    **\*equal contribution**

26

27 **ABSTRACT**

28

29 **Background and objectives**

30 *Mycobacterium bovis* and *Mycobacterium caprae* are two of the most important agents of

31 tuberculosis (TB) in livestock and the most important causes of zoonotic TB in humans.

32 However, little is known about the global population structure, phylogeography and

33 evolutionary history of these pathogens.

34 **Methodology**

35 We compiled a global collection of 3364 whole-genome sequences from *M. bovis* and *M.*

36 *caprae* originating from 35 countries and inferred their phylogenetic relationships, geographic

37 origins and age.

38 **Results**

39 Our results resolved the phylogenetic relationship among the four previously defined clonal

40 complexes of *M. bovis*, and another eight newly described here. Our phylogeographic analysis

41 showed that *M. bovis* likely originated in East Africa. While some groups remained restricted

42 to East- and West Africa, others have subsequently dispersed to different parts of the world.

43 **Conclusions and implications**

44 Our results allow a better understanding of the global population structure of *M. bovis* and its

45 evolutionary history. This knowledge can be used to define better molecular markers for

46 epidemiological investigations of *M. bovis* in settings where whole genome sequencing

47 cannot easily be implemented.

48

49

## BACKGROUND AND OBJECTIVES

Tuberculosis (TB) remains an important burden for global health and the economy [1]. TB is the number one cause of human death due to infection globally, with an estimated 10.0 million new cases and 1.5 million deaths occurring every year [1]. TB is caused by members of the *Mycobacterium tuberculosis* complex (MTBC), which includes seven human-adapted lineages, and several animal-adapted ecotypes including *M. bovis* and *M. caprae.* Animal TB complicates the control of human TB due to the zoonotic transfer of TB bacilli from infected animals to exposed human populations e.g. through the consumption of unpasteurized milk or handling of contaminated meat [2]. *M. bovis* and *M. caprae* are the most important agents of TB in livestock and the most important agents of zoonotic TB in humans, causing an estimated 147 000 new human cases and 12 500 human deaths yearly [1, 3]. Zoonotic TB caused by *M. bovis* also poses a challenge for patient treatment, due to its natural resistance to pyrazinamide (PZA), one of the four first-line drugs used in the treatment of TB. In addition, TB in livestock accounts for an estimated loss of three billion US dollars per year [4]. In Africa, the prevalence of *M. bovis* is highest in peri-urban dairy belts of larger cities and remains at low levels in rural areas [5], often also threatening wildlife populations [6].

During the last few years, analyses of large globally representative collections of whole genome sequences (WGS) from the human-adapted MTBC lineages have enhanced our understanding of the global population structure, phylogeography and evolutionary history of these pathogens [7]. By contrast, little corresponding data exist for the various animal-adapted ecotypes of the MTBC such as *M. bovis*.

Current knowledge about global *M. bovis* populations stems mostly from spoligotyping [8, 9]. This method has been highly valuable for showing that *M. bovis* populations vary by geography, and defining strain families based on the presence or absence of spacers in the Direct Repeat region of the MTBC genome [8]. However, the discriminatory capacity of spoligotyping is limited since diversity is measured at a single locus prone to convergent evolution and phylogenetic distances cannot be reliably inferred [10].

In addition to spoligotyping, other genomic markers such as deletions [11-14] and single nucleotide polymorphisms (SNPs) [14], have given insights into the biogeography of *M. bovis*. These markers have been used to define four major groups of genotypes within *M. bovis*, known as clonal complexes European 1 and 2 (Eu1, Eu2) and African 1 and 2 (Af1 and Af2) [11-14]. Bovine TB in West Africa and East Africa is mainly caused by the clonal complexes Af1 and Af2, respectively [11, 12]. Bovine TB in Europe and in the Americas is caused by clonal complex Eu1, which affects mostly the British Islands and former trading

84    countries of the UK [13], while Eu2 is prevalent mostly in the Iberian Peninsula and Brazil

85    [14].

86    More recently, studies based on WGS have brought deeper insights into the population

87    dynamics of *M. bovis* and showed that unlike *M. tuberculosis*, wild animals can act as *M.*

88    *bovis* reservoirs in different regions of the world [15-18]. However, most studies using WGS

89    have aimed at investigating local epidemics, and little is known about the global population

90    structure and evolutionary history of *M. bovis*. Recently, we suggested a scenario for the

91    evolution of the animal-adapted MTBC, in which we propose that *M. caprae* and *M. bovis*

92    might have originally come out of Africa [19]. Here we gathered 3356 *M. bovis* and *M.*

93    *caprae* WGS from the public domain, to which we added eight new *M. bovis* sequences from

94    strains isolated in East Africa. Our results provide a phylogenetic basis to better understand

95    the global population structure of *M. bovis*. Moreover, they point to East Africa as the most

96    likely origin of contemporary *M. bovis*.

97

98

99  **METHODS**

100

101  **Data collection**

102  A total of 3929 *M. bovis* genomes were retrieved from EBI: 3834 BioSamples were registered

103  on EBI with the taxon id 1775 (corresponding to "*Mycobacterium tuberculosis* variant bovis")

104  and downloaded on the 11th of March 2019 and 95 *M. bovis* genomes were registered under

105  taxon id 1765 (corresponding to "*Mycobacterium tuberculosis*").

106  Of these, 457 were excluded because they were part of pre-publications releases from the

107  Wellcome Trust Sanger Institute, 130 were excluded because they were registered as BCG –

108  Bacille Calmette Guérin, the vaccine strain derived from *M. bovis,* one genome was excluded

109  because it was wrongly classified as *M. bovis*, and three samples were excluded because they

110  corresponded to RNA-seq libraries.

111  In addition, we added 81 publically available *M. caprae* genomes and eight previously

112  unpublished sequences from *M. bovis* isolated in Ethiopia (n=7) and Burundi (n=1). The

113  sequencing data has been deposited in the European Nucleotide Archive (EMBL-EBI) under

114  the study ID PRJEB33773.

115  From this total of 3427 genomes, 63 sequences were excluded because they did not meet our

116  criteria for downstream analyses (average whole-genome coverage below 7, ratio of

117  heterogeneous SNPs to fixed SNPs above 1), yielding a final dataset of 3364 genomes (Fig.

118  S1, Table S1). Geographical origin of the isolates, date of isolation and host metadata were

119  recovered from EBI (Table S1).

120

121  **Whole genome sequence analysis**

122  All samples were subject to the same whole-genome sequencing analysis pipeline, as

123  described in [20]. In brief, reads were trimmed with Trimmomatic v0.33 [21]. Only reads

124  larger than 20 bp were kept for the downstream analysis. The software SeqPrep

125  (https://github.com/jstjohn/SeqPrep) was used to identify and merge any overlapping paired-

126  end reads. The resulting reads were aligned to the reconstructed ancestral sequence of the

127  MTBC [22] using the MEM algorithm of BWA v0.7.13 [23] with default parameters.

128  Duplicated reads were marked using the MarkDuplicates module of Picard v2.9.1

129  (https://github.com/broadinstitute/picard). The RealignerTargetCreator and IndelRealigner

130  modules of GATK v 3.4.0 were used to perform local realignment of reads around InDels

131  [24]. Finally, SNPs were called with Samtools v1.2 mpileup [25] and VarScan v2.4.1 [26]

132  using the following thresholds: minimum mapping quality of 20, minimum base quality at a

5

133 position of 20, minimum read depth at a position of 7X, maximum strand bias for a position

134 90%. Only SNPs considered to have reached fixation within an isolate (frequency within-

135 isolate ≥90%) were considered. For SNPs with ≤10% frequency, the ancestor state was called.

