



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

# Unraveling the epidemiology of *Mycobacterium bovis* using whole-genome sequencing combined with environmental and demographic data

### Citation for published version:

Rossi, G, Shih, B, Egbe, FN, Motta, P, Duchatel, F, Kelly, R, Ndip, L, Sander, M, Tanya, V, Lycett, S, Bronsvort, M & Muwonge, A 2023, 'Unraveling the epidemiology of *Mycobacterium bovis* using whole-genome sequencing combined with environmental and demographic data', *Frontiers in Veterinary Science*, vol. 10, pp. 1-16. <https://doi.org/10.3389/fvets.2023.1086001>

### Digital Object Identifier (DOI):

[10.3389/fvets.2023.1086001](https://doi.org/10.3389/fvets.2023.1086001)

### Link:

[Link to publication record in Edinburgh Research Explorer](#)

### Document Version:

Peer reviewed version

### Published In:

Frontiers in Veterinary Science

### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1 **Title**

2 ***Unravelling the epidemiology of Mycobacterium bovis using whole genome sequencing***  
3 ***combined with environmental and demographic data***

4

5 **Authors**

6 *Gianluigi Rossi<sup>1,2+</sup>, Barbara Shih<sup>1</sup>, Franklyn N. Egbe<sup>3</sup>, Paolo Motta<sup>4</sup>, Florian Duchatel<sup>1</sup>, Robert*  
7 *Kelly<sup>5</sup>, Lucy Ndip<sup>6,7</sup>, Melissa Sander<sup>8</sup>, Vincent Tanya<sup>9</sup>, Samantha J. Lycett<sup>1,2</sup>, Barend Mark de*  
8 *Clare Bronsvort<sup>1,2</sup>, Adrian Muwonge<sup>1</sup>*

9

10 **Affiliations**

11 <sup>1</sup>The Roslin Institute, R(D)SVS, University of Edinburgh, Easter Bush Campus, Midlothian,  
12 EH25 9RG, Scotland

13 <sup>2</sup>Centre of Expertise on Animal Diseases Outbreaks, EPIC, Scotland

14 <sup>3</sup>School of Life Sciences, University of Lincoln, Brayford Pool, Lincoln LN6 7TS United  
15 Kingdom

16 <sup>4</sup>The European Commission for the Control of Foot-and-Mouth Disease (EuFMD), Food and  
17 Agriculture Organization of the United Nations, Rome, Italy

18 <sup>5</sup>Royal (Dick) School of Veterinary Studies and The Roslin Institute, University of Edinburgh,  
19 Easter Bush, Midlothian, EH25 9RG, United Kingdom

20 <sup>6</sup>Laboratory for Emerging Infectious Diseases, University of Buea, Buea, Cameroon

21 <sup>7</sup>Department of Biomedical Sciences, Faculty of Health Sciences, University of Buea, Buea,  
22 Cameroon

23 <sup>8</sup>Tuberculosis Reference Laboratory, Bamenda, P.O. Box 586, Cameroon

24 <sup>9</sup>Cameroon Academy of Sciences, P.O. Box 1457, Yaoundé, Cameroon

25 <sup>+</sup>Corresponding author: [g.rossi@ed.ac.uk](mailto:g.rossi@ed.ac.uk)

26

27 **Author Contributions**

28 BCMB, LN, and VT conceived the original project. BCMB, RK, LN, VT, MS, and NE designed  
29 the field study, the databases, and the survey instrument. RK, NE, VT, and BB developed the  
30 field SOPs and collected the data. PM collected the cattle movement data. GR, AM, FD and  
31 FE cleaned the data. GR, BCMB, SJL and AM conceived the quantitative analysis. GR and SJL  
32 run the phylogenetic analysis. BS implemented the bioinformatical pipelines. GR performed

33 the machine learning analysis. GR was responsible for writing the initial drafts. All authors  
34 contributed comments for the final draft. All authors contributed to the article and  
35 approved the submitted version.

36

### 37 **Acknowledgements**

38 We want to thank Dr. Isaac D. Otchere for providing the isolation year of the *M. bovis*  
39 sequences sampled in Ghana, Dr. Stella Mazeri and Dr. William Harvey, Roslin Institute  
40 University of Edinburgh, for data collection and analysis support, and Prof. Kim Vanderwaal,  
41 University of Minnesota for the useful comments.

42

### 43 **Fundings**

44 The primary data used in this work was generated with funding from Wellcome Trust  
45 (WT094945), with BCMB as the Principal investigator. BS, SJL, MB, BMCB, and AM are  
46 supported by Biotechnology and Biological Sciences Research Council (BBSRC) programme  
47 grant to Roslin Institute (Award Numbers BBS/E/D/20002172 and BBS/E/D/200021723), and  
48 AM later by his BBSRC Future leader Fellowship and current Chancellor's fellowship. GR, SJL,  
49 BCMB additional received support from received support from the Scottish Government  
50 Rural and Environment Science and Analytical Services Division as part of Centre of  
51 Expertise on Animal Disease Outbreaks (EPIC). We are also grateful to the staff of MENIPIA,  
52 especially the veterinarians and delegates who diligently supported the primary fieldwork  
53 that generated this data. BS was partially funded by a BBSRC Core Capability Grant  
54 BB/CCG1780/1 awarded to The Roslin Institute. *M. bovis* sequencing was carried out by  
55 Edinburgh Genomics, The University of Edinburgh, which is partly supported through core  
56 grants from NERC (R8/H10/56), MRC (MR/K001744/1) and BBSRC (BB/J004243/1).

57

58

59 **Abstract**

60 When studying the dynamics of a pathogen in a host population, one crucial question is  
61 whether it transitioned from an epidemic (i.e. the pathogen population and the number of  
62 infected hosts are increasing) to an endemic stable state (i.e. the pathogen population  
63 reached an equilibrium). For slow-growing and slow-evolving clonal pathogens like  
64 *Mycobacterium bovis*, the causative agent of bovine (or animal) and zoonotic tuberculosis, it  
65 can be challenging to discriminate between these two states. This is a result of the  
66 combination of suboptimal detection tests, so that the actual extent of the pathogen  
67 prevalence is often unknown, as well as of the low genetic diversity, which can hide the  
68 temporal signal provided by the accumulation of mutations in the bacteria DNA.

69

70 In recent years, the increased availability, efficiency and reliability of genomic reading  
71 techniques, such as whole-genome sequencing (WGS), has significantly increased the  
72 amount of information we can use to study infectious diseases, and therefore it has  
73 improved the precision of epidemiological inferences for pathogens like *M. bovis*.

74 In this study, we use WGS to gain insights into the epidemiology of *M. bovis* in Cameroon, a  
75 developing country where the pathogen has been reported for decades. Ninety-one high-  
76 quality sequences were obtained from tissue samples collected in four abattoirs, 64 of  
77 which with complete metadata. We combined these with environmental, demographic,  
78 ecological and cattle movement data to generate inferences using phylodynamic models.

79

80 Our findings suggest *M. bovis* in Cameroon is slowly expanding its epidemiological range  
81 over time, therefore endemic stability is unlikely. This suggests that animal movement plays  
82 an important role in transmission. The simultaneous prevalence of *M. bovis* in co-located  
83 cattle and humans highlights the risk of such transmission being zoonotic. Therefore, using  
84 genomic tools as part of surveillance would vastly improve our understanding of disease  
85 ecology and control strategies.

86

87 **Keywords**

88 *Mycobacterium bovis*; whole-genome sequencing; phylodynamic analysis; genomic  
89 surveillance; livestock epidemics; zoonotic tuberculosis; One-health

90

## 91 **1. Introduction**

92 In the last two decades, the increased availability, efficiency and reliability of genomic  
93 reading techniques, such as whole-genome sequencing (WGS) techniques, have ignited a  
94 profound transformation in understanding disease ecology and epidemiology. This, coupled  
95 with improved statistical methodologies and high-performance computing, has enhanced  
96 our understanding of pathogen dynamics and evolution (1).

97 Techniques such as WGS can identify polymorphisms in the genetic material, which is  
98 generated by transcription errors that can occur to the pathogen while replicating within  
99 their host (2). As the pathogen is transmitted through the host population, the  
100 accumulation of polymorphisms in its DNA/RNA can be used as a “transmission signature”.  
101 Therefore, by tracking these mutations across bacterial genomes sampled in a host  
102 population, we are now able to infer transmission events between individual hosts, sub-  
103 populations, geographical areas or species, while at the same time gather insights about the  
104 evolutionary trajectory of a pathogen (2). Furthermore, when accurate spatial information  
105 on the sampled isolates is available, we can combine it with pathogen genetic data to  
106 disentangle the spatio-temporal dynamics of outbreaks, particularly in natural or other  
107 scarcely sampled animal populations (3).

108 Despite these advances, many challenges still exist, including the reconciliation between the  
109 temporal signal of outbreaks with pathogen mutations (4). *Mycobacterium tuberculosis*  
110 Complex (MTBC) members are clonal species, and therefore recombination has been  
111 considered rare (although a recent publication showed otherwise (5)). A few mutations are  
112 expected to occur for these species per year, generating little diversity during outbreaks in  
113 host populations. Consequently, there is inherent uncertainty in establishing infection  
114 patterns within the infected population and their associated infections. Therefore,  
115 combining genomic information with metadata is essential for accurate transmission chain  
116 estimation (6).

117 *Mycobacterium bovis*, a member of the MTBC group, is the aetiological agent of animal or  
118 bovine tuberculosis (bTB) in bovids and other mammals and of zoonotic tuberculosis (TB)  
119 in humans (7). Its infections are characterised by chronic disease, with or without a latent  
120 period, where infected cattle are hard to identify, making it hard to quantify potential  
121 infectious contacts (8). The estimation of *M. bovis* prevalence is often affected by several  
122 factors, including the inaccuracy of diagnostic tests (9), and the potential co-infection with

123 other pathogens (10). Such challenges explain why *M. bovis* has only been successfully  
124 eliminated or controlled in a few countries. Yet, it still represents a significant threat to  
125 cattle industries and human health in many other countries. For example, zoonotic  
126 tuberculosis due to *M. bovis* is a major public health problem in low and medium-income  
127 countries (LMICs), where close interaction between people and livestock is common and the  
128 limited access to pasteurized milk (7,11). Indeed, the magnitude of this burden is likely  
129 underestimated since human-animal transmission is predominantly via ingestion of infected  
130 products and presenting with a range of non-specific symptoms (12).

131 In Cameroon *M. bovis* is circulating in the cattle population, both in the southern areas (13)  
132 and, in particular, in the northern regions, where a previous study on cattle sampled at four  
133 regional abattoirs showed a sampled population prevalence ranging from 2.75% (31 positive  
134 over 1'129 cattle inspected, Northwest) to 21.25% (34 over 160, North (14)). Abattoirs  
135 surveillance, where carcasses are inspected for TB-like lesions, is the only surveillance  
136 strategy regularly implemented in the country; in Bamenda (Northwest region), Awah-  
137 Ndukum and colleagues (15) showed that the TB-like lesion in cattle increased in the period  
138 from 1994 to 2010.

139 Commonly to many LMICs, bTB control in Cameroon is also made difficult by the absence of  
140 detailed records on cattle population, by local rearing practices such as pastoralism which  
141 expose animals to contacts with other herds and potential reservoir wildlife species, and by  
142 the transhumance cattle movements westward towards Nigeria, where the demand of meat  
143 is driven by a fast human population increase (16).

144 In a previous study, Egbe et al. (16) employed two molecular typing techniques to  
145 understand the relatedness of *M. bovis* strains circulating in the region. These are  
146 spoligotyping and MIRU-VNTR typing: the former is based on the presence of multiple  
147 spacer oligonucleotides in the genome Direct Repeat region, while the latter is based on 12  
148 loci containing variable numbers of tandem repeats of mycobacterial interspersed repetitive  
149 units (17,18). Compared to WGS, these techniques consider a limited genome region and  
150 can be more subject to homoplasmy (19). The results reported by Egbe (16) showed that most  
151 of the isolates belonged to the Af1 clonal complex ( $n = 250$ /total  $n = 255$ ), while the  
152 remaining ones had an unidentified clonal complex. They also highlighted an unexpectedly  
153 high genetic diversity, as showed by the 37 sampled spoligotypes, of which 19 newly

154 observed, and a total of 97 genotypes, obtained by combining spoligotypes with MIRU-  
155 VNTR (16).

