1    **Phylogenomics Sheds Light on the population structure of *Mycobacterium bovis* from a Multi-**

2    **Host Tuberculosis System**

3

4    Ana C. Reis[1,2], Liliana C.M. Salvador[3,4,5], Suelee Robbe-Austerman[6], Rogério Tenreiro[2], Ana Botelho[7],

5    Teresa Albuquerque[7], Mónica V. Cunha[1,2*]

6

7    [1]Centre for Ecology, Evolution and Environmental Changes (cE3c), Faculdade de Ciências da

8    Universidade de Lisboa, Lisboa, Portugal

9    [2]Biosystems & Integrative Sciences Institute (BioISI), Faculdade de Ciências da Universidade de

10   Lisboa, Lisboa, Portugal

11   [3]Department of Infectious Diseases, College of Veterinary Medicine, University of Georgia, Athens,

12   Georgia, USA

13   [4]Institute of Bioinformatics, University of Georgia, Athens, Georgia, USA

14   [5]Center for the Ecology of Infectious Diseases, University of Georgia, Athens, Georgia, USA

15   [6]USDA/APHIS National Veterinary Services Laboratories, USA

16   [7]INIAV, IP - National Institute for Agrarian and Veterinary Research

17

18

19   *Correspondence: mscunha@fc.ul.pt; Centre for Ecology, Evolution and Environmental Changes

20   (cE3c), Faculdade de Ciências, Universidade de Lisboa, Campo Grande, C2, Room 2.4.11, 1749-016

21   Lisboa. Phone +351 217 500 000

22 **Abstract**

23 Molecular analyses of *Mycobacterium bovis* based on spoligotyping and Variable Number Tandem

24 Repeat (MIRU-VNTR) brought insights into the epidemiology of animal tuberculosis (TB) in Portugal,

25 showing high genotypic diversity of circulating strains that mostly cluster within the European 2 clonal

26 complex. The genetic relatedness of *M. bovis* isolates from cattle and wildlife have also suggested

27 sustained transmission within this multi-host system. However, while previous surveillance highlighted

28 prevalent genotypes in areas where livestock and wild ungulates are sympatric and provided valuable

29 information on the prevalence and spatial occurrence of TB, links at the wildlife-livestock interfaces

30 were established mainly via genotype associations. Therefore, evidence at a local fine scale of

31 transmission events linking wildlife hosts and cattle remains lacking. Here, we explore the advantages

32 of whole genome sequencing (WGS) applied to cattle, red deer and wild boar isolates to reconstruct

33 the evolutionary dynamics of *M. bovis* and to identify putative pathogen transmission events. Whole

34 genome sequences of 44 representative *M. bovis* isolates, obtained between 2003 and 2015 from

35 three TB hotspots, were compared through single nucleotide polymorphism (SNP) variant calling

36 analyses. Consistent with previous results combining classical genotyping with Bayesian population

37 admixture modelling, SNP-based phylogenies support the branching of this *M. bovis* population into

38 five genetic clades, three with geographic specificities, as well as the establishment of a SNPs

39 catalogue specific to each clade, which may be explored in the future as phylogenetic markers. The

40 core genome alignment of SNPs was integrated within a spatiotemporal metadata framework to

41 reconstruct transmission networks, which together with inferred secondary cases, further structured

42 this *M. bovis* population by host species and geographic location.

43 WGS of *M. bovis* isolates from Portugal is reported for the first time, refining the spatiotemporal

44 context of transmission events and providing further support to the key role of red deer and wild boar

45 on the persistence of animal TB in this Iberian multi-host system.

46

47 **Keywords:** animal tuberculosis, multi-host transmission, *Mycobacterium bovis*, transmission network,

48 whole-genome sequencing

49 **1. Introduction**

50 *Mycobacterium bovis* is an important pathogen, responsible for causing animal tuberculosis (TB) in

51 livestock and wildlife vertebrates, as well as in humans (Brites et al., 2018; Gagneux, 2018). Cattle

52 (*Bos taurus*) is the main livestock affected species, while several reports evidence the importance of

53 the livestock-wildlife interface for disease maintenance (Corner, Murphy, & Gormley, 2011; Fitzgerald

54 & Kaneene, 2012; M. V. Palmer, 2007; Mitchell V. Palmer, Thacker, Waters, Gortázar, & Corner,

55 2012). In the Iberian Peninsula, red deer (*Cervus elephus*) and wild boar (*Sus scrofa*) have both been