136 SNPs were annotated using snpEff v4.1144 [27], using the *M. tuberculosis* H37Rv reference

137 annotation (NC_000962.3) as the genome of *M. bovis* (AF2122/97) has no genes absent from

138 H37Rv except for TbD1 (contains *mmpS6* and the 5′ region of *mmpL6*) [28].

139

140 ***In silico* spoligotyping, genomic deletions and previously defined clonal complexes**

141

142 The spoligotype pattern of the 3364 genomes was determined *in silico* using KvarQ [29]. The

143 results were submitted to the *Mycobacterium bovis* spoligotype database

144 https://www.mbovis.org/ [30] and SB numbers obtained.

145 All 3364 genomes were screened *in silico* for the presence of molecular markers defining the

146 previously described *M. bovis* clonal complexes; i.e. for the presence or absence of the

147 genomic deletions RDAf1, RDAf2, RDEu1 (also known as RD17) [11-13], and in the case of

148 Eu2 [14], for the presence of SNP 3813236 G to A with respect to the H37Rv

149 (NC_000962.3). Other deletions, RD4, RDpan and N-RD17, previously used to genotype *M.*

150 *bovis* lineages were also screened for [31-33]. The genomic coordinates in H37Rv

151 (NC_000962.3) used to determine each deletion were the following; RDAf1 (664254-

152 669601); RDAf2 (680337-694429); RDEu1 (1768074-1768878); RD4 (1696017-1708748);

153 RDpan (4371020-4373425); N-RD17 (3897069-3897783). A genomic region was considered

154 deleted if the average coverage over the region was below two.

155

156 **Phylogenetic analyses**

157 All phylogenetic trees were inferred with RAxML (v.8.2.12) using alignments containing

158 only polymorphic sites. A position was considered polymorphic if at least one genome had a

159 SNP at that position with a minimum percentage of reads supporting the call of 90%.

160 Deletions and positions not called according to the minimum threshold of 7, were encoded as

161 gaps. We excluded positions with more than 10% missing data, positions falling in PE/PPE

162 genes, phages, insertion sequences and in regions with at least 50 bp identity to other regions

163 in the genome [34]. Positions falling in drug resistance-related genes were also excluded. The

164 alignment used to produce Figure 2 comprised 22 492 variable positions and the alignment

165 used to produce Figure S2 comprised 45 981 variable positions.

166 Maximum likelihood phylogenies were computed using the general time-reversible model of

167 sequence evolution (-m GTRCAT -V options), 1,000 rapid bootstrap inferences, followed by

168 a thorough maximum-likelihood search performed through CIPRES [35]. All phylogenies

169 were rooted using a *M. africanum* Lineage (L) 6 genome from Ghana (SAMEA3359865).

170

171 **Obtaining a representative dataset of *M. bovis* genomes - Subsampling 1**

172

173 Our phylogenetic reconstruction indicated that sequences belonging to clonal complex Eu1

174 and Eu2 were over-represented in the initial 3364 genome dataset, particularly from the USA,

175 Mexico, New Zealand and the UK. To obtain a smaller dataset with a more even

176 representation of the different phylogenetic groups, we pruned the 3364 genomes using the

177 following criteria: 1) we removed all genomes with non-available country metadata (n=739),

178 which resulted in 2625 genomes; 2) we used Treemmer v0.2 [20] with the option *-RTL 99* to

179 keep 99% of the original tree length and the option *–lm* to include a list of taxa to protect from

180 pruning. This list included all genomes belonging to clonal complexes Af1 and Af2, as well

181 as any genome belonging to any unclassified clade; 3) we visually identified monophyletic

182 clades with all taxa from the same country and used Treemmer v0.2 [20], using options *-lmc*

183 and *–lm*, to only keep a few representatives of each of these clades. To have representatives of

184 the BCG clade, we kept 11 BCG genomes from [36]. This selection process rendered a

185 dataset of 476 genomes.

186

187 **Ancestral reconstruction of geographic ranges - Subsampling 2**

188

189 To infer the geographic origin of the ancestors of the main groups of *M. bovis* and *M. caprae*,

190 we used the 476 genomes dataset (see subsampling 1) and excluded all BCG genomes and all

191 *M. bovis* from human TB cases or from unknown hosts,  if the strains were isolated in a low

192 incidence TB country (Europe, North America, Oceania). This is justified by the fact that the

193 majority of such cases correspond to immigrants from high incidence countries that were

194 infected in their country of origin, i.e. country of isolation does not correspond to the native

195 geographic range of the strain and is thus not informative for the geographic reconstruction.

196 *M. bovis* from patients in high incidence countries were kept (Table S1). The resulting dataset

197 was composed of 392 genomes.

198 For the ancestral reconstruction of geographic ranges, we used the geographic origin of the

199 strains and the phylogenetic relationships of the 392 genomes. Geographic origin was treated

200    as a discrete character to which 13 states, corresponding to UN-defined regions, were

201    assigned. To select the best model of character evolution, the function fitMk from the package

202    phytools 0.6.60 in R 3.5.0 [37] was used to obtain the likelihoods of the models ER (equal-

203    rates), SYM (symmetrical) and ARD (all rates different) [38]. A Likelihood Ratio Test (LRT)

204    was used to compare the different log-Likelihoods obtained. According to the former, the best

205    fitting model was SYM, a model that allows states to transition at different rates in a

206    reversible way, i.e. reverse and forward transitions share the same parameters (Table S2). The

207    function *make.simmap* in phytools package 0.6.60 in R 3.5.0 [37, 39] was used to apply

208    stochastic character mapping as implemented in SIMMAP [40] on the 392 genomes

209    phylogeny inferred from the best-scoring ML tree rooted on L6, using the SYM model with

210    100 replicates. We summarized the results of the 100 replicates using the function *summary* in

211    phytools package 0.6.60 in R [37].

212

213

214    **Molecular Dating of *M. bovis* and *M. caprae* – Subsampling 3**

215

216    For the molecular clock analyses, we considered only genomes for which the date of isolation

217    was known (n=2058). For the eight genomes sequenced in this study, the date of isolation was

218    retrieved at a later point, and these strains were not included in the dating analysis (Table S1).

219    We used a pipeline similar to that reported in [20]. We built SNP alignments including

220    variable positions with less than 10% of missing data (alignment length; 24 828 variable

221    positions). We added an L6 strain as outgroup (SAMEA3359865) and inferred the Maximum

222    Likelihood tree as described above. Since the alignment contained only variable positions, we

223    rescaled the branch lengths of the trees: rescaled_branch_length = ((branch_length *

224    alignment_length) / (alignment_length + invariant_sites)). To evaluate the strength of the

225    temporal signal, we performed root-to-tip regression using the R package ape [41].

226    Additionally, we used the least square method implemented in LSD v0.3-beta [42] to estimate

227    the molecular clock rate in the observed data and performed a date randomization test with

228    100 randomized datasets. To do this, we used the quadratic programming dating (QPD)

229    algorithm and calculated the confidence interval (options -f 100 and -s).

230    We also estimated the molecular clock rates using a Bayesian analysis. For this, we reduced

231    the dataset to 300 strains with Treemmer v0.2 in the following way: we randomly subsampled

232    strains, maintaining the outgroup and at least one representative of four small clades of the

233    tree that would have disappeared with simple random subsampling strategy (Af2 clonal

8

234    complex: G42133; Af1 clonal complex: G02538; *PZA_ sus_unknown1*: G04143, G04145,

235    G04147; *M. caprae*: G42152, G42153, G37371, G37372, G41838; Table S1). The resulting

236    alignment included 13 012 variable sites (subset1).