156 While those techniques are instrumental to investigating potential infection clusters at a  
157 broader level, they can be limited for a more in-depth understanding of the spatio-temporal  
158 dynamics of the disease. This study aimed to fill these gaps and enhance our understanding  
159 of the *M. bovis* epidemiology and spatial dynamics in Cameroon using WGS. We applied  
160 novel phylogenetic techniques to determine whether there was endemic stability across  
161 Cameroon's cattle-rearing regions while examining the role of environmental and ecological  
162 variables and animal movements in the pathogen spread.

163 We used 91 high-quality *M. bovis* sequences obtained from cattle's tissues sampled at  
164 regional abattoirs as described by Egbe et al. (14). After determining the single nucleotide  
165 polymorphisms (SNPs), we built a tree by joining the Cameroonian WGSs with other African  
166 sequences obtained from publicly available repositories, in order to understand how the  
167 sampled population fit in the continent context. Then, we ran a continuous space  
168 phylogeographical analysis with *BEAST* (20) on the Cameroonian sequences while testing  
169 different random walk diffusion models (21). This was possible because the origin village of  
170 the cattle tested at the abattoir was known for 64 *M. bovis* cattle isolates, allowing us to  
171 associate spatial coordinates to these sequences. Further, we tested the association  
172 between the spatial pathogen distribution obtained with the georeferenced phylogenetic  
173 tree and environmental, anthropic and ecological factors (22), and we finally ran a machine  
174 learning analysis to test whether the empirical cattle movement network (23) or other  
175 variables could explain the genetic diversity across isolates.

176 Our findings strengthen the call for an improved *M. bovis* molecular surveillance in  
177 underrepresented regions and countries, so to gather insights on potential patterns that can  
178 be missed when limiting the studies to areas of low genetic diversity, consequence of strict  
179 control practices such as test-and-cull.

180

181 **2. Materials and methods**

182 **2.1. Data collection**

183 Four regional abattoirs were sampled between 2012 and 2013, in the Northwest  
184 (Bamenda), Adamawa (Ngaoundere), North (Garoua) and Extreme North (Maroua) regions  
185 of Cameroon (Figure S1). As part of the regular operations, cattle carcasses were inspected  
186 for the presence of TB-like lesions. The tissues, including lymph nodes, of all animals with  
187 lesions and of some randomly chosen without lesions were collected to be cultured, and  
188 information about the animal (age, breed, village of provenance, among others) were taken.  
189 A detailed description of the data collection and bacterial isolation can be found in Egbe et  
190 al. (14). The DNA extraction was conducted in BSL 3 facilities (Tuberculosis Reference  
191 Laboratories in Bamenda, Cameroon), and the procedure is fully described in Egbe et al.  
192 (16). Sequencing was also attempted for *M. bovis* isolates sampled in human hosts at the  
193 Bamenda hospital (Northwest region) during a cross-sectional study within the wider  
194 project. We reported a summary of the number of sampled animals and the number of *M.*  
195 *bovis* positive ones in Table S1.

196

197 **2.2. Whole genome sequencing processing**

198 The sequencing was carried out at Edinburgh Genomic facilities (University of Edinburgh).  
199 Samples were prepared with 1 TruSeq Nano 550 bp insert, 76 Pippin selected library from  
200 the supplied genomic DNA, while MiSeq v2 (Illumina) was used to generate 250 base paired-  
201 end sequence from library to yield at least 11M+11M reads (1 run) at 30x coverage. The  
202 output was read from a 4 lane Miseq. A total of 124 *M. bovis* WGSs were obtained (two  
203 from human hosts), while for nine isolates (one from human) the sequencing failed.  
204 We used adapted *BovTB*-nf pipeline (24) for quality control. Reads were deduplicated using  
205 *fastuniq*, trimmed using *Trimmomatics* (25) (-phred33 ILLUMINACLIP:\$adapters:2:30:10  
206 SLIDINGWINDOW:10:20 MINLEN:36), and mapped to the reference genome using *bwa*-  
207 *mem2* (26). The mapped reads were filtered (-Shuf 2308 -) and sorted using *Samtools* (27),  
208 and then classified using *Kraken2* (28) (--quick) against a prebuilt Kraken 2 database  
209 (*Minikraken* v2 (28)). The Kraken2 output was summarized with Braken (29) (-r 150 -l S), and  
210 the top 20 list of species from the Bracken output was used to determine if the sample was  
211 contaminated with other microorganisms. Variants were called using *bcftools* (30) (--  
212 IndelGap 5 -e 'DP<5 && AF<0.8') and strain-specific SNPs were used for classifying whether



213 the samples were *M. bovis* or not (custom script and differentiating SNPs taken from (24)).  
214 The percentage of coverage (>60%) on the reference, read depth (>10), and number of  
215 reads (> 600,000) were used to identify and remove samples with insufficient data. To  
216 curate aligned core-variants for the downstream phylogenetic analysis, variants were called  
217 and filtered using *Snippy* v4.6.0 (31) using the default settings (minimum coverage = 10,  
218 minimum VCF variant call quality = 100), with the *M. bovis* AF2122/97 genome (GeneBank:  
219 LT708304.1) as the reference genome. Variants from repeated regions were removed (mask  
220 for repetitive regions taken from (24)). Core-SNPs were determined by *snippy-core* function  
221 within *Snippy*, where a genomic position was considered to be a core-site when present in  
222 all samples. We defined as “high-quality” sequences the ones with genome coverage > 90%  
223 and reading depth > 10 (32), and we renamed the sequences with a string composed by the  
224 following information: host species, location (administrative subdivision, or country, see  
225 Section 2.3), sequential number, and date.

226 For each sequence, the spoligotype and the clonal complex were retrieved from Egbe et al.  
227 (16). In one case a sequence was missing the spoligotype number, however, it was assessed  
228 with the *vSNP* pipeline (33). For all bioinformatics tools we used the default settings, unless  
229 stated otherwise.

230 We checked if divergent sequences belonged to other mycobacteria species. We tested the  
231 presence of RD (regions of difference) 1, 4, 9 and 12 patterns (34) in the outlier samples,  
232 raw reads from each sample were aligned to *M. tuberculosis* (NC\_000962.3) with *Burrows-*  
233 *Wheeler Aligner* v0.7.17 (35), and sorted and indexed with *SAMtools* v1.10 (36). Primer  
234 flanking regions for the RDs on *M. tuberculosis* were determined through querying the  
235 sequences using NCBI web nucleotide BLAST with the default parameters (37), while the  
236 presence of RDs were manually determined by examining the read alignment in *Integrative*  
237 *Genomics Viewer* v2.14.1 (38).

238

### 239 **2.3. Cameroonian *M. bovis* sequences in the African context**

240 We obtained other *M. bovis* genomes from online repositories: first, from the *Patric* (now  
241 BV-BRC) dataset (39), and second, selecting the appropriate genomes among the ones listed  
242 by Loiseau et al. (40) and obtained from the EBI dataset (for details and references, see  
243 Table S2). We selected all the available sequences sampled in Africa, in order to qualitatively  
244 detect potential genetic similarities between the sampled Cameroonian *M. bovis* population

245 and other isolates from the African continent, and thus provide a broader context to our  
246 analysis.

247 When analysing sequences from *Patric*, genomes were shredded into pseudo by *Snippy*  
248 followed by the process of alignment and SNP identification described above. The core-SNP  
249 alignments were made with and without the other African genomes. We used *iqtree* web  
250 server (41,42) to compute a phylogenetic tree ( $n = 212$ ) which included all the Cameroonian  
251 high-quality sequences ( $n = 91$ ) and the other African ones plus the 1997 reference from UK  
252 ( $n = 121$ ).

253

#### 254 **2.4. Cameroonian sequences phylogenetic analysis**

255 The quantitative analyses were performed on a subsample of the Cameroonian sequences,  
256 obtained after removing the non-cattle ones, the ones missing the geographical  
257 coordinates, and potential outliers, i.e., isolates not clustering within the main Cameroonian  
258 clade. We initially joined the remaining sequences ( $n = 64$ ) tree using the TN93 genetic  
259 distance model and Neighbour-Joining (NJ) algorithm *ape* package (43) in *R* v4.0.5 (44)  
260 with the sole purpose of estimating a temporal signal within the sample in *TempEst* v1.5.3  
261 (45). We then used the sequences SNP alignment completed with sampling dates, to infer  
262 time-scaled phylogenetic trees using *BEAST* v1.10.4 (20) with the *BEAGLE* library (46), and  
263 evaluated the results with *Tracer* v1.7.2 (47). Since the sequences had associated  
264 geographical location metadata, we included latitudes and longitudes as an additional  
265 continuous space variable for phylogeographic inference.

266 To select the best model, we ran a series of exploratory models using a HKY (48)  
267 substitution model, similar to other studies (49–51), and a strict molecular clock. We  
268 sequentially selected the best continuous trait model first, then the best bacterial  
269 population size model (tree prior). We tested the Brownian random walk, Cauchy Relaxed  
270 Random Walk (RRW), lognormal RRW and Gamma RRW for the former, and constant  
271 population, exponential growth and Bayesian *Skygrid* (52,53) for the latter. In the  
272 exploratory *BEAST* runs, we chose a truncated (between 0 and 0.1) normally distributed  
273 clock rate prior, with mean and standard deviation set as the slope in the root-to-tip  
274 obtained in *Tempest*; the chain length was set to  $10^8$ , sampled every  $10^4$  steps. The models  
275 were compared using marginal likelihood estimation (MLE), with path sampling (PS) and  
276 stepping-stone sampling (SS), if they reached a satisfactory effective sample size ( $>200$ ).

277 Once the model features were selected, we ran a final one setting the chain length to  $10^9$   
278 steps, sampled every  $10^5$  steps. In this case, we used the clock rate posterior of the selected  
279 exploratory model as a prior for the final model. The maximum clade credibility (MCC) tree  
280 was extracted with *TreeAnnotator* v1.10.4 (part of the *BEAST* suite), and clades were visually  
281 defined within the MCC tree branches. The MCC tree was plotted against the sequences  
282 spoligotype and MIRU-VNTR typing to visually assess the correspondence between  
283 molecular typing and clades.

284

### 285 **2.5. Spatial statistics and environmental factors analysis**

286 From the final *BEAST* run, we extracted a set of 100 trees from the posterior distribution  
287 and further analysed using *seraphim* v1.0 (22,54) to obtain the spatial spread statistics:  
288 branch velocity and epidemic wavefront. The former was calculated for each branch dividing  
289 the geographical distance from the origin to the destination nodes by the time branch time  
290 duration. The epidemic wavefront shows the geographical range of the epidemic over time:  
291 at each time it is calculated as the geographical distance between the positions of the tree  
292 estimated root and the most distant node (spatial distance wavefront), or accounting for  
293 the distance of nodes closer to the root (patristic distance wavefront).

294 Additionally, *seraphim* allows to statistically test hypothesis on the effect of environmental  
295 layers on the epidemic dynamics; the effect can either be of “conductance”, when the layer  
296 favours the pathogen diffusion, or “resistance”, when it hampers it. We tested nine layers:  
297 elevation, cattle population density, human population density, two describing the roads  
298 infrastructure (number of intersections and total road length), and four land cover types  
299 (waterbodies, forest, grassland and grazeland, and other vegetation types: mosaic, shrub,  
300 sparse vegetation). The original raster layers were downloaded from online repositories (see  
301 Table S3 for the sources) and adapted to a 5km x 5km grid using *QGIS* v3.26.1. For each cell,  
302 elevation, cattle and human populations were averaged for the 5x5km grids, while roads  
303 intersections were counted, and roads length were measured starting from the same road  
304 original raster. For land cover, each value represents the percentage of that cell covered by  
305 each land cover type. The original land cover raster included 38 different cover types. To  
306 ease computation, we selected the most relevant for the study and merged them in four  
307 layers: waterbodies, forest, cropland/grassland, and other vegetation, including mosaic,  
308 shrub, and partial cover (Table S4).