237    We used jModelTest 2.1.10 v20160303 [43] to identify the best fitting nucleotide substitution

238    model among 11 possible schemes, including unequal nucleotide frequencies (total models =

239    22, options -s 11 and -f ). We performed Bayesian inference with BEAST2 [44]. We corrected

240    the xml file to specify the number of invariant sites as indicated here:

241    https://groups.google.com/forum/#!topic/beast-users/QfBHMOqImFE, and used the tip

242    sampling years to calibrate the molecular clock.

243    We used the uncorrelated lognormal relaxed clock model [45], the best fitting nucleotide

244    substitution model according to the results of jModelTest (all criteria selected the

245    transversional model (TVM) as the best model), and three different coalescent priors: constant

246    population size, exponential population growth and the Bayesian Skyline [46]. We chose a 1/x

247    prior for the population size [0- $10^9$], a 1/x prior for the mean of the lognormal distribution of

248    the clock rate [$10^{-10} - 10^{-5}$], and the standard Gamma distribution as prior for the standard

249    deviation of the lognormal distribution of the clock rate [0 – infinity]. For the exponential

250    growth rate prior, we used the standard Laplace distribution [-infinity – infinity]. For all

251    analyses, we ran two runs and used Tracer 1.7.1 [47] to evaluate convergence among runs and

252    to calculate the estimated effective sample size (ESS). We stopped the runs when they reached

253    convergence, and the ESS of the posterior and of all parameters were larger than 200. The

254    number of generations ranged from 150 to 300 million depending on the run. We used Tracer

255    [48] to identify and exclude the burn-in, which ranged from 4 to 30 million generations,

256    depending on the run.

257    Since the BEAST analysis was based on a sub-sample of the data, we tested the robustness of

258    the sub-sampling by repeating twice the sub-sampling with Treemmer [20]. This resulted in

259    two alignments of 13 272 and 12 820 SNPs, respectively (subset2 and subset3). We then

260    repeated all the BEAST analyses described above on these two additional datasets. For the

261    BEAST analyses, all trees were summarized in a maximum clade credibility tree with the

262    software Treeannotator (part of the BEAST package), after removing the burn-in and sub-

263    sampling one tree every 10 000 generations.

264

265

266     **RESULTS AND DISCUSSION**

267

268      **Phylogenetic inference of *M. bovis* and *M. caprae* populations**

269

270     The phylogenetic reconstruction of all *M. bovis* and *M. caprae* sequences obtained (n=3364)

271     confirmed that these two ecotypes correspond to two monophyletic groups, despite infecting

272     similar hosts [49, 50] (Fig. S3). The range of host species from which *M. bovis* was isolated is

273     broad, confirming that *M. bovis* can cause infection in many different mammalian species

274     (Table S1). Our collection of *M. caprae* included genomes from Japan (isolated in elephants

275     from Borneo) [51], China (isolated in primates, Table S1) and Peru (host information

276     unavailable, but possibly human [52]), suggesting that the host and geographic distribution of

277     this ecotype ranges well beyond Southern- and Central Europe [53, 54].  One group of *M.*

278     *caprae* genomes with origin in Germany contained a deletion of 36-38 kb, which

279     encompasses the region of difference RD4 (Table S1). This is in agreement with a previous

280     study reporting Alpine *M. caprae* isolates as RD4 deleted, using conventional RD typing [55].

281     For all *M. bovis* genomes, we determined *in silico* clonal complexes Eu1, Eu2, Af1 and Af2

282     and spoligotypes, and mapped them on the phylogenetic tree and onto a world map (Fig. 1,

283     Fig. S2-S3, Table S1). All previously described clonal complexes corresponded to

284     monophyletic groups in our genome-based phylogeny (Fig. S3). The phylogenetic tree also

285     revealed *M. bovis* representatives that did not fall into any of the previously described clonal

286     complexes (n=175, 5.3%, Fig. 1, Fig. S2-S3, Table S1). These belonged to eight

287     monophyletic clades with unknown classification and to a few singleton branches (Fig. S3).

288     The tree topology showed a strong bias towards closely related strains, in particular among

289     Eu1, which reflects the different sampling and WGS efforts in the different geographic

290     regions (Fig. 1, Fig. S3). Closely related genomes inform the local epidemiology but not the

291     middle/long term evolutionary history of the strains and were thus excluded from further

292     analysis (Subsampling 1, see Methods).

293     Two deep divergence events in *M. bovis* populations were notorious: one giving rise to an

294     unclassified lineage we named *M. bovis PZA_sus_unknown1* (RD4 deleted as other *M. bovis*),

295     which included five samples from Uganda (isolated from *Bos taurus* cattle, Table S1), three

296     from Malawi (isolated from humans, Table S1) and one isolated from an antelope in Germany

297     (Table S1). These *M. bovis* isolates lacked the *PncA* H57D mutation that is responsible for the

298     intrinsic pyrazinamide resistance of canonical *M. bovis* as reported previously [56] (Fig. 2).

299     This clade, retained the region of difference N-RD17, unlike all remaining *M. bovis*, and is

300  probably related to a group of strains previously isolated from cattle in Tanzania and reported

301  as "ancestral" by [31] (Fig. 2). In agreement with that, our *PZA_sus_unknown1* clade had

302  other deletions reported as specific to these Tanzanian isolates; a larger deletion

303  encompassing RDpan (RDbovis(a)_Δpan) and RDbovis(a)_kdp [31]. The second deep

304  branching lineage included all other *M. bovis* strains descendent from an ancestor that

305  acquired the *PncA* H57D mutation and therefore encompasses all previously described clonal

306  complexes [8, 14], as well as the other previously unclassified clades we describe here.

307

308  From the *M. bovis* PZA resistant ancestor strains, two main splits occurred; one split led to the

309  ancestor of Af2 and its previously unclassified sister clade which we called *unknown2* and

310  which contains the BCG vaccine strains (Fig. 2). *M. bovis* strains with spoligotyping patterns

311  similar to BCG have previously been referred to as "BCG-like". However, our genome-based

312  phylogeny shows that BCG-like spoligotyping patterns are present in several clades and have

313  thus little discriminatory power [10] (Fig. 2, Table S4). In Af2 and *unknown2* the region of

314  difference RDpan is present [31] (Fig. 2). Otherwise, RDpan is deleted in all other *M. bovis*

315  except for the *unknown3* clade, in which it is polymorphic (Fig. 2). The other split led to the

316  ancestor, from which all remaining *M. bovis* strains evolved, i.e. Af1, Eu2 and Eu1 as well as

317  other groups (Fig. 2). Interestingly, Af1 does not share a MRCA with Af2 but with Eu1 and

318  Eu2 as well as with another unclassified group, which we called *unknown3* (Fig. 2). Clonal

319  complexes Eu1 and Eu2 share a MRCA together with five other *unknown* clades (*unknown 4*

320  *– unknown 8*). Eu2 is more closely related to clades *unknown4* and *5,* than to Eu1 (Fig. 2).

321  Eu1 in turn shares a common ancestor with three other clades *unknown6*, *7* and *8* (Fig. 2).

322

323  **The temporal and geographic origin of *M. bovis***

324

325  Our reconstruction of ancestral geographical ranges points to East Africa as the most likely

326  origin for the ancestor of all *M. bovis* (Fig. 2, Fig. S4). This is supported by the fact that the

327  basal clade *M. bovis - PZA_sus_unknown1* has an exclusively East African distribution and is

328  pyrazinamide susceptible. Pyrazinamide susceptibility in *M. bovis* is probably an ancestral

329  character given that all other lineages of the MTBC are pyrazinamide susceptible.