309 First, we ran a preliminary analysis on each variable, to determine if it could have played a  
310 role as conductance or resistance in the pathogen spread. For each of the 100 extracted  
311 trees, we estimated the correlation between dispersal duration and environmental distance.  
312 Results are summarised by two statistics: the number of positive variable's coefficient of  
313 determination out of the 100 trees, and the number of positive  $Q$  statistic, calculated as  
314  $Q = R_{var}^2 - R_{null}^2$ , that is the difference between the correlation  $R^2$  for the variable's raster  
315 and for a null raster, again calculated for each tree (54). For the analysis, we used two path  
316 models: straight line (where the branch "weight" is calculated as the by summing the cells  
317 values through which the straight-line passes), and least cost path (where the branch  
318 "weight" is calculated by summing the values between adjacent cells along the least-cost  
319 path).  
320 Once we identified the potential resistance or conductance factor, we performed ten tree  
321 randomisations and calculated the statistics again. In this case, we used the Bayes Factor  
322 ( $BF_e$ ), calculated as  $BF_e = p_e / (1 - p_e)$ , where  $p_e$  is the probability that  $Q_{observed} > Q_{randomised}$ . We  
323 used two criteria for trees randomisations: 1) randomisations of nodes positions while  
324 maintaining the branches lengths, the tree topology and the location of the most ancestral  
325 node; and 2) randomisations of nodes positions while maintaining only the branches  
326 lengths.

327

## 328 **2.6. Genetic distance regression and role of the cattle movements**

329 We finally tested which variables can better explain the genetic distances between the  
330 sampled *M. bovis* isolates, so to understand the signatures of temporal, spatial, and  
331 demographic factors (56,57). We ran this analysis using a Boosted Regression Trees (BRT)  
332 regression model (58) in *R* (packages *dismo*(59) and *gbm*(60)), a very flexible tool which  
333 combines decision trees and boosting techniques (61). In this model, the dependent  
334 variable was the genetic distance between *M. bovis* strains, expressed as SNPs. We tested a  
335 total of 28 relational variables, calculated for each pair of isolates (Table S5). Except for the  
336 temporal and spatial distance (which were calculated from the original isolates metadata),  
337 and for a binary variable indicating whether two sequences have the same spoligotype,  
338 MIRU-VNTR and clade (yes/no), the other variables are associated to the *M. bovis* isolates  
339 administrative subdivision.

340 We built two subdivision-level contact networks. The first one is a spatial network where  
341 nodes represent subdivisions and edges between them are positive if they share a border.  
342 This network is undirected (edges are not directional) and unweighted (all edges values are  
343 set to one). For this network we computed six variables to be associated with each pair of  
344 isolates: degree and betweenness centrality (62) of both isolates' subdivisions; shortest path  
345 and a binary variable indicating whether the two subdivisions belonged to the same  
346 network's community.

347 The second network represented the cattle movements, and edges correspond to the  
348 number of animals moved between subdivisions over a year. We built this network by  
349 aggregating the empirical data collected by Motta et al. (23), which originally reported the  
350 monthly number of cattle exchanged between markets. For this network we computed  
351 eight variables: degree, strength and betweenness centrality of both isolates' subdivisions;  
352 shortest path and the same community binary variable. The degree counts the number of  
353 each subdivision's connections, while the strength is the sum of the number of cattle moved  
354 to and from each subdivision. All networks' metrics were computed using the *R* package  
355 *igraph* (63).

356 Once we computed all the variables (the full list is reported in Table S5), we trained the BRT  
357 model using 75% of the observations, while the remaining 25% were used for testing. We  
358 evaluated the models based on pseudo- $R^2$  and Root Mean Squared Error (RMSE) on the test  
359 dataset. These were both calculated using the package *caret* (64). For BRT the relative  
360 influence of the variables is determined by the times each variable is selected to split the  
361 data in a decision tree, which in turn is weighted by the improvement in the model fit that  
362 resulted from that variable being used at each split (58). All models were fitted with a 10-  
363 fold cross validation. The BRT algorithm has two main parameters: the learning rate, which  
364 controls the contribution of each tree to the final model, and the tree complexity, which  
365 corresponds to the number of nodes in the tree. We ran some preliminary tests to tune the  
366 BRT in order to improve the predictions. Finally, we set the learning rate to 0.05 and the  
367 tree complexity to 8.

368

### 369 3. Results

#### 370 3.1. Cameroonian sequences in the African context

371 We analysed 124 *M. bovis* sequences (nine of the original 133 failed), with 91 having enough  
372 read depth and genome coverage to allow further analyses (see Table S6 for further details).  
373 Two of these sequences came from isolates sampled humans, while for a third the  
374 sequencing failed. One of the excluded sequences was marked as not-*M. bovis*, and based  
375 on the presence of the four RD1, 4, 9 and 12 patterns (34), it was likely *M. tuberculosis*. All  
376 the high-quality *M. bovis* sequences were merged in a tree with other 22 obtained from the  
377 *Patric* dataset, 99 from EBI, and the 1997 UK Reference to provide a continental context.  
378 The qualitative phylogenetic tree in Figure 1 shows that most of the Cameroonian  
379 sequences (two of which obtained from human tissue samples) cluster with the Ghanaian  
380 human samples, and two Nigerians ones recovered from unreported hosts. All human  
381 samples from West Africa cluster with cattle sequences except for the Malian human  
382 sequence. Most sequences ( $n = 89$ ) belonged to Af1 clonal complex and except one, the  
383 spoligotypes were already known; for the other, we identified a new pattern (hex code: 6F-  
384 1F-5F-7F-BF-40). Being characterised by the absence of spacer 30, this spoligotype was  
385 considered as Af1 (65). The dominant spoligotype was SB0944 ( $n = 32/89$ ).  
386 Two outlier sequences did not cluster with the rest of the sampled Cameroonian population.  
387 Their average distance from the rest of the Cameroonian population (respectively 235 and  
388 231 SNPs) was slightly higher than the average distance of the 1997 UK reference from the  
389 Cameroonian isolates (222 SNPs), and they did not cluster with any other WGS sequence  
390 sampled in Africa (Figure 1). For both outlier sequences, the spoligotype was SB2332, found  
391 for the first time in Cameroon and submitted for classification at [www.Mbovis.org](http://www.Mbovis.org) by Egbe  
392 et al. (16). Following Warren et al. (34), we tested the presence of RD1, 4, 9 and 12 patterns,  
393 finding only the first one, confirming that they are likely *M. bovis*. We compared this  
394 spoligotype pattern with all the others from the [www.Mbovis.org](http://www.Mbovis.org) database, and we  
395 identified four patterns differing by two spacers: SB0858 sampled in Spain (66) (different  
396 spacers 20 and 22), SB1102 sampled in Chad (65) and Cameroon (13) (different spacers 33  
397 and 34), SB2333 reported by Egbe et al. (16) (different spacers 22 and 34) and SB2691  
398 sampled in France (not found in publications, different in spacers 20 and 34). We also  
399 identified eleven patterns differing by three spacers, sampled in France (67), Portugal (68),  
400 and Spain (66).

401

### 402 **3.2. *M. bovis* evolutionary time scale in Cameroon**

403 A total of 1'540 SNPs were determined from the *Snippy* core-SNP analysis on Cameroonian  
404 *M. bovis* genomes (Figure S2). This reduced to 1'106 SNPs when the dataset was reduced to  
405 the 64 samples with complete metadata and excluding the non-cattle ones (two sampled in  
406 humans), which were used for the downstream quantitative analysis. The median SNP  
407 distance among the remaining high-quality sequences was 70 SNPs (mean 68, range from 0  
408 to 144, 2.5<sup>th</sup> and 97.5<sup>th</sup> quantiles 14 and 118). For two cattle (one from Bibemi, the other  
409 from Touboro), two *M. bovis* isolates sequenced were available (obtained from different  
410 tissues). In both cases, the two strains were identical (Bibemi 3 and 4, Touboro 7 and 8,  
411 Figure 2), which suggests a single infection disseminated in different organs, rather than two  
412 separate infections.

413 The analysis in *Tempest* showed a slightly positive temporal signal (coefficient of  
414 determination 0.11, and correlation coefficient 0.33) and a slope of  $1.267 \times 10^{-2}$  (Figure S3).  
415 We used a sequential approach in *BEAST* to select the best spatial model and bacterial  
416 population models. Based on the MLE estimation of the exploratory models (Table S7) we  
417 determined the best model included a Gamma Relaxed Random Walk (RRW) spatial model  
418 (first step of the sequential analysis) and the SkygGrid population model (second step). The  
419 final *BEAST* model was run with 10 bins and a cut-off of 400 years. The population trend is  
420 shown in Figure S4. The model estimates suggest the mean age of the root was in July 1950  
421 (95<sup>th</sup> high-posterior density, HPD, April 1938 – August 1961), while the average clock rate  
422 was  $1.32 \times 10^{-7}$  substitution/site/year (95<sup>th</sup> HPD  $1.20 \times 10^{-7} - 1.44 \times 10^{-7}$ ). The maximum  
423 clade credibility (MCC) tree is reported in Figure 2, which also shows the division in four  
424 clades: clade 1 (green, 22 isolates), clade 2 (blue, 17 isolates), clade 3 (purple, 19 isolates)  
425 and clade 4 (red, 5 isolates). One sequence was excluded from all clades (Belel 4, Figure 2,  
426 reported as “no clade” in the figures). The geographical distribution of the clades is reported  
427 in Figure 3, showing the number of *M. bovis* isolates per administrative subdivision, which  
428 ranged from 1 to 17 (see Table S8 for the number of isolates per clade by regional abattoir).  
429 In Figure 4, we superimposed the MCC tree with spoligotypes; the most prevalent  
430 spoligotype, SB0944, occurred 26 times (out of 64 sequences) and was present in three of  
431 the four clades. The second most prevalent spoligotypes were SB0953 and SB2312, the first  
432 occurring five times in two clades, the latter occurring five times in one clade only (clade 2).

433 We also superimposed the MIRU-VNTR types as shown in Figure S5. The most prevalent  
434 MIRU-VNTR type in the sampled population was V89, which occurred nine times; V82 and  
435 V37 respectively occurred six and four times; and V81, V76 and V100 all occurred three  
436 times. Seven MIRU-VNTR types occurred twice, while 39 types occurred only once.

437

### 438 **3.3. Spatio-temporal pathogen expansion**

439 The estimated mean branch velocity was 53.1 km/year (95<sup>th</sup> CI 18.4 – 219.0, temporal trend  
440 reported in Figure S7). The wavefront statistics in Figure 5 suggests that the pathogen  
441 expansion was slow until the mid 1960s, but accelerated thereafter to reach the entire  
442 study area, with a slow but constant expansion in the following period. This is reflected in an  
443 increase of the branch velocity at the same time (Figure S7), which is approximately the  
444 period when the branches formed the observed clades (Figure 2). The timing of the different  
445 branches in space is reported in Figure 6 (95<sup>th</sup> HPD in Figure S8, with nodes coloured by  
446 estimated/observed date).

447 We tested the association between nine geographical variables with the dispersal duration.  
448 Table 1 shows the results obtained using the straight line and the least cost path models,  
449 the latter run considering the variables as potential conductance or resistance factor. Six  
450 variables resulted in a significant association (positive coefficients for all at least 95 out of  
451 100 trees, and above 75% of positive Q): mosaic, shrub and other vegetation cover (with  
452 both path models, as resistance in the least cost one); forest cover, elevation and  
453 waterbodies cover (all as conductance); and cattle density (as resistance). However, when  
454 their statistical significance was tested through the randomisation, only forest cover and  
455 elevation (both as conductance) showed a Bayes Factor significant ( $\geq 3$  (69)). The result was  
456 robust against two different trees randomisation algorithms for the forest layer, while for  
457 the elevation this was true only when maintaining only the branches length and excluding  
458 the other tree topological characteristics.

459

### 460 **3.4. Factors associated with genetic distance**

461 The RMSE of the boosted regression trees BRT model ran using all 28 variables was 20.23,  
462 while the  $R^2$  was 0.450. We simplified the model using the *dismo* package, which tests the  
463 performance of the model by dropping the less important variables with a procedure similar  
464 to backward selection in regression (58). The algorithm brought to eliminating 12 variables



465 (see Table S5), nonetheless the model run using the remaining 16 variables performed very  
466 similarly to the original one (RMSE = 20.22 and  $R^2 = 0.452$ ). Therefore, we used the latter to  
467 calculate the variable importance (Figure 7).