330  Alternatively, the ancestral *M. bovis* pyrazinamide susceptible populations could have had a

331  much broader geographic distribution, which later became restricted to East Africa. For *M.*

332  *caprae,* the sampling was too small and biased (Table S1) and no conclusions can be

333  confidently drawn. We performed tip-dating calibration using the isolation dates of the strains

11

334 with both Bayesian methods and LSD (see methods). Both the tip-to-root regression and the

335 randomization tests performed indicated a temporal signal in the data (Fig. S5). We estimated

336 a clock rate of between $6.66 \times 10^{-8}$ and $1.26 \times 10^{-7}$ for the BEAST analyses (Table S3), and

337 between $6.10 \times 10^{-8}$ and $8.29 \times 10^{-8}$ for the LSD analysis. These results are in line with the

338 results of previous studies [15, 57]. The common ancestor of *M. bovis* was estimated to have

339 evolved between the years 256 and 1125 AD by the Bayesian analysis (cumulative range of

340 the 95% Highest Posterior Density (HPD) of all nine BEAST analyses) and in the year 388

341 AD by LSD (Fig. 3, TreeS1-S10). Together, these estimates suggest that *M. bovis* has

342 emerged in East Africa sometime during the period spanning the 3rd to the 12th century AD

343 (Fig. 3). However, the credibility intervals of the different analysis spanned several centuries.

344 This was due to the intrinsic uncertainty of dating ancestral nodes when the clock calibration

345 is based on the sampling time of recently sampled tips (all strains considered in this study

346 have been sampled in the last 40 years). Moreover, by relying exclusively on recent

347 calibration points to date older nodes, we ignored the issue of the time-dependency of the

348 estimated clock rates [58]. According to the time-dependency hypothesis, evolutionary rate

349 estimates depend on the age of the calibration points, with older calibration points resulting in

350 lower rates. In MTBC, this topic was discussed at length in other publications, and the

351 available data to date does not allow to confidently accept or reject that hypothesis [57, 59-

352 62]. Our analysis assumes that rates of evolution do not depend on the age of the calibration

353 points. Therefore, we potentially underestimate the age of the older nodes of the tree.

354 Molecular archaeological evidence suggests indeed that our dating analyses possibly

355 underestimate the age of the MRCA of *M. bovis*. In particular, 2000 years old *M. bovis* DNA

356 reported as having the RD4 and RD17 deletions, was found in human remains in Siberia [63].

357 The region of difference RD17 corresponds to RDEu1 [13], and defines the clonal complex

358 Eu1 of *M. bovis*, which we estimate has evolved between the years 1236 and 1603 AD (Fig.

359 3).

360 As discussed elsewhere [57], both tip-dating and the analysis of ancient DNA have potential

361 pitfalls, and these discrepancies cannot be reconciled without additional data. Nevertheless,

362 the tip-dating calibration provided accurate results for the emergency of BCG strains and for

363 the introduction of *M. bovis* to New Zealand [64, 65](TreeS1-S10), indicating that the method

364 can reliably infer divergence times at least for events occurred in the last 200 years.

365

366

367 **Insights into the detailed population structure of *M. bovis* around the world**

12

368 Understanding the evolutionary history of the *M. bovis* populations requires understanding

369 their geographic distribution at a continental scale. Our WGS data set has limited

370 geographical resolution due to the biased sampling of certain regions of the world, and to the

371 partial unavailability of associated metadata such as the origin of foreign-born TB patients

372 from Western countries. To get more insights into the geographical ranges of the different *M.*

373 *bovis* clades, we used the spoligotype patterns inferred from the WGS data and searched for

374 references describing the prevalence of those in different regions of the world (Table S4).

375 Patterns SB0120 and SB0134, known as "BCG-like" and reported to be relatively prevalent

376 [9], as well as SB0944, are phylogenetically uninformative; SB0120 is present in several

377 clades, and SB0134 and SB0944 have evolved independently in two different *M. bovis*

378 populations (Fig. 2, Table S4).

379

380 Our results suggest that the sister clade of all contemporary pyrazinamide resistant *M. bovis*,

381 *PZA_sus_unknown1,* is restricted to East Africa. The same holds true for Af2, which is in

382 accordance with previous reports [8, 12]. Our findings further suggest that the geographical

383 distribution of the Af2 sister clade *unknown2* includes East Africa (Eritrea, Ethiopia), but also

384 Southern Europe (Spain and France). Informative spoligotypes of the *unknown2* clade show

385 that it also circulates in North Africa (Fig. 2, Table S4). Of note, the original strain, from

386 which all BCG vaccine strains were derived, was isolated in France [66]. Our inferences

387 suggest that a common ancestor of Af2 and *unknown2* evolved in East Africa, and while Af2

388 remained geographically restricted, its sister clade *unknown2* has subsequently dispersed (Fig.

389 2).

390 All remaining *M. bovis* descended from a common ancestor, for which the geographical

391 origin was impossible to infer reliably with our data. However, the tree topology showed that

392 from this ancestor several clades have evolved which are important causes of bovine TB

393 today in different regions of the world (i.e., the clonal complexes Eu1, Eu2 and Af1; Fig.2).

394

395 The most basal clade within this group is *unknown3,* which contained 25 genomes mostly

396 isolated from humans (Table S1). The *in silico* derived spoligotypes suggest that the

397 geographical spread of *unknown3* ranges from Western Asia to Eastern Europe, but also

398 includes East Africa (Fig. 2, Table S4). The next split in our phylogeny corresponds to Af1,

399 which has been characterized extensively using the deletion RDAf1 and spoligotyping, and

400 shown to be most prevalent in countries from West- and Central Africa [11]. Here, we could

401 only compile nine Af1 genomes, of which five originated in Ghana [67], and the remaining

402    had either a European or an unknown origin. The small diversity of Af1 spoligotypes found in

403    our WGS dataset [11] indicates strong undersampling (Fig. 2, Table S4). Nevertheless, it was

404    possible to estimate the divergence of the Af1 clade from the remaining *M. bovis* to a period

405    ranging from the year 921 to 1449 AD (Fig. 3), making it unlikely that Af1 was originally

406    brought to West Africa by Europeans [68].

407

408    The next split comprises clades *unknown4*, *unknown5* and Eu2. Clade *unknown4* was

409    composed of 33 genomes with little geographic information and for which the most common

410    spoligotyping pattern was the uninformative SB0120 (n=19). Additional *unknown4*

411    spoligotypes indicate that strains belonging to this clade circulate in Southern Europe,

412    Northern and Eastern Africa (Fig. 2, Table S4), supporting dispersion events between Africa

413    and Southern Europe. Clade *unknown5* comprised only nine genomes isolated mostly from

414    Zambian cattle. Its corresponding spoligotype is also SB0120, limiting further geographical

415    inferences.

416    In contrast to the strains from clades *unknown4* and *unknown5*, among the 323 Eu2 genomes,

417    no genomes of East African origin were found, and Africa was only represented by nine

418    South African genomes [69]. By far, most Eu2 were isolated in the Americas. Previous

419    studies have shown that Eu2 dominates in Southern Europe, particularly in the Iberian

420    Peninsula [14], thus possibly the source of Eu2 in the Americas. There were no

421    representatives of Eu2 from the Iberian Peninsula in our dataset. However, our molecular

422    dating analysis revealed that the common ancestor of Eu2 evolved during the period 1416 to

423    1705 AD (Fig. 3), which would be compatible with an introduction from Europe into the

424    Americas.

425

426    Clonal complex Eu1, *unknown6*, *unknow7* and *unknown8*, form a sister group to the

427    previously described. Eu1 has previously been characterized based on the RDEu1 deletion

428    and spoligotyping, showing that it is highly prevalent in regions of the world that were former

429    trading partners of the UK [8, 13]. That geographic range is well represented in our dataset,

430    including many genomes from the UK (n=215) and Ireland (n=45) (Table S1). The latter were

431    very closely related, suggesting that there was probably fixation of just a few genotypes in

432    this region as previously proposed [8]. In contrast, most branching events within Eu1

433    correspond to *M. bovis* isolated in North- and Central America as well as New Zealand,

434    resulting from the expansion of clonal families not seen in the British Islands. Consequently,

435    most of the genetic diversity of Eu1 exists outside of its putative region of origin. Our

436  molecular dating is compatible with this view, indicating that the ancestor of Eu1 is likely to

437  have emerged between the years 1236 to 1603 AD (Fig. 3), with several Eu1 sub-clades

438  evolving in the last 200-300 years (TreeS1-S10).