468 As expected, the most relevant variables were the temporal distance between the samples  
469 (1<sup>st</sup>) and the binary variable indicating whether the two *M. bovis* isolates belonged to the  
470 same clade in the MCC tree (2<sup>nd</sup>). The variables describing the subdivisions' population were  
471 also relevant in the model (population.y, 3<sup>rd</sup>, and population.x, 5<sup>th</sup>), as well as whether two  
472 isolates shared the same MIRU-VNTR (4<sup>th</sup>). This was more relevant than if two isolates  
473 shared the same spoligotypes (11<sup>th</sup>), suggesting the former as more useful to discriminate  
474 closer *M. bovis* strains. The markets movement network strength (i.e. the number of cattle  
475 moved from/to a subdivision) was the most important (6<sup>th</sup> and 9<sup>th</sup>) among network-related  
476 variables, while the betweenness (8<sup>th</sup> and 10<sup>th</sup>) was the only spatial network variable  
477 retained in the simplified model. Interestingly, when both variables were selected for the  
478 same metric, the one related to the youngest isolate (marked by y) was always preferred to  
479 the one related to the oldest isolate (marked by x). The partial dependency plots, showing  
480 the relationship between SNP distance and variables, are reported in Figure S9.

481

## 482 **4. Discussion**

483 We sought to unravel the characteristics of the spread of a pathogen with zoonotic  
484 potential in time and space to improve our understanding and inform control and  
485 preparedness strategies. Our basic premise is that the accumulation of mutations in the  
486 pathogen's genome can be used as signatures of transmission events from host to host  
487 across time and space. Within space, the environment can create barriers which influence  
488 the population dynamics of diseases, i.e., altering host-to-host and pathogen-host  
489 interactions has direct effects on the genetic structure of the pathogen (70). The availability  
490 of high-throughput genomic techniques means we can interrogate the structural changes  
491 linked to the environment over time to gain critical insights into how the epidemic has  
492 evolved. In this study, we aimed to characterise *M. bovis* sampled from cattle in Cameroon  
493 using genetic and demographic data to understand whether the pathogen is in a stable  
494 endemic state and the influence on the spread dynamic of environmental and ecological  
495 factors and cattle movements.

496

### 497 **4.1. Evidence of dynamic endemicity**

498 An important question was whether the *M. bovis* outbreak in North Cameroon was in a  
499 steady state, at an endemic equilibrium, or if it was expanding. Determining whether a  
500 pathogen is endemic has implications on risk perception and, consequently, on resource  
501 allocation. At the same time, the chances of zoonotic transmission are likely to be higher in  
502 the case of endemicity. In our analysis, the Bayesian model estimation with SkyGrid as a  
503 population model showed an increasing pathogen effective population size, corresponding  
504 to a constant increase in the disease velocity after the sudden jump during the mid-to-late  
505 1960s. This suggests that the pathogen is not in a state of endemic stability, instead it has  
506 been expanding at various rates over the years. This is in agreement with a previous  
507 publication using spoligotypes and MIRU-VNTR (16) and with the work by Awah-Ndukum et  
508 al. (15). The expansion of *M. bovis* might represent an issue for livestock and humans,  
509 particularly as we showed that the bacterium is circulating in both. At the moment, disease  
510 control in the area is absent, while on the other hand, the dairy industry in Africa is  
511 generally expanding. A lack of widespread milk pasteurization could lead to an increase in  
512 zoonotic TB cases, which already represent a problematic issue in the region (12).

513

#### 514 **4.2. Genetic diversity of *M. bovis* in Cameroon**

515 We observed a high diversity of *M. bovis*, confirming earlier observations with molecular  
516 typing techniques providing less granular information (16), considering the short time span  
517 of the sampling campaign and the small sample size. This contrasts with areas such as Great  
518 Britain and other European countries, where strict control measures, such as routine testing  
519 and stamping out of positive individuals, have been in place for decades. This can act as a  
520 bottleneck with a consequent reduction in the pathogen's genetic variability by reducing the  
521 time a pathogen has to develop inside a domestic host, therefore, the likelihood of  
522 substitutions in the DNA. As an example, Crispell et al. (57) reported a similar SNP distance  
523 range (0 to 150), albeit across a much bigger sample ( $n = 230$ ), with a lower median (20  
524 SNPs) and with isolates dating back two decades, while in a similar size monophyletic  
525 outbreak ( $n = 64$ ), Rossi et al. reported a maximum SNP distance of only 6 SNPs (56). In  
526 Spain, Pozo et al. (71) found a similar SNP distance average and range (62, and 0 to 150) in a  
527 bigger *M. bovis* population, sampled in both cattle and wildlife over 13 years. It is  
528 noteworthy that high diversity can be associated with dynamic epidemiology and not with  
529 endemic stability.

530 All 64 core isolates belonged to the clonal complex Af1, which was observed in the region in  
531 previous studies (65). The most common spoligotype, SB0944, was found by Müller et al.  
532 (65) as the most prevalent in West Africa and considered as the original of the Af1 clonal  
533 complex. Our findings also suggest zoonotic transmission in West Africa, as sequences  
534 recovered from humans in Cameroon and Ghana clustered with Cameroonian cattle *M.*  
535 *bovis* isolates (72). Because it is known that zoonotic TB represents a minoritarian but still  
536 crucial part of all TB cases in Africa, these results strengthen the case for One Health  
537 approaches to control, that involve humans, livestock, wildlife and environmental health  
538 (12,73). Except for the one sequence in Mali and the two Cameroonian outliers, all the  
539 sequences from West Africa clustered together, hinting to a high connectivity likely caused  
540 by cattle movements throughout the area, as previously showed by another study (74). Our  
541 results showed that the areas with the highest *M. bovis* diversity were in the Adamawa and  
542 North regions, both reporting all the clades identified by the maximum clade credibility  
543 (MCC) tree. All clades were also sampled in the towns of Touboro and Tchollire, both  
544 located in the North region but close to the Adamawa border. Previous studies reported  
545 that this area receive cattle from neighbouring country as part of the transhumance

546 migration, suggesting that cattle movements and markets play an important role in defining  
547 the dynamics of the pathogen, and therefore influencing its genetic diversity (16,23). The  
548 Northwest region was underrepresented in the sample, with only five high-quality  
549 sequences on 31 infected cattle detected at the abattoir. This inherently reduces the level of  
550 diversity, which is far lower than reported using spoligotypes and MIRU-VNTR (16).  
551 Despite covering a smaller portion of the genome and the higher occurrence of homoplasmy  
552 with respect to WGS, in other contexts spoligotypes have been used as a proxy cluster, or to  
553 narrow down potential transmission within the study population (57,75). Our results  
554 showed that this cannot be done for areas with high diversity such as the one we  
555 considered, as we observed little correspondence between the MCC tree branches and the  
556 spoligotypes. Similarly, other studies pointed out the limitations of such typing techniques  
557 (19), in case of an expanding infection where transmission is steadily ongoing, compared to  
558 point-source ones (76). The high SNP distances among the sampled isolates also precluded  
559 the use of methods to infer direct transmission between hosts (8,77).

560 When considering the entire sampled population, therefore including the sequences with  
561 incomplete metadata, we found two of the 91 sequences not belonging to the clonal  
562 complex Af1. In their spoligotype pattern (SB2332), we noted the absence of spacer 21 (78),  
563 and the closest relatives analysed by Loiseau et al. (40) were identified as part of the clonal  
564 complex Eu2, including isolates sampled both in South-western Europe (SB0837, SB1090,  
565 SB1308) and West Africa (SB1102, isolated in Cameroon as well (13)). We can then  
566 speculate that these sequences likely belong to Eu2 as well, although we could not exclude  
567 one of the “unknown” clonal complexes identified by other studies (40,79). Further  
568 development on this point was beyond the scope of this study, as we focused on the 64 core  
569 sequences to gather insights on the pathogen dynamics in the area.

570

### 571 **4.3. Tracking the spread of *M. bovis* in Cameroon**

572 We acknowledge that our estimates for the most recent common ancestor (MRCA) have a  
573 wide credible interval around it (23 years). This uncertainty is likely due to the short  
574 duration of the samples collection campaign, which also generated a weak temporal signal,  
575 although the coefficient of determination was similar to other *M. bovis* studies in highly  
576 sampled populations (56,57). Nonetheless, our estimates coalesce around 1950, suggesting  
577 that the pathogen has been spreading in the area for at least six decades at the time of

578 sampling. For the same reason, the estimated clock rate was higher than others in the  
579 literature but in the same order of magnitude ( $0.67-1.26 \times 10^{-7}$ ,  $n = 2625$  (40)).

580 The estimated MCC tree located the most recent common ancestor (MRCA) in Touboro  
581 (North region) and, from there, a rapid expansion of the outbreak reaching most of the  
582 study area by the early 1970s. From the estimated origin, the pathogen likely spread first  
583 northward to Garoua (in the same region) and westward, to the North West region, and later  
584 to the Extreme North and Adamawa regions and again to the Northwest.

585 The results of the spatial factors analysis showed that forest cover and elevation were the  
586 only significant ones, both acting as “conductance”. Forest cover could be a proxy for  
587 potential wildlife interactions, as *M. bovis* is known to be quite effective in spreading at the  
588 wildlife-livestock (and humans) interface (73,80). The elevation as conductance was  
589 counter-intuitive, however, this could be linked to cattle movements in pastoralist  
590 communities within the plateau located in the study area. This is important because, if  
591 confirmed, altitude could be used as a proxy for the missing pastoralist movements.

592 Our regression model performed reasonably well, although the amount of variability  
593 explained was below 50%. However, our objective was to understand which variables could  
594 better explain the genetic distance between *M. bovis* isolates, expressed as SNP distance.  
595 Except for the between isolates temporal distance and clade, the demographic variables  
596 were the most effective in explaining SNP distance, particularly the administrative  
597 subdivision human population size. These variables had a negative effect on the SNP  
598 distance, meaning that smaller population was associated to a close relatedness of the *M.*  
599 *bovis* strains. This could be an effect of the population distribution in the country because  
600 the northern regions, where cattle are most concentrated, are less populated compared to  
601 the cities in the south. The simplified model performed similarly to the full model,  
602 suggesting some variables were not important in explaining the genetic distance. Beyond  
603 the human population size, also the other demographic variables (population and cattle  
604 density) were all retained. Conversely, only five network related variables were retained,  
605 three for the cattle movement network (out of eight) and two for the spatial network (out  
606 of six). All network related variables had a positive effect on the SNP distance, with the  
607 number of cattle moved in or out a subdivision (i.e., strength) having the higher predictive  
608 effect. Interestingly, this result was similar to other studies where cattle movements alone  
609 could not fully capture *M. bovis* genetic diversity (56,57).

610

#### 611 **4.4. Limitations**

612 The major limitation of this dataset was the short data collection time window, less than a  
613 year and a half, which resulted in uncertainty in the MRCA estimate and a weak temporal  
614 signal. While we can speculate the sampled bacterial population already reached the entire  
615 study area before the 1970s, a wider sampling time window would likely allow a stronger  
616 temporal signal and improve our estimate of the MRCA, which might be prior with respect  
617 to the current estimate. In turn, this affected the pathogen's expansion patterns, including  
618 the branch velocity and wavefront, which are also limited by the sampled area size. The  
619 spatial uncertainty might also be affected by the absence of dense cattle movements  
620 records, so the known spatial coordinates associated with each sequence correspond to the  
621 last village the animal lived in. The Adamawa and Northwest regions are home to 1.25  
622 million and 450.000 cattle respectively (23), and while this abattoir-based study provides a  
623 very informative snapshot of the *M. bovis* population in North Cameroon, it adds to the calls  
624 to improve cattle records and movements routine data collections in LMICs (81), as well as  
625 bTB detection efforts.