439  The closest relative to Eu1 is a genome from Ethiopia (*unknown8*) with the spoligotyping

440  pattern SB1476, commonly found in Ethiopia [12]. *Unknown6* comprised seven genomes

441  from North America (Fig. 2, Table S1, Table S4), whereas *unknown7* included eight genomes,

442  four of which were isolated in Western Europe and another four without country of origin

443  available. Spoligotyping patterns indicate that identical strains are common in Southern

444  Europe, Northern and Eastern Africa expanding the geographic range of *unknown7*.

445

446

447

448  **CONCLUSIONS AND IMPLICATIONS**

449  We screened the public repositories and compiled 3364 genome sequences of *M. bovis* and *M.*

450  *caprae* from 35 countries. Despite the biased geographic distribution of our samples, our

451  results provide novel insights into the phylogeography of *M. bovis* and *M. caprae.* Our whole-

452  genome based phylogeny showed that although certain spoligotypes are associated with

453  specific monophyletic groups, prevalent patterns such as the so-called "BCG-like" should not

454  be used to infer phylogenetic relatedness. Moreover, our data extend the previously known

455  phylogenetic diversity of *M. bovis* by eight previously uncharacterized clades in addition to

456  the four clonal complexes described previously. Among those, *Pza_sus_unknown 1* shares a

457  common ancestor with the rest of *M. bovis,* has an exclusively East African distribution and

458  does not share the PncA mutation H57D, conferring intrinsic resistance to PZA.

459  Our further inferences suggest that *M. bovis* evolved in East Africa. The evolutionary success

460  of *M. bovis* is linked to the fact that it can infect and transmit very efficiently in cattle. Cattle

461  have been domesticated twice independently; once in the Near East (*Bos taurus*) and once in

462  the Indus Valley (*Bos indicus*) approximately 10 000 year ago, and both were introduced to

463  Africa at different time points and various locations, subsequently  interbreeding with local

464  wild species [70]. Whereas *B. taurus* was introduced probably during the 6 millennium BC

465  possibly through Egypt, *B. indicus* was most likely introduced twice, first during the second

466  millennium BC and later during the Islamic conquests [71]. *M. bovis* could have emerged

467  after the introduction of cattle, benefiting from the development of African pastoralism and

468  expanding within the continent. The timing of these events is difficult to estimate; the initial

469  introductions of cattle predate by several thousands of years our inferred temporal origin of

15

470 *M. bovis*. But as discussed, our estimates are possibly affected by the current uncertainty in

471 dating deeper evolutionary events within the MTBC. Alternatively, *M. bovis* could have

472 emerged in the Near East and been introduced to Africa together with cattle. We cannot test

473 this hypothesis formally, as the Near East is poorly represented in our dataset. However, this

474 scenario is difficult to reconcile with the restricted East African distribution of the

475 *Pza_sus_unknown_1* clade, as the Near East was also the origin of taurine cattle both in

476 Europe and Asia, where no clade retaining the ancestral characteristic of pyrazinamide

477 susceptibility was found.

478 While some *M. bovis* groups remained restricted to East Africa, others have dispersed to

479 different parts of the world. The contemporary geographic distribution of *M. bovis* clades

480 suggest that East- and North Africa, Southern Europe and Western Asia have played an

481 important role in shaping the population structure of these pathogens. However, these regions

482 were not well represented in our dataset. Thus, more *M. bovis* genomes from these regions are

483 necessary to generate better insights, particularly given the central role of these regions in the

484 history of cattle domestication [72]. From a more applied perspective, our work provides a

485 global phylogenetic framework that can be further exploited to find better molecular markers

486 for studying *M. bovis* in settings where genome sequencing cannot be easily implemented.

487

488

494

495

**REFERENCES**

497

498 1.     WHO; Global tuberculosis report 2018. Geneva: World Health Organization, 2018.

499 2.     Olea-Popelka F, Muwonge A, Perera A, et al.; Zoonotic tuberculosis in human beings
500 caused by Mycobacterium bovis-a call for action. *Lancet Infect Dis* 2017;**17**(1):e21-e25. doi:
501 10.1016/S1473-3099(16)30139-6.

502 3.     Muller B, Durr S, Alonso S, et al.; Zoonotic Mycobacterium bovis-induced
503 tuberculosis in humans. *Emerg Infect Dis* 2013;**19**(6):899-908. doi: 10.3201/eid1906.120543.

16

504   4.      Waters WR, Palmer MV, Buddle BM, et al.; Bovine tuberculosis vaccine research:
505   historical perspectives and recent advances. *Vaccine* 2012;**30**(16):2611-22.

506   5.      Tschopp R, Hattendorf J, Roth F, et al.; Cost estimate of bovine tuberculosis to
507   Ethiopia. *Current topics in microbiology and immunology* 2013;**365**:249-68.

508   6.      Michel AL, Muller B, van Helden PD; Mycobacterium bovis at the animal-human
509   interface: a problem, or not? *Veterinary microbiology* 2010;**140**(3-4):371-81.

510   7.      Gagneux S; Ecology and evolution of Mycobacterium tuberculosis. *Nat Rev Microbiol*
511   2018;**16**(4):202-213. doi: 10.1038/nrmicro.2018.8.

512   8.      Smith NH; The global distribution and phylogeography of Mycobacterium bovis
513   clonal complexes. *Infect Genet Evol* 2012;**12**(4):857-65. doi: 10.1016/j.meegid.2011.09.007.

514   9.      Ghavidel M, Mansury D, Nourian K, et al.; The most common spoligotype of
515   Mycobacterium bovis isolated in the world and the recommended loci for VNTR typing; A
516   systematic review. *Microbial Pathogenesis* 2018;**118**:310-315.

517   10.     Comas I, Homolka S, Niemann S, et al.; Genotyping of genetically monomorphic
518   bacteria: DNA sequencing in mycobacterium tuberculosis highlights the limitations of current
519   methodologies. *PLoS ONE* 2009;**4**(11):e7815. doi: 10.1371/journal.pone.0007815.

520   11.     Muller B, Hilty M, Berg S, et al.; African 1; An Epidemiologically Important Clonal
521   Complex of Mycobacterium bovis Dominant in Mali, Nigeria, Cameroon and Chad. *J
522   Bacteriol* 2009;**191**(6): 1951-1960. doi: 10.1128/JB.01590-08.

523   12.     Berg S, Garcia-Pelayo MC, Muller B, et al.; African 2, a clonal complex of
524   Mycobacterium bovis epidemiologically important in East Africa. *J Bacteriol*
525   2011;**193**(3):670-8. doi: 10.1128/JB.00750-10.

526   13.     Smith NH, Berg S, Dale J, et al.; European 1: a globally important clonal complex of
527   Mycobacterium bovis. *Infect Genet Evol* 2011;**11**(6):1340-51. doi:
528   10.1016/j.meegid.2011.04.027.