626 The low-quality WGSs disproportionately affected the Northwest region, as reported in Table  
627 S7. This could have hampered the representativeness of the *M. bovis* diversity in that  
628 region, reducing the number of clades observed. The Adamawa region was the most  
629 represented, despite most of the sequences excluded from the quantitative analysis  
630 because of missing coordinates, came from the Ngaoundere abattoir. The bacterium  
631 diversity in the Northwest might also be affected by the demographic of the slaughtered  
632 cattle in the region (14): because the region is highly populated by humans and more  
633 isolated in the trade network (23), local animals of both sexes and at any age are  
634 slaughtered. Conversely, young male calves from the Adamawa, North and Extreme North  
635 regions are often sent to richer southern regions to maximise their economic values, leaving  
636 the older cows to be slaughtered. By being exposed to the *M. bovis* for longer, the latter  
637 have more chances to develop lesions. On the other hand, these trends likely reduce the  
638 impact of missing information on the previous location of the animals, because these  
639 animals have more chances of being reared locally.

640 In agreement with many studies, and with the *vSNP* analysis result, we used AF2122/97 as  
641 reference genome (50,51,56,57,79,82,83). In order to account for genes, absent in *M. bovis*,

642 Loiseau et al. (40) used *M. tuberculosis* H37Rv, a choice driven by the different purpose of  
643 their work compared to ours (define the origin and the global population structure of *M.*  
644 *bovis*). Generally, the pipelines used to call the SNPs differed in many of the aforementioned  
645 studies, contributing to the estimates uncertainty and potentially generating biases the  
646 analysis results and the clock rate calculations.

647

## 648 **5. Conclusion**

649 In conclusion, our study indicates endemic stability of *M. bovis* is unlikely in North  
650 Cameroon, but rather the disease is slowly expanding over time. Our findings highlight the  
651 importance of collecting data in underrepresented areas to enrich insights in the current  
652 body of literature, predominantly from developed countries. Moreover, our results pave the  
653 way for future research aimed to understand whether the observed *M. bovis* high genetic  
654 diversity affects the spread dynamics.

655 Our findings underscore the need to adopt a one-health surveillance strategy for *M. bovis*  
656 control (12). More work on combining tools such as phylogeography, statistical modelling,  
657 landscape and ecology will be beneficial to map spread patterns and effectively inform  
658 control and preparedness strategies (56).

659

660 **References**

- 661 1. Pybus OG, Rambaut A. Evolutionary analysis of the dynamics of viral infectious  
662 disease. *Nat Rev Genet* (2009) 10:540–550. doi: 10.1038/nrg2583
- 663 2. Kao RR, Haydon DT, Lycett SJ, Murcia PR. Supersize me: How whole-genome  
664 sequencing and big data are transforming epidemiology. *Trends Microbiol* (2014)  
665 22:282–291. doi: 10.1016/j.tim.2014.02.011
- 666 3. Pybus OG, Suchard MA, Lemey P, Bernardin FJ, Rambaut A, Crawford FW, Gray RR,  
667 Arinaminpathy N, Stramer SL, Busch MP, et al. Unifying the spatial epidemiology and  
668 molecular evolution of emerging epidemics. *Proc Natl Acad Sci* (2012) 109:15066–  
669 15071. doi: 10.1073/pnas.1206598109
- 670 4. Campbell F, Strang C, Ferguson N, Cori A, Jombart T. When are pathogen genome  
671 sequences informative of transmission events? *PLoS Pathog* (2018) 14:1–21. doi:  
672 10.1371/journal.ppat.1006885
- 673 5. Patané JSL, Martins J, Castelão AB, Nishibe C, Montera L, Bigi F, Zumárraga MJ, Cataldi  
674 AA, Junior AF, Roxo E, et al. Patterns and Processes of *Mycobacterium bovis* Evolution  
675 Revealed by Phylogenomic Analyses. *Genome Biol Evol* (2017) 9:521–535. doi:  
676 10.1093/gbe/evx022
- 677 6. Campbell F, Strang C, Ferguson N, Cori A, Jombart T. When are pathogen genome  
678 sequences informative of transmission events? *PLoS Pathog* (2018) 14:1–21. doi:  
679 10.1371/journal.ppat.1006885
- 680 7. WHO. Zoonotic tuberculosis factsheet. (2017) [http://www.who.int/tb/areas-of-](http://www.who.int/tb/areas-of-work/zoonotic-tb/en)  
681 [work/zoonotic-tb/en](http://www.who.int/tb/areas-of-work/zoonotic-tb/en) [Accessed October 26, 2022]
- 682 8. Rossi G, Crispell J, Balaz D, Lycett S, Delahay R, Kao R. Identifying likely transmission  
683 pairs with pathogen sequence data using Kolmogorov Forward Equations; an  
684 application to *M. bovis* in cattle and badgers. *Sci Rep* (2020)1–13. doi:  
685 10.1101/2020.06.11.146894
- 686 9. Bernitz N, Kerr TJ, Goosen WJ, Chileshe J, Higgitt RL, Roos EO, Meiring C, Gumbo R, de  
687 Waal C, Clarke C, et al. Review of Diagnostic Tests for Detection of *Mycobacterium*  
688 *bovis* Infection in South African Wildlife. *Front Vet Sci* (2021) 8:1–11. doi:  
689 10.3389/fvets.2021.588697
- 690 10. Kelly RF, Gordon LG, Egbe NF, Freeman EJ, Mazeri S, Ngwa VN, Tanya V, Sander M,  
691 Ndip L, Muwonge A, et al. Bovine Tuberculosis Antemortem Diagnostic Test



- 692 Agreement and Disagreement in a Naturally Infected African Cattle Population.  
693 (2022) 9: doi: 10.3389/fvets.2022.877534
- 694 11. Sichewo PR, Vander Kelen C, Thys S, Michel AL. Risk practices for bovine tuberculosis  
695 transmission to cattle and livestock farming communities living at wildlife-livestock-  
696 human interface in northern Kwazulu Natal, South Sfrica. *PLoS Negl Trop Dis* (2020)  
697 14:1–18. doi: 10.1371/journal.pntd.0007618
- 698 12. Olea-Popelka F, Muwonge A, Perera A, Dean AS, Mumford E, Erlacher-Vindel E,  
699 Forcella S, Silk BJ, Ditiu L, El Idrissi A, et al. Zoonotic tuberculosis in human beings  
700 caused by *Mycobacterium bovis*—a call for action. *Lancet Infect Dis* (2017) 17:e21–  
701 e25. doi: 10.1016/S1473-3099(16)30139-6
- 702 13. Koro Koro F, Ngatchou AF, Portal JL, Gutierrez C, Etoa FX, Eyangoh SI. The genetic  
703 population structure of *Mycobacterium bovis* strains isolated from cattle slaughtered  
704 at the Yaoundé and Douala abattoirs in Cameroon. *OIE Rev Sci Tech* (2015) 34:1001–  
705 1010. doi: 10.20506/rst.34.3.2390
- 706 14. Egbe NF, Muwonge A, Ndip L, Kelly RF, Sander M, Tanya V, Ngwa VN, Handel IG,  
707 Novak A, Ngandalo R, et al. Abattoir-based estimates of mycobacterial infections in  
708 Cameroon. *Sci Rep* (2016) 6:1–14. doi: 10.1038/srep24320
- 709 15. Awah-Ndukum J, Kudi AC, Bradley G, Ane-Anyangwe I, Titanji VPK, Fon-Tebug S,  
710 Tchoumboue J. Prevalence of bovine tuberculosis in cattle in the highlands of  
711 Cameroon based on the detection of lesions in slaughtered cattle and tuberculin skin  
712 tests of live cattle. *Vet Med (Praha)* (2012) 57:59–76. doi: 10.17221/5252-VETMED
- 713 16. Egbe NF, Muwonge A, Ndip L, Kelly RF, Sander M, Tanya V, Ngwa VN, Handel IG,  
714 Novak A, Ngandalo R, et al. Molecular epidemiology of *Mycobacterium bovis* in  
715 Cameroon. *Sci Rep* (2017) 7:1–17. doi: 10.1038/s41598-017-04230-6
- 716 17. Kamerbeek J, Schouls L, Kolk A, Van Agterveld M, Van Soolingen D, Kuijper S,  
717 Bunschoten A, Molhuizen H, Shaw R, Goyal M, et al. Simultaneous detection and  
718 strain differentiation of *Mycobacterium tuberculosis* for diagnosis and epidemiology. *J*  
719 *Clin Microbiol* (1997) 35:907–914. doi: 10.1128/jcm.35.4.907-914.1997
- 720 18. Supply P, Allix C, Lesjean S, Cardoso-Oelemann M, Rüsç-Gerdes S, Willery E, Savine  
721 E, De Haas P, Van Deutekom H, Roring S, et al. Proposal for standardization of  
722 optimized mycobacterial interspersed repetitive unit-variable-number tandem repeat  
723 typing of *Mycobacterium tuberculosis*. *J Clin Microbiol* (2006) 44:4498–4510. doi:

- 724 10.1128/JCM.01392-06
- 725 19. Reis AC, Salvador LCM, Robbe-Austerman S, Tenreiro R, Botelho A, Albuquerque T,  
726 Cunha M V. Article whole genome sequencing refines knowledge on the population  
727 structure of *Mycobacterium bovis* from a multi-host tuberculosis system.  
728 *Microorganisms* (2021) 9: doi: 10.3390/microorganisms9081585
- 729 20. Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian  
730 phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol* (2018)  
731 4:1–5. doi: 10.1093/ve/vey016
- 732 21. Lemey P, Rambaut A, Welch JJ, Suchard MA. Phylogeography takes a relaxed random  
733 walk in continuous space and time. *Mol Biol Evol* (2010) 27:1877–1885. doi:  
734 10.1093/molbev/msq067
- 735 22. Dellicour S, Rose R, Faria NR, Lemey P, Pybus OG. SERAPHIM: Studying environmental  
736 rasters and phylogenetically informed movements. *Bioinformatics* (2016) 32:3204–  
737 3206. doi: 10.1093/bioinformatics/btw384
- 738 23. Motta P, Porphyre T, Handel I, Hamman SM, Ngu Ngwa V, Tanya V, Morgan K,  
739 Christley R, Bronsvoort BMD. Implications of the cattle trade network in Cameroon  
740 for regional disease prevention and control. *Sci Rep* (2017) 7:1–13. doi:  
741 10.1038/srep43932
- 742 24. MMM O. BovTB-nf-docker. (2022) <https://github.com/oxfordmmm/BovTB-nf-docker>  
743 [Accessed October 31, 2022]
- 744 25. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina  
745 sequence data. *Bioinformatics* (2014) 30:2114–2120. doi:  
746 10.1093/bioinformatics/btu170
- 747 26. Md V, Misra S, Li H, Aluru S. Efficient architecture-aware acceleration of BWA-MEM  
748 for multicore systems. *Proceedings - 2019 IEEE 33rd International Parallel and  
749 Distributed Processing Symposium, IPDPS 2019*. IEEE (2019). p. 314–324 doi:  
750 10.1109/IPDPS.2019.00041
- 751 27. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,  
752 Durbin R. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* (2009)  
753 25:2078–2079. doi: 10.1093/bioinformatics/btp352
- 754 28. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome  
755 Biol* (2019) 20:1–13. doi: 10.1186/s13059-019-1891-0

- 756 29. Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: Estimating species abundance in  
757 metagenomics data. *PeerJ Comput Sci* (2017) 2017:1–17. doi: 10.7717/peerj-cs.104
- 758 30. Li H. A statistical framework for SNP calling, mutation discovery, association mapping  
759 and population genetical parameter estimation from sequencing data. *Bioinformatics*  
760 (2011) 27:2987–2993. doi: 10.1093/bioinformatics/btr509
- 761 31. Seemann T. snippy: Rapid haploid variant calling and core genome alignment. (2022)  
762 <https://github.com/tseemann/snippy> [Accessed October 31, 2022]
- 763 32. Jiang Y, Jiang Y, Wang S, Zhang Q, Ding X. Optimal sequencing depth design for whole  
764 genome re-sequencing in pigs. *BMC Bioinformatics* (2019) 20:1–12. doi:  
765 10.1186/s12859-019-3164-z
- 766 33. USDA. vSNP. <https://github.com/USDA-VS/vSNP> [Accessed January 16, 2023]
- 767 34. Warren RM, Pittius NCG Van, Barnard M, Hesselting A, Engelke E, Kock M De,  
768 Gutierrez MC, Chege GK, Victor TC, Hoal EG, et al. Differentiation of *Mycobacterium*  
769 *tuberculosis* complex by PCR amplification of genomic regions of difference. *Int J*  
770 *Tuberc Lung Dis* (2006) 10:818–822.
- 771 35. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
772 transform. *Bioinformatics* (2009) 25:1754–1760. doi: 10.1093/bioinformatics/btp324
- 773 36. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane  
774 T, McCarthy SA, Davies RM, et al. Twelve years of SAMtools and BCFtools.  
775 *Gigascience* (2021) 10:1–4. doi: 10.1093/gigascience/giab008
- 776 37. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search  
777 tool. *J Mol Biol* (1990) 215:403–410. doi: 10.1016/S0022-2836(05)80360-2
- 778 38. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov  
779 JP. Integrative genomics viewer. *Nat Biotechnol* (2011) 29:24–26. doi:  
780 10.1038/nbt.1754
- 781 39. Wattam AR, Davis JJ, Assaf R, Boisvert S, Brettin T, Bun C, Conrad N, Dietrich EM, Disz  
782 T, Gabbard JL, et al. Improvements to PATRIC, the all-bacterial bioinformatics  
783 database and analysis resource center. *Nucleic Acids Res* (2017) 45:D535–D542. doi:  
784 10.1093/nar/gkw1017
- 785 40. Loiseau C, Menardo F, Aseffa A, Hailu E, Gumi B, Ameni G, Berg S, Rigouts L, Robbe-  
786 Austerman S, Zinsstag J, et al. An African origin for *Mycobacterium bovis*. *Evol Med*  
787 *Public Heal* (2020) 2020:49–59. doi: 10.1093/EMPH/EOAA005