529   14.     Rodriguez-Campos S, Schurch AC, Dale J, et al.; European 2--a clonal complex of
530   Mycobacterium bovis dominant in the Iberian Peninsula. *Infect Genet Evol* 2012;**12**(4):866-
531   72. doi: 10.1016/j.meegid.2011.09.004.

532   15.     Crispell J, Zadoks RN, Harris SR, et al.; Using whole genome sequencing to
533   investigate transmission in a multi-host system: bovine tuberculosis in New Zealand. *BMC
534   genomics* 2017;**18**(1):180.

535   16.     Orloski K, Robbe-Austerman S, Stuber T, et al.; Whole Genome Sequencing of
536   Mycobacterium bovis Isolated From Livestock in the United States, 1989-2018. *Front Vet Sci*
537   2018;**5**:253. doi: 10.3389/fvets.2018.00253.

538    17.    Salvador LCM, O'Brien DJ, Cosgrove MK, et al.; Disease management at the wildlife-
539    livestock interface: Using whole-genome sequencing to study the role of elk in
540    Mycobacterium bovis transmission in Michigan, USA. *Mol Ecol* 2019;**28**(9):2192-2205. doi:
541    10.1111/mec.15061.

542    18.    Price-Carter M, Brauning R, de Lisle GW, et al.; Whole Genome Sequencing for
543    Determining the Source of Mycobacterium bovis Infections in Livestock Herds and Wildlife
544    in New Zealand. *Front Vet Sci* 2018;**5**:272. doi: 10.3389/fvets.2018.00272.

545    19.    Brites D, Loiseau C, Menardo F, et al.; A New Phylogenetic Framework for the
546    Animal-Adapted Mycobacterium tuberculosis Complex. *Frontiers in microbiology*
547    2018;**9**:2820.

548    20.    Menardo F, Loiseau C, Brites D, et al.; Treemmer: a tool to reduce large phylogenetic
549    datasets with minimal loss of diversity. *BMC Bioinformatics* 2018;**19**(1):164. doi:
550    10.1186/s12859-018-2164-8.

551    21.    Bolger AM, Lohse M, Usadel B; Trimmomatic: a flexible trimmer for Illumina
552    sequence data. *Bioinformatics* 2014;**30**(15):2114-20. doi: 10.1093/bioinformatics/btu170.

553    22.    Comas I, Chakravartti J, Small PM, et al.; Human T cell epitopes of Mycobacterium
554    tuberculosis are evolutionarily hyperconserved. *Nat Genet* 2010;**42**(6):498-503. doi:
555    10.1038/ng.590.

556    23.    Li H, Handsaker B, Wysoker A, et al.; The Sequence Alignment/Map format and
557    SAMtools. *Bioinformatics* 2009;**25**(16):2078-9. doi: btp352 [pii]
558    10.1093/bioinformatics/btp352.

559    24.    McKenna A, Hanna M, Banks E, et al.; The Genome Analysis Toolkit: a MapReduce
560    framework for analyzing next-generation DNA sequencing data. *Genome Res*
561    2010;**20**(9):1297-303. doi: 10.1101/gr.107524.110.

562    25.    Li H; A statistical framework for SNP calling, mutation discovery, association
563    mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*
564    2011;**27**(21):2987-93. doi: 10.1093/bioinformatics/btr509.

565    26.    Koboldt DC, Zhang Q, Larson DE, et al.; VarScan 2: somatic mutation and copy
566    number alteration discovery in cancer by exome sequencing. *Genome Res* 2012;**22**(3):568-76.
567    doi: 10.1101/gr.129684.111.

568    27.    Cingolani P, Platts A, Wang le L, et al.; A program for annotating and predicting the
569    effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila
570    melanogaster strain w1118; iso-2; iso-3. *Fly (Austin)* 2012;**6**(2):80-92. doi:
571    10.4161/fly.19695.

572    28.    Garnier T, Eiglmeier K, Camus JC, et al.; The complete genome sequence of
573    Mycobacterium bovis. *Proc Natl Acad Sci U S A* 2003;**100**(13):7877-82.

574  29.    Steiner A, Stucki D, Coscolla M, et al.; KvarQ: targeted and direct variant calling from
575       fastq reads of bacterial genomes. *BMC Genomics* 2014;**15**:881. doi: 10.1186/1471-2164-15-
576       881.

577  30.    Smith NH, Upton P; Naming spoligotype patterns for the RD9-deleted lineage of the
578       Mycobacterium tuberculosis complex; www.Mbovis.org. *Infection, genetics and evolution :*
579       *journal of molecular epidemiology and evolutionary genetics in infectious diseases*
580       2012;**12**(4):873-6.

581  31.    Mostowy S, Inwald J, Gordon S, et al.; Revisiting the evolution of Mycobacterium
582       bovis. *J Bacteriol* 2005;**187**(18):6386-95.

583  32.    Brosch R, Gordon SV, Garnier T, et al.; Genome plasticity of BCG and impact on
584       vaccine efficacy. *Proc Natl Acad Sci U S A* 2007.

585  33.    Salamon H, Kato-Maeda M, Small PM, et al.; Detection of deleted genomic DNA
586       using a semiautomated computational analysis of GeneChip data. *Genome Res*
587       2000;**10**(12):2044-54.

588  34.    Stucki D, Brites D, Jeljeli L, et al.; Mycobacterium tuberculosis lineage 4 comprises
589       globally distributed and geographically restricted sublineages. *Nat Genet* 2016;**48**(12):1535-
590       1543. doi: 10.1038/ng.3704.

591  35.    Miller MA, Pfeiffer W, Schwartz T; Creating the CIPRES Science Gateway for
592       inference of large phylogenetic trees. *Proceedings of the Gateway Computing Environments*
593       *Workshop (GCE)*. New Orleans, LA, 2010, 1-8.

594  36.    Copin R, Coscolla M, Efstathiadis E, et al.; Impact of in vitro evolution on antigenic
595       diversity of Mycobacterium bovis bacillus Calmette-Guerin (BCG). *Vaccine*
596       2014;**32**(45):5998-6004. doi: 10.1016/j.vaccine.2014.07.113.

597  37.    Revell LJ; phytools: an R package for phylogenetic comparative biology (and other
598       things). *Methods in Ecology and Evolution* 2011;**3**(2):217-223.

599  38.    Lewis PO; A likelihood approach to estimating phylogeny from discrete
600       morphological character data. *Syst Biol* 2001;**50**(6):913-25. doi:
601       10.1080/106351501753462876.

602  39.    Team RC; R: A language and environment for statistical computing. R Foundation for
603       Statistical Computing. Vienna, Austria, 2019.

604  40.    Bollback JP; SIMMAP: stochastic character mapping of discrete traits on phylogenies.
605       *BMC Bioinformatics* 2006;**7**:88. doi: 10.1186/1471-2105-7-88.

606  41.    Paradis E, Schliep K; ape 5.0: an environment for modern phylogenetics and
607       evolutionary analyses in R. *Bioinformatics* 2019;**35**(3):526-528. doi:
608       10.1093/bioinformatics/bty633.

609    42.    To TH, Jung M, Lycett S, et al.; Fast Dating Using Least-Squares Criteria and
610    Algorithms. *Systematic Biology* 2016;**65**(1):82-97. doi: 10.1093/sysbio/syv068.

611    43.    Darriba D, Taboada GL, Doallo R, et al.; jModelTest 2: more models, new heuristics
612    and parallel computing. *Nature Methods* 2012;**9**(8):772-772. doi: DOI 10.1038/nmeth.2109.

613    44.    Bouckaert R, Heled J, Kuhnert D, et al.; BEAST 2: a software platform for Bayesian
614    evolutionary analysis. *PLoS Comput Biol* 2014;**10**(4):e1003537. doi:
615    10.1371/journal.pcbi.1003537.

616    45.    Drummond AJ, Ho SYW, Phillips MJ, et al.; Relaxed phylogenetics and dating with
617    confidence. *Plos Biology* 2006;**4**(5):699-710. doi: ARTN e88
618    10.1371/journal.pbio.0040088.

619    46.    Drummond AJ, Rambaut A, Shapiro B, et al.; Bayesian coalescent inference of past
620    population dynamics from molecular sequences. *Mol Biol Evol* 2005;**22**(5):1185-92. doi:
621    msi103 [pii]
622    10.1093/molbev/msi103.

623    47.    Rambaut A, Drummond AJ, Xie D, et al.; Posterior Summarization in Bayesian
624    Phylogenetics Using Tracer 1.7. *Systematic Biology* 2018;**67**(5):901-904. doi:
625    10.1093/sysbio/syy032.

626    48.    Rambaut A, Drummond AJ, Xie D, et al.; Posterior Summarization in Bayesian
627    Phylogenetics Using Tracer 1.7. *Syst Biol* 2018;**67**(5):901-904. doi: 10.1093/sysbio/syy032.

628    49.    Rodriguez S, Bezos J, Romero B, et al.; Mycobacterium caprae Infection in Livestock
629    and Wildlife, Spain. *Emerging Infectious Diseases* 2011;**17**(3):532-535. doi:
630    10.3201/eid1703.100618.

631    50.    Prodinger WM, Brandstatter A, Naumann L, et al.; Characterization of
632    Mycobacterium caprae isolates from Europe by mycobacterial interspersed repetitive unit
633    genotyping. *J Clin Microbiol* 2005;**43**(10):4984-92.

634    51.    Yoshida S, Suga S, Ishikawa S, et al.; Mycobacterium caprae Infection in Captive
635    Borneo Elephant, Japan. *Emerg Infect Dis* 2018;**24**(10):1937-1940. doi:
636    10.3201/eid2410.180018.

637    52.    Consortium C, the GP, Allix-Beguec C, et al.; Prediction of Susceptibility to First-
638    Line Tuberculosis Drugs by DNA Sequencing. *N Engl J Med* 2018;**379**(15):1403-1415. doi:
639    10.1056/NEJMoa1800474.

640    53.    Aranaz A, Cousins D, Mateos A, et al.; Elevation of Mycobacterium tuberculosis
641    subsp. caprae Aranaz et al. 1999 to species rank as Mycobacterium caprae comb. nov., sp.
642    nov. *Int J Syst Evol Microbiol* 2003;**53**(Pt 6):1785-9.

643  54.     Broeckl S, Krebs S, Varadharajan A, et al.; Investigation of intra-herd spread of
644  Mycobacterium caprae in cattle by generation and use of a whole-genome sequence. *Vet Res*
645  *Commun* 2017;**41**(2):113-128. doi: 10.1007/s11259-017-9679-8.

646  55.     Domogalla J, Prodinger WM, Blum H, et al.; Region of difference 4 in alpine
647  Mycobacterium caprae isolates indicates three variants. *J Clin Microbiol* 2013;**51**(5):1381-8.
648  doi: 10.1128/JCM.02966-12.

649  56.     Loiseau C, Brites D, Moser I, et al.; Revised Interpretation of the Hain Lifescience
650  GenoType MTBC To Differentiate Mycobacterium canettii and Members of the
651  Mycobacterium tuberculosis Complex. *Antimicrobial Agents and Chemotherapy* 2019;**63**(6).
652  doi: ARTN e00159-19
653  10.1128/AAC.00159-19.

654  57.     Menardo F, Duchene S, Brites D, et al.; The molecular clock of Mycobacterium
655  tuberculosis. *PLoS Pathog* 2019;**15**(9):e1008067. doi: 10.1371/journal.ppat.1008067.

656  58.     Ho SY, Phillips MJ, Cooper A, et al.; Time dependency of molecular rate estimates
657  and systematic overestimation of recent divergence times. *Mol Biol Evol* 2005;**22**(7):1561-8.
658  doi: 10.1093/molbev/msi145.

659  59.     Comas I, Coscolla M, Luo T, et al.; Out-of-Africa migration and Neolithic
660  coexpansion of Mycobacterium tuberculosis with modern humans. *Nat Genet*
661  2013;**45**(10):1176-82. doi: 10.1038/ng.2744.

662  60.     Pepperell CS, Casto AM, Kitchen A, et al.; The role of selection in shaping diversity
663  of natural M. tuberculosis populations. *PLoS Pathog* 2013;**9**(8):e1003543. doi:
664  10.1371/journal.ppat.1003543.

665  61.     Bos KI, Harkins KM, Herbig A, et al.; Pre-Columbian mycobacterial genomes reveal
666  seals as a source of New World human tuberculosis. *Nature* 2014;**514**(7523):494-7. doi:
667  10.1038/nature13591.

668  62.     Eldholm V, Pettersson JH, Brynildsrud OB, et al.; Armed conflict and population
669  displacement as drivers of the evolution and dispersal of Mycobacterium tuberculosis. *Proc*
670  *Natl Acad Sci U S A* 2016;**113**(48):13881-13886. doi: 10.1073/pnas.1611283113.

671  63.     Taylor GM, Murphy E, Hopkins R, et al.; First report of Mycobacterium bovis DNA
672  in human remains from the Iron Age. *Microbiology* 2007;**153**(Pt 4):1243-9.

673  64.     Behr MA, Small PM; A historical and molecular phylogeny of BCG strains. *Vaccine*
674  1999;**17**(7-8):915-922.

675  65.     Binney BM, Biggs PJ, Carter PE, et al.; Quantification of historical livestock
676  importation into New Zealand 1860-1979. *N Z Vet J* 2014;**62**(6):309-14. doi:
677  10.1080/00480169.2014.914861.

678    66.      Oettinger T, Jorgensen M, Ladefoged A, et al.; Development of the Mycobacterium
679    bovis BCG vaccine: review of the historical and biochemical evidence for a genealogical tree.
680    *Tuber Lung Dis* 1999;**79**(4):243-50. doi: 10.1054/tuld.1999.0206.

681    67.      Otchere ID, van Tonder AJ, Asante-Poku A, et al.; Molecular epidemiology and whole
682    genome sequencing analysis of clinical Mycobacterium bovis from Ghana. *PLoS One*
683    2019;**14**(3):e0209395. doi: 10.1371/journal.pone.0209395.

684    68.      Muwonge A, Franklyn E, Mark B, et al.; Molecular Epidemiology of Mycobacterium
685    bovis in Africa. In: B. DA, J. KNP, O. TCs (eds). *Tuberculosis in Animals: An African
686    Perspective*. Switzerland: Springer, 2019, 127-170.

687    69.      Dippenaar A, Parsons SDC, Miller MA, et al.; Progenitor strain introduction of
688    Mycobacterium bovis at the wildlife-livestock interface can lead to clonal expansion of the
689    disease in a single ecosystem. *Infect Genet Evol* 2017;**51**:235-238. doi:
690    10.1016/j.meegid.2017.04.012.

691    70.      Loftus RT, MacHugh DE, Bradley DG, et al.; Evidence for two independent
692    domestications of cattle. *Proc Natl Acad Sci U S A* 1994;**91**(7):2757-61.

693    71.      Verdugo MP, Mullin VE, Scheu A, et al.; Ancient cattle genomics, origins, and rapid
694    turnover in the Fertile Crescent. *Science* 2019;**365**(6449):173-176. doi:
695    10.1126/science.aav1002.