- 788 41. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von Haeseler A,  
789 Lanfear R, Teeling E. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic  
790 Inference in the Genomic Era. *Mol Biol Evol* (2020) 37:1530–1534. doi:  
791 10.1093/molbev/msaa015
- 792 42. Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ. W-IQ-TREE: a fast online  
793 phylogenetic tool for maximum likelihood analysis. *Nucleic Acids Res* (2016)  
794 44:W232–W235. doi: 10.1093/NAR/GKW256
- 795 43. Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and  
796 evolutionary analyses in {R}. *Bioinformatics* (2018) 35:526–528.
- 797 44. R Core Team. R: A Language and Environment for Statistical Computing. (2021)  
798 <https://www.r-project.org/>
- 799 45. Rambaut A, Lam TT, Carvalho LM, Pybus OG. Exploring the temporal structure of  
800 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* (2016)  
801 2:1–7. doi: 10.1093/ve/vew007
- 802 46. Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO, Huelsenbeck JP,  
803 Ronquist F, Swofford DL, Cummings MP, et al. BEAGLE: An application programming  
804 interface and high-performance computing library for statistical phylogenetics. *Syst*  
805 *Biol* (2012) 61:170–173. doi: 10.1093/sysbio/syr100
- 806 47. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior summarization in  
807 Bayesian phylogenetics using Tracer 1.7. *Syst Biol* (2018) 67:901–904. doi:  
808 10.1093/sysbio/syy032
- 809 48. Hasegawa M, Kishino H, Yano T aki. Dating of the human-ape splitting by a molecular  
810 clock of mitochondrial DNA. *J Mol Evol* (1985) 22:160–174. doi: 10.1007/BF02101694
- 811 49. Crispell J, Zadoks RN, Harris SR, Paterson B, Collins DM, De-Lisle GW, Livingstone P,  
812 Neill MA, Biek R, Lycett SJ, et al. Using whole genome sequencing to investigate  
813 transmission in a multi-host system: Bovine tuberculosis in New Zealand. *BMC*  
814 *Genomics* (2017) 18:180. doi: 10.1186/s12864-017-3569-x
- 815 50. Duault H, Michelet L, Boschioli M-L, Durand B, Canini L. A Bayesian evolutionary  
816 model towards understanding wildlife contribution to F4-family *Mycobacterium bovis*  
817 transmission in the South-West of France. *Vet Res* (2022) 53:1–12. doi:  
818 10.1186/s13567-022-01044-x
- 819 51. Salvador LCM, O’Brien DJ, Cosgrove MK, Stuber TP, Schooley A, Crispell J, Church S,

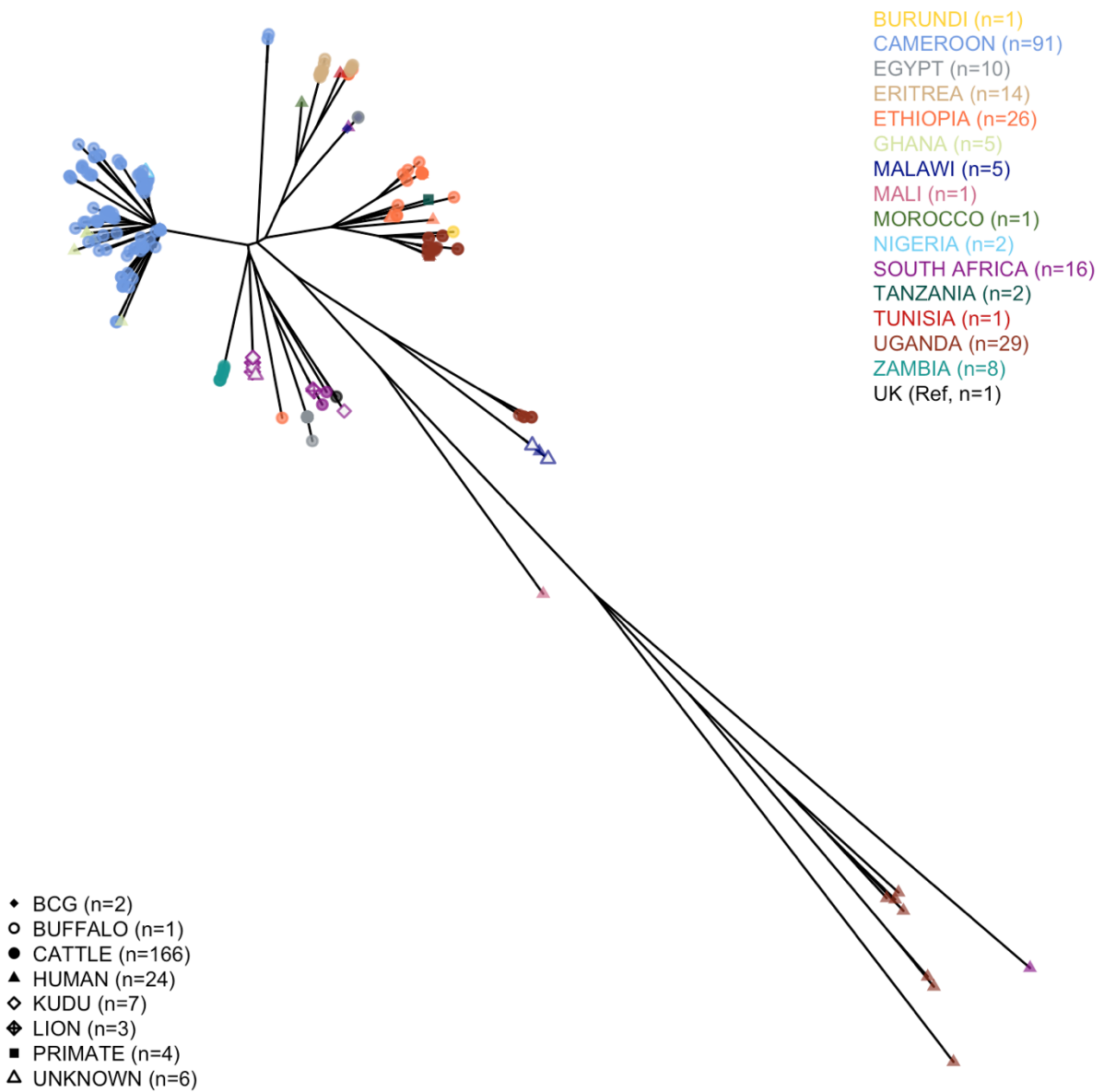
- 820 Grohn YT, Robbe-Austerman S, Kao RR. Disease management at the wildlife-livestock  
821 interface: using whole-genome sequencing to study the role of elk in *Mycobacterium*  
822 *bovis* transmission in Michigan, USA. *Mol Ecol* (2019)1–14. doi: 10.1111/mec.15061
- 823 52. Gill MS, Lemey P, Faria NR, Rambaut A, Shapiro B, Suchard MA. Improving bayesian  
824 population dynamics inference: A coalescent-based model for multiple loci. *Mol Biol*  
825 *Evol* (2013) 30:713–724. doi: 10.1093/molbev/mss265
- 826 53. Hill V, Baele G. Bayesian Estimation of Past Population Dynamics in BEAST 1.10 Using  
827 the Skygrid Coalescent Model. *Mol Biol Evol* (2019) 36:2620–2628. doi:  
828 10.1093/molbev/msz172
- 829 54. Dellicour S, Rose R, Pybus OG. Explaining the geographic spread of emerging  
830 epidemics: A framework for comparing viral phylogenies and environmental  
831 landscape data. *BMC Bioinformatics* (2016) 17:1–12. doi: 10.1186/s12859-016-0924-x
- 832 55. McRae BH. Isolation By Resistance. *Evolution (N Y)* (2006) 60:1551. doi: 10.1554/05-  
833 321.1
- 834 56. Rossi G, Crispell J, Brough T, Lycett SJ, White PCL, Allen A, Ellis RJ, Gordon S V.,  
835 Harwood R, Palkopoulou E, et al. Phylodynamic analysis of an emergent  
836 *Mycobacterium bovis* outbreak in an area with no previously known wildlife  
837 infections. *J Appl Ecol* (2022) 59:210–222. doi: 10.1111/1365-2664.14046
- 838 57. Crispell J, Benton CHCH, Balaz D, De Maio N, Ahkmetova A, Allen A, Biek R, Presho  
839 ELEL, Dale J, Hewinson G, et al. Combining genomics and epidemiology to analyse bi-  
840 directional transmission of *Mycobacterium bovis* in a multi-host system. *Elife* (2019)  
841 8:1–36. doi: <https://doi.org/10.7554/eLife.45833.001>
- 842 58. Elith J, Leathwick JR, Hastie T. A working guide to boosted regression trees. *J Anim*  
843 *Ecol* (2008) 77:802–813. doi: 10.1111/j.1365-2656.2008.01390.x
- 844 59. Hijmans RJ, Phillips S, Leathwick J, Elith J. dismo: Species Distribution Modeling.  
845 (2021) <https://cran.r-project.org/package=dismo>
- 846 60. Greenwell B, Boehmke B, Cunningham J, Developers GBM. gbm: Generalized Boosted  
847 Regression Models. (2022) <https://cran.r-project.org/package=gbm>
- 848 61. Brock PM, Fornace KM, Grigg MJ, Anstey NM, William T, Cox J, Drakeley CJ, Ferguson  
849 HM, Kao RR. Predictive analysis across spatial scales links zoonotic malaria to  
850 deforestation. *Proc R Soc B Biol Sci* (2019) 286: doi: 10.1098/rspb.2018.2351
- 851 62. Newman M. *Networks: An Introduction*. New York, NY, USA: Oxford University Press,