696    72.      Decker JE, McKay SD, Rolf MM, et al.; Worldwide Patterns of Ancestry, Divergence,
697    and Admixture in Domesticated Cattle. *Plos Genetics* 2014;**10**(3). doi: ARTN e1004254
698    10.1371/journal.pgen.1004254.

699    73.      Rambaut A; FigTree. Edinburgh: Institute of Evolutionary Biology, University of
700    Edinburgh, 2010.
701
702

703    **FIGURE LEGENDS**

704

705    **Figure 1 –** Geographic distribution of the *M. bovis* samples used in this study according to

706    isolation country. The circles correspond to pie charts and are coloured according to clonal

707    complexes.

708

709    **Figure 2** – Maximum likelihood phylogeny of 476 of the 3364 genomes included in this

710    study (redundant genomes were removed), and inferred from 22 492 variable positions. The

711    scale bar indicates the number of substitutions per polymorphic site. The phylogeny is rooted

712    on a *M. tuberculosis* Lineage 6 genome from Ghana (not shown) and bootstrap values are

713   shown for the most important splits. The coloured bars on the side of the phylogeny show the

714   different clonal complexes. Other "unknown" monophyletic clades are coloured in black and

715   additionally the branches of the eight clades are coloured to show their phylogenetic position

716   more precisely. The pie charts mapped on the tree represent the summary posterior

717   probabilities (from 100 runs) of the reconstructed ancestral geographic states and are coloured

718   according to geographical UN region. Inferred spoligotype patterns from WGS described in

719   *M. bovis* spoligotype database [30] are indicated for the unknown clades. The red circles at

720   the tips correspond to the eight newly sequenced genomes. Regions of difference (RD) as in

721   [31] are indicated; superscript + and - refers to presence of the region or its deletion,

722   respectively.

723

724

725   **Figure 3** – The inferred age of main monophyletic clades according to LSD and BEAST

726   dating analyses. For BEAST we report the results of the two analyses that resulted in the

727   lowest clock rate (subset1 Bayesian Skyline) and in the highest clock rate (subset3

728   exponential population growth). The confidence intervals reported correspond to the merged

729   HPD interval of the two BEAST analyses mentioned above. The BEAST analysis was based

730   on 300 genomes and the LSD analysis was based on 2058 genomes (see methods section for

731   subsampling strategy). Only one genome from the Af1 clonal complex was included in the

732   dating analyses and therefore the dates reported correspond to the node where Af1 diverged.

733

734   **Supplementary Figures:**

735

736   **Figure S1** – Flow chart showing the selection of genomes.

737

738   **Figure S2** - Geographic distribution of the *M. bovis* samples with unknown classification

739   used in this study according to isolation country.

740

741   **Figure S3** – Maximum likelihood phylogeny of all 3364 genomes, based on 45 981 variable

742   positions. The scale bar indicates the number of substitutions per polymorphic site. The

743   phylogeny is rooted on a *M. tuberculosis* Lineage 6 genome from Ghana. The outer ring

744   indicates the geographical region from which the strains were isolated. The four clonal

745   complexes are highlighted on the tree. Branches corresponding to BCG genomes are coloured

746   in grey and the *PncA* mutation H57D is indicated by a yellow star.

23

747

748     **Figure S4** – Phylogeographic reconstruction of *M. bovis* and *M. caprae*, inferred from 392

749     genomes. Thirteen UN-defined geographic regions were assigned to the discrete character

750     geographic origin, and mapped onto the phylogeny. Pie charts at internal nodes represent the

751     summary posterior probabilities (from 100 runs) of the reconstructed ancestral geographic

752     states and are coloured according to geographical UN region**.**

753

754     **Figure S5** – A) Tip-to-root regression and B) Date randomization tests (DTR). The

755     confidence interval of the clock rate estimate for the observed data does not overlap with the

756     confidence intervals of the clock rate estimates obtained from the randomized sets.

757

758     **Supplemental Tables:**

759

760     **Table S1 -** List of genomes included in this study along with metadata used for the analyses.

761

762     **Table S2 -** Comparison of models for discrete character evolution using likelihood ratio tests.

763

764     **Table S3 –** Results of all BEAST analyses.

765

766     **Table S4 -** Spoligotype patterns determined *in silico* for different clonal complex groups with

767     reference to other studies.

768

769     **Supplementary files:**

770     **TreeS1-S10** Ten time-calibrated trees resulted from the molecular clock analyses. The file

771     names indicate the software used (LSD or BEAST), the subsample, and the coalescent

772     population prior (BSP: Bayesian Skyline; exponential: exponential population growth;

773     constant: constant population size).  Tip labels are present in Table S1.  Ages in years before

774     present can be visualized as well as the 95% High Posterior Density (HPD) (BEAST trees)
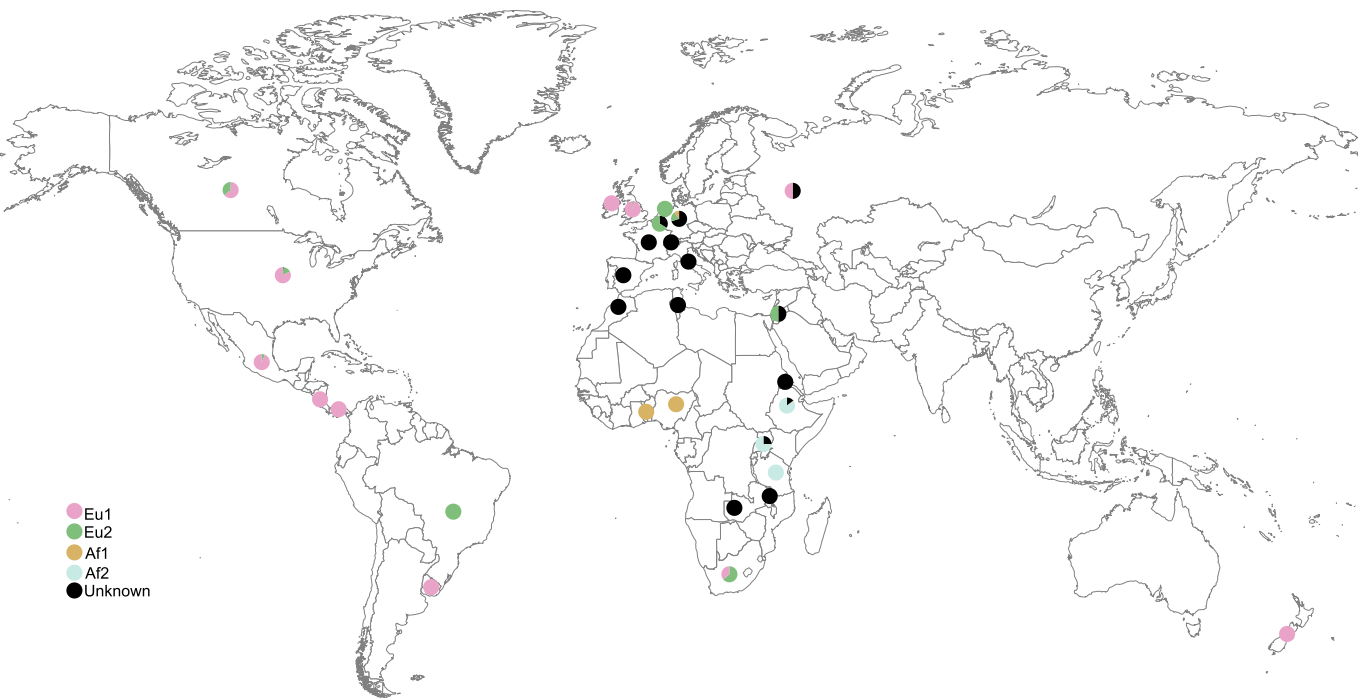
775     using FigTree [73].

776

777

778

779

780