- 852 Inc. (2010).
- 853 63. Csárdi G, Nepusz T. The igraph software package for complex network research. *J*  
854 *Comput Appl* (2014) Complex Sy:9. doi: 10.3724/SP.J.1087.2009.02191
- 855 64. Kuhn M, Weston S, Keefer C, Engelhardt A, Cooper T, Mayer Z, Kenkel B, Team RC,  
856 Benesty M, Lescarbeau R, et al. Classification and Regression Training. (2016)198.  
857 <https://cran.r-project.org/package=caret>
- 858 65. Müller B, Hilty M, Berg S, Garcia-Pelayo MC, Dale J, Boschioli ML, Cadmus S,  
859 Ngandolo BNR, Godreuil S, Diguimbaye-Djaibé C, et al. African 1, an epidemiologically  
860 important clonal complex of *Mycobacterium bovis* dominant in Mali, Nigeria,  
861 Cameroon, and Chad. *J Bacteriol* (2009) 191:1951–1960. doi: 10.1128/JB.01590-08
- 862 66. Rodríguez S, Romero B, Bezos J, de Juan L, Álvarez J, Castellanos E, Moya N, Lozano F,  
863 González S, Sáez-Llorente JL, et al. High spoligotype diversity within a *Mycobacterium*  
864 *bovis* population: Clues to understanding the demography of the pathogen in Europe.  
865 *Vet Microbiol* (2010) 141:89–95. doi: 10.1016/j.vetmic.2009.08.007
- 866 67. Hauer A, De Cruz K, Cochard T, Godreuil S, Karoui C, Henault S, Bulach T, Bañuls AL,  
867 Biet F, Boschioli ML. Genetic evolution of *Mycobacterium bovis* causing tuberculosis  
868 in livestock and wildlife in France since 1978. *PLoS One* (2015) 10:1–17. doi:  
869 10.1371/journal.pone.0117103
- 870 68. Matos F, Cunha M V., Canto A, Albuquerque T, Amado A, Botelho A. Snapshot of  
871 *Mycobacterium bovis* and *Mycobacterium caprae* infections in livestock in an area  
872 with a low incidence of bovine tuberculosis. *J Clin Microbiol* (2010) 48:4337–4339.  
873 doi: 10.1128/JCM.01762-10
- 874 69. Kass RE, Raftery AE. Bayes Factors. *J Am Stat Assoc* (1995) 90:773–795. doi:  
875 10.1080/01621459.1995.10476572
- 876 70. Real LA, Biek R. Spatial dynamics and genetics of infectious diseases on  
877 heterogeneous landscapes. *J R Soc Interface* (2007) 4:935–948. doi:  
878 10.1098/rsif.2007.1041
- 879 71. Pozo P, Lorente-Leal V, Robbe-Austerman S, Hicks J, Stuber T, Bezos J, de Juan L, Saez  
880 JL, Romero B, Alvarez J. Use of Whole-Genome Sequencing to Unravel the Genetic  
881 Diversity of a Prevalent *Mycobacterium bovis* Spoligotype in a Multi-Host Scenario in  
882 Spain. *Front Microbiol* (2022) 13: doi: 10.3389/fmicb.2022.915843
- 883 72. Otchere ID, van Tonder AJ, Asante-Poku A, Sánchez-Busó L, Coscollá M, Osei-Wusu S,

- 884 Asare P, Aboagye SY, Ekuban SA, Yahayah AI, et al. Molecular epidemiology and  
885 whole genome sequencing analysis of clinical *Mycobacterium bovis* from Ghana. *PLoS*  
886 *One* (2019) 14:1–13. doi: 10.1371/journal.pone.0209395
- 887 73. Mohamed A. Bovine tuberculosis at the human–livestock–wildlife interface and its  
888 control through one health approach in the Ethiopian Somali Pastoralists: A review.  
889 *One Heal* (2020) 9:100113. doi: 10.1016/j.onehlt.2019.100113
- 890 74. Valerio VC, Walther OJ, Eilittä M, Cissé B, Muneeppeerakul R, Kiker GA. Network  
891 analysis of regional livestock trade in West Africa. *PLoS One* (2020) 15:1–20. doi:  
892 10.1371/journal.pone.0232681
- 893 75. Trewby H, Wright DM, Skuce RA, McCormick C, Mallon TR, Presho EL, Kao RR, Haydon  
894 DT, Biek R. Relative abundance of *Mycobacterium bovis* molecular types in cattle: a  
895 simulation study of potential epidemiological drivers. *BMC Vet Res* (2017) 13:268. doi:  
896 10.1186/s12917-017-1190-5
- 897 76. Rodriguez-Campos S, Aranaz A, De Juan L, Sáez-Llorente JL, Romero B, Bezos J,  
898 Jiménez A, Mateos A, Domínguez L. Limitations of spoligotyping and variable-number  
899 tandem-repeat typing for molecular tracing of *Mycobacterium bovis* in a high-  
900 diversity setting. *J Clin Microbiol* (2011) 49:3361–3364. doi: 10.1128/JCM.00301-11
- 901 77. Campbell F, Cori A, Ferguson N, Jombart T. Bayesian inference of transmission chains  
902 using timing of symptoms, pathogen genomes and contact data. *PLoS Comput Biol*  
903 (2019) 15:1–20. doi: 10.1371/journal.pcbi.1006930
- 904 78. Rodriguez-Campos S, Schürch AC, Dale J, Lohan AJ, Cunha M V., Botelho A, Cruz K De,  
905 Boschiroli ML, Boniotti MB, Pacciarini M, et al. European 2 - A clonal complex of  
906 *Mycobacterium bovis* dominant in the Iberian Peninsula. *Infect Genet Evol* (2012)  
907 12:866–872. doi: 10.1016/j.meegid.2011.09.004
- 908 79. Zimpel CK, Patané JSL, Guedes ACP, de Souza RF, Silva-Pereira TT, Camargo NCS, de  
909 Souza Filho AF, Ikuta CY, Neto JSF, Setubal JC, et al. Global Distribution and Evolution  
910 of *Mycobacterium bovis* Lineages. *Front Microbiol* (2020) 11:1–19. doi:  
911 10.3389/fmicb.2020.00843
- 912 80. Renwick AR, White PCL, Bengis RG. Bovine tuberculosis in southern African wildlife: A  
913 multi-species host-pathogen system. *Epidemiol Infect* (2007) 135:529–540. doi:  
914 10.1017/S0950268806007205
- 915 81. Chaters GL, Johnson PCD, Cleaveland S, Crispell J, De Glanville WA, Doherty T,

916 Matthews L, Mohr S, Nyasebwa OM, Rossi G, et al. Analysing livestock network data  
917 for infectious disease control: An argument for routine data collection in emerging  
918 economies. *Philos Trans R Soc B Biol Sci* (2019) 374: doi: 10.1098/rstb.2018.0264  
919 82. Reis AC, Cunha M V. The open pan-genome architecture and virulence landscape of  
920 *Mycobacterium bovis*. *Microb Genomics* (2021) 7: doi: 10.1099/MGEN.0.000664  
921 83. Crispell J, Cassidy S, Kenny K, McGrath G, Warde S, Cameron H, Rossi G, Macwhite T,  
922 White PCL, Lycett S, et al. *Mycobacterium bovis* genomics reveals transmission of  
923 infection between cattle and deer in Ireland. *Microb Genomics* (2020) 6:1–8. doi:  
924 10.1099/mgen.0.000388  
925  
926

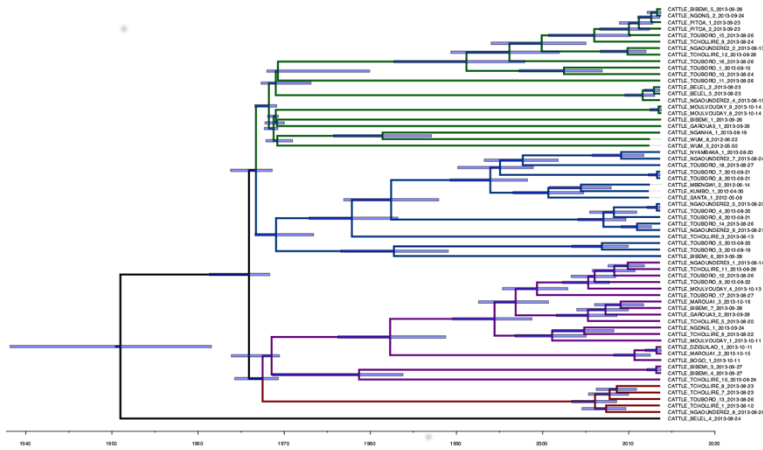




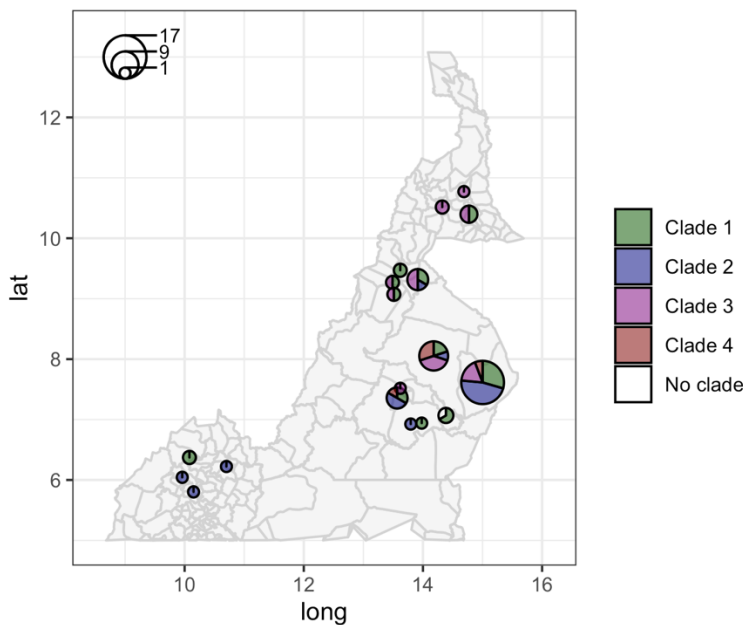
928

929 **Figure 1:** Phylogenetic tree of the African *Mycobacterium bovis* whole-genome sequences  
 930 considered in the study. The tree includes 91 high-quality Cameroonian sequences, 101  
 931 from the EBI dataset, 20 from Patric and the 1997 UK *M. bovis* reference.

932

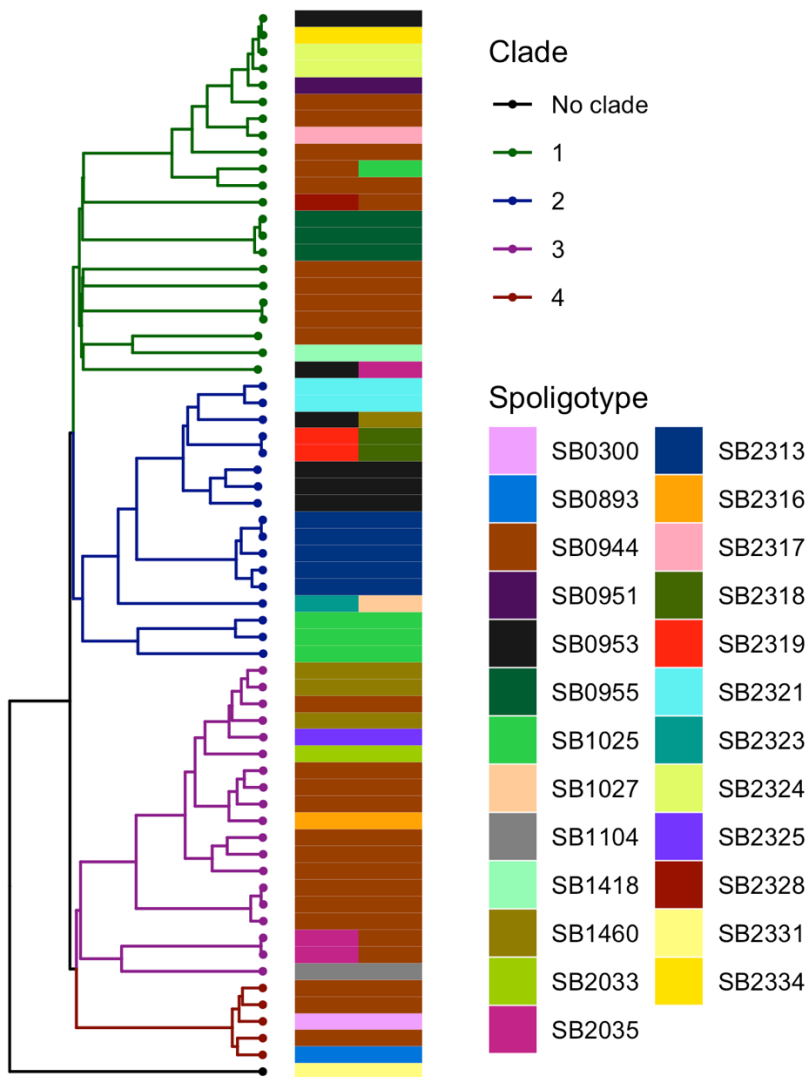


933  
 934 **Figure 2:** Phylogenetic time scaled MCC tree of the 64 high-quality *M. bovis* whole-genome  
 935 sequences sampled in Cameroon in 2012 and 2013. The thin lines represent the 95<sup>th</sup> HPD of  
 936 the internal node dates, while the branch colours represent different clades: 1 (green), 2  
 937 (blue), 3 (purple) and 4 (red). A non-time scaled tree showing the genetic distance between  
 938 the 64 sequences is reported in Figure S6.  
 939



940  
 941 **Figure 3:** Geographic distribution of the 64 high-quality *M. bovis* whole-genome sequences  
 942 in Cameroon. Circle sizes correspond to the number of sequences per administrative  
 943 subdivision, and colours represent different clades (clade 1 green, clade 2 blue, clade 3  
 944 purple and clade 4 red).

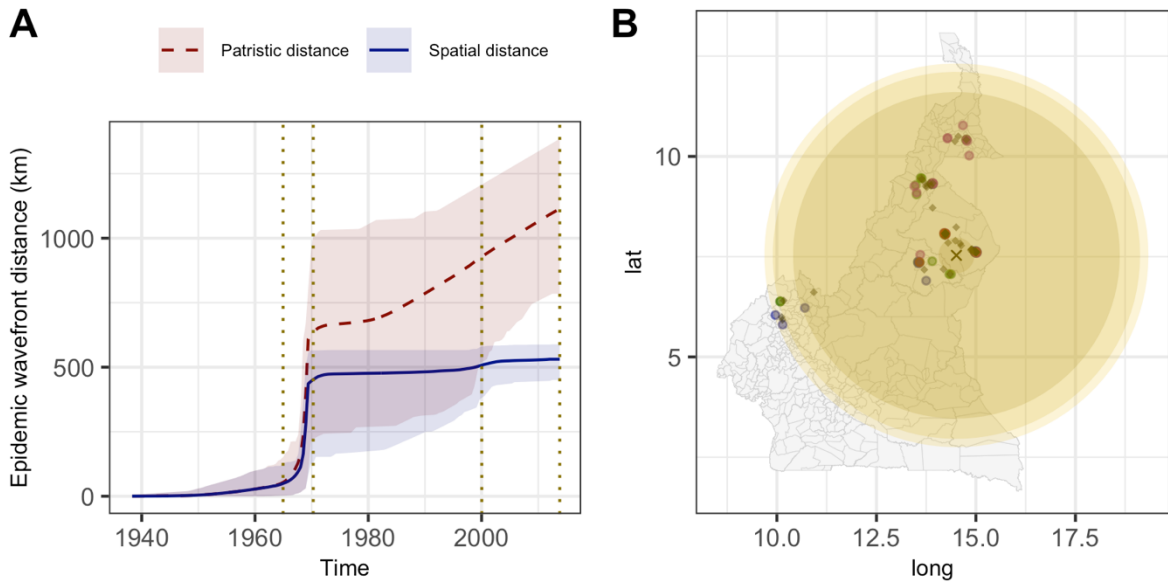
945



946

947 **Figure 4:** Visual comparison between the *M. bovis* phylogenetic time scaled MCC tree and  
948 the spoligotypes obtained by Egbe et al. (16). Ten sequences were associated with two  
949 spoligotypes, because multiple samples from the same animal (up to three) were submitted  
950 for spoligotyping.

951



952

953

**Figure 5:** The estimated epidemic wavefront over time (panel A) and the expansion of the

954

epidemic wavefront on the map (panel B). A: mean (lines) and 95<sup>th</sup> HPD (shades) of the

955

epidemic wavefront spatial distance (blue) and patristic distance (red) over time. B:

956

different yellow shades represent the epidemic wavefront at sequential point in time

957

(marked by vertical dotted lines in panel A), and lighter shades of yellow correspond to

958

more recent expansion; the estimated tree's root location is indicated by the black cross,

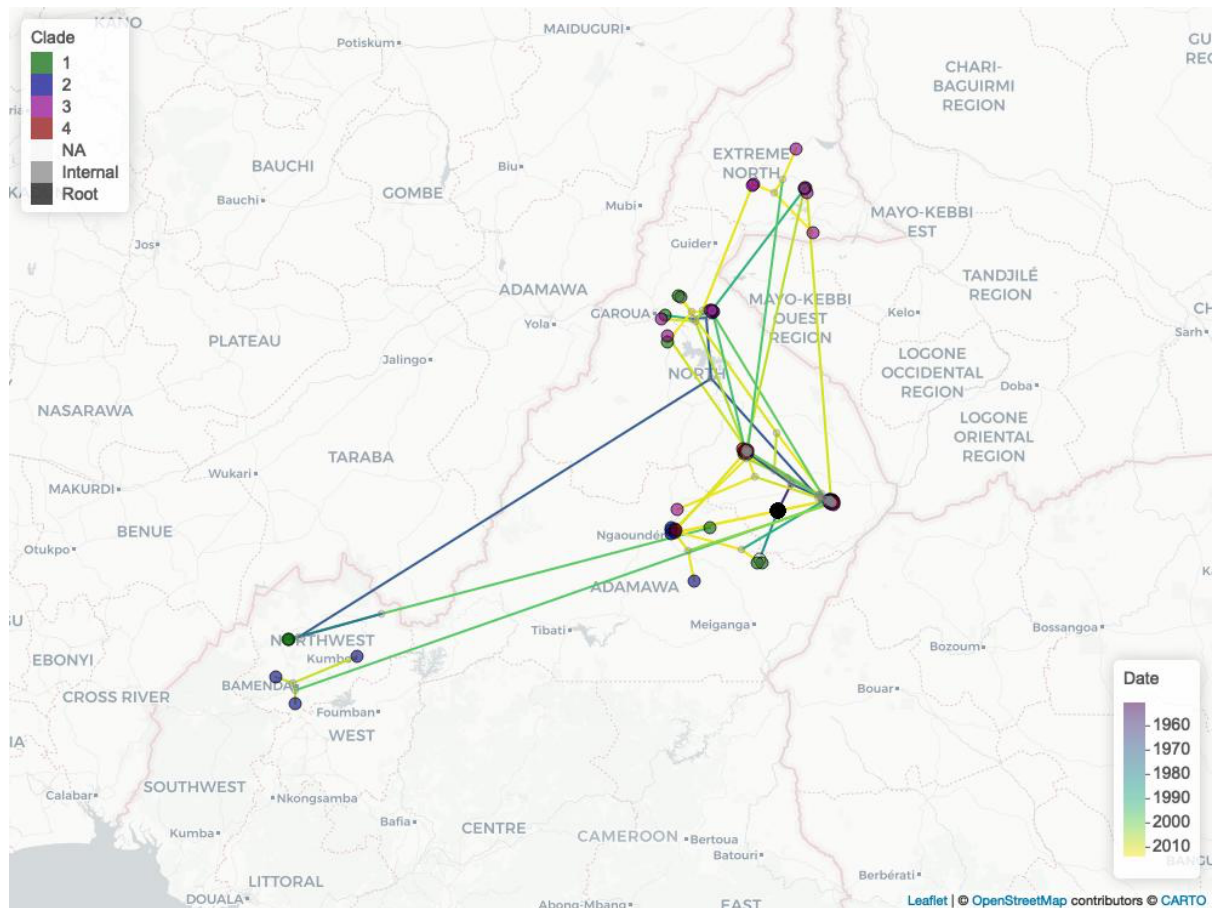
959

diamonds represent the internal nodes estimated locations, and circles the sampled isolates

960

(coloured by clade: 1 green, 2 blue, 3 purple and 4 red).

961

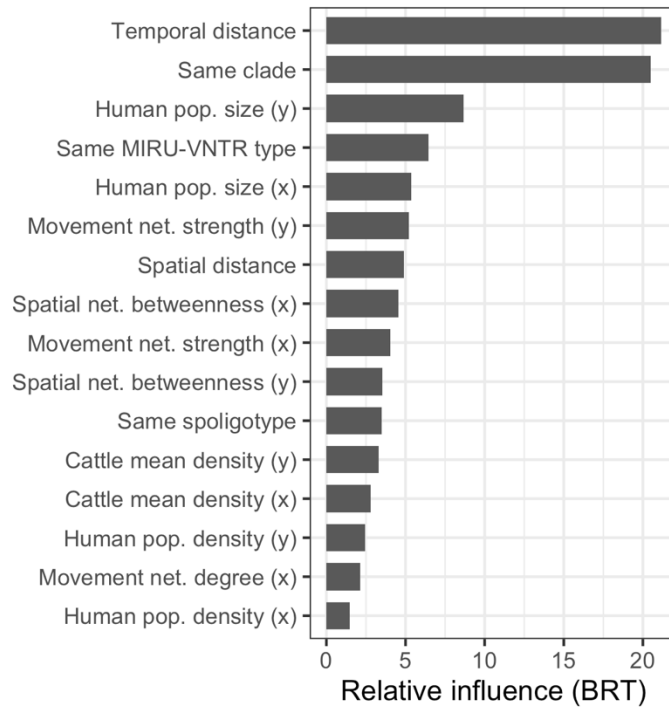


962

963 **Figure 6:** The Cameroonian *M. bovis* epidemic estimated expansions in space and time.

964 Nodes are coloured by clade (1, green; 2, blue; 3 purple; 4, red; no clade, light grey; internal  
 965 nodes, dark grey; tree root, black), while the branches are coloured by estimated movement  
 966 date, from 2007 (purple) to 2013 (yellow).

967



968

969

**Figure 7:** Relative influence of the most relevant variables in the simplified boosted

970

regression tree (BRT) model. The purpose of the BRT model was to explain the SNP distance

971

between the 64 high-quality *M. bovis* isolates. Many variables are calculated between

972

isolates pairs, x refers to the oldest isolate's subdivision, and y to the youngest one.

973

Variable	Type	Path model	Number of positive coefficients	Number of positive Q statistic	Mean Bayes Factor	
					(Randomisation #1)	(Randomisation #2)
Mosaic_shrub_otherv	Resistance	Least cost	100	99	1.89	1.83
<b>Forest</b>	<b>Conductance</b>	<b>Least cost</b>	<b>100</b>	<b>96</b>	<b>3.66</b>	<b>3.66</b>
Mosaic_shrub_otherv	NA	Straight line	100	89	0.62	1.32
<b>Elevation</b>	<b>Conductance</b>	<b>Least cost</b>	<b>100</b>	<b>88</b>	<b>2.39</b>	<b>3.00</b>
Waterbodies	Conductance	Least cost	100	87	1.10	0.97
Cattle_density	Resistance	Least cost	99	77	2.49	2.88
Cattle_density	NA	Straight line	100	73	Not run	Not run
Cattle_density	Conductance	Least cost	100	70	Not run	Not run
Grassland_cropland	Resistance	Least cost	100	66	Not run	Not run
Grassland_cropland	NA	Straight line	100	56	Not run	Not run
Mosaic_shrub_otherv	Conductance	Least cost	100	44	Not run	Not run
Roads_intersections	Conductance	Least cost	100	42	Not run	Not run
Waterbodies	Resistance	Least cost	100	38	Not run	Not run
Waterbodies	NA	Straight line	100	27	Not run	Not run
Grassland_cropland	Conductance	Least cost	100	15	Not run	Not run
Forest	Resistance	Least cost	100	12	Not run	Not run
Elevation	NA	Straight line	100	7	Not run	Not run
Forest	NA	Straight line	100	3	Not run	Not run
Pop_density	Conductance	Least cost	99	40	Not run	Not run
Elevation	Resistance	Least cost	99	15	Not run	Not run
Roads_length	NA	Straight line	97	0	Not run	Not run
Pop_density	NA	Straight line	96	16	Not run	Not run
Pop_density	Resistance	Least cost	96	7	Not run	Not run
Roads_length	Conductance	Least cost	92	11	Not run	Not run
Roads_length	Resistance	Least cost	59	0	Not run	Not run
Roads_intersections	NA	Straight line	56	1	Not run	Not run
Roads_intersections	Resistance	Least cost	6	1	Not run	Not run

975

976 **Table 1:** Results of the analysis on nine spatial variables, assuming two path models, straight  
977 line and least cost, and for the least cost path, whether the variable worked as a  
978 conductance or resistance. Results show the number of positive coefficients for the 100  
979 sampled trees, the number of positive Q statistics, and the mean Bayes factor calculated  
980 over 10 randomisations, testing two algorithms: 1) randomisations of nodes positions while  
981 maintaining branches lengths, tree topology and location of the most ancestral node; and 2)  
982 randomisations of nodes positions while maintaining only the branches lengths.

983