56 implicated in the transmission of *M. bovis* to cattle via direct and indirect routes and in pathogen

57 persistence across ecosystems, depending on the specificities of the epidemiological scenario and the

58 ecological relationships established by the hosts (Barasona, Torres, Aznar, Gortázar, & Vicente, 2017;

59 Cunha et al., 2012; Gortázar et al., 2008; Naranjo, Gortázar, Vicentea, & de la Fuente, 2008; Santos

60 et al., 2009; Vieira-Pinto et al., 2011). The presence of maintenance hosts in the wild is associated

61 with difficulties in the success of test and slaughter schemes implemented in the cattle population, but

62 it also brings concerns regarding wildlife welfare, biodiversity and public health.

63 To date, works with reference to *M. bovis* molecular characterization in Portugal have been based on

64 the analysis of repetitive genomic regions, namely spoligotyping (spacer oligonucleotide typing) and

65 MIRU-VNTR (*Mycobacterial Interspersed Repetitive Units-Variable Number Tandem Repeats*). The

66 focus has been placed over isolates from cattle, red deer and wild boar from TB hotspot areas located

67 in the central and south regions of the country (Cunha et al., 2012; Duarte, Domingos, Amado, Cunha,

68 & Botelho, 2010; Reis, Tenreiro, Albuquerque, Botelho, & Cunha, 2020). These works provided

69 evidences for *M. bovis* population diversity and structure, highlighting the main genotypes across host

70 species and regions, as well as intra- and inter-specific transmission (Cunha et al., 2012; Duarte et al.,

71 2010; Reis et al., 2020). However, these molecular typing approaches have explored epidemiological

72 links via genotype associations, which are not sufficiently discriminatory to accurately assess

73 transmission at a fine-scale, nor to gain insights on the roles exerted by different species in a multi-

74 host system. But understanding the evolutionary processes driving transmission among sympatric

75 wildlife reservoirs and livestock populations is crucial for an effective management of animal TB in an

76 endemic system.

77 The progressive application of whole genome sequencing (WGS) to infectious disease systems has

78 resulted in unprecedented advances in the ability to resolve epidemiological information at different

79    scales. WGS data provides higher discriminatory power than classical molecular approaches for

80    resolving complex outbreak situations, allowing a finer definition of the spatiotemporal context in which

81    pathogen spread and persistence occurs. WGS also aids in the identification of the infection source,

82    the establishment of epidemiological links, and the reconstruction of transmission chains (Biek et al.,

83    2012; Glaser et al., 2016; Price-Carter et al., 2018).

84    When considering the livestock-wildlife interface, WGS has been used to demonstrate the close

85    genetic relationship among *M. bovis* isolates recovered from sympatric cattle and wildlife populations,

86    in different epidemiological settings, including UK (Biek et al., 2012; Crispell et al., 2019; Trewby et al.,

87    2016), Ireland (Crispell et al., 2020), New Zealand (Crispell et al., 2017; Price-Carter et al., 2018) and

88    United States of America (Glaser et al., 2016; Salvador et al., 2019). In this context, single nucleotide

89    polymorphisms (SNPs) emerged as good phylogenetic markers, helping in the definition of *M. bovis*

90    population structure, having been recently used to define four *M. bovis* lineages, and to inform

91    transmission models (Biek et al., 2012; J. Guerra-Assunção et al., 2015; Price-Carter et al., 2018).

92    When placed together with data concerning the time needed for this slow growing bacterium to

93    accumulate new SNPs, this information can provide temporal clues on the emergence and divergence

94    of specific genotypes.

95    With the aim to improve knowledge on *M. bovis* population structure and transmission within and

96    across TB hotspots in Portugal, WGS of 44 *M. bovis* obtained between 2003 and 2015 from cattle, red

97    deer and wild boar was completed. These isolates were selected as being representative of the *M.*

98    *bovis* population diversity in those areas, which was previously assessed by a large-scale genotyping

99    study that combined standard genotyping techniques (spoligotyping and MIRU-VNTR) with Bayesian

100   clustering (Reis et al., 2020). The methodological framework aimed to: (1) identify phylogenetic clades

101   and build a catalogue of SNPs that may be used as specific molecular markers of each clade; (2)

102   explore how specific nucleotide differences are associated with distinct host and/or geographic

103   regions; (3) disclose transmission networks for global and local datasets from Portugal; and (4) infer

104   infectivity in distinct host species and geographic locations.

105    **2. Methods**

106    **2.1. *M. bovis* isolates dataset**

107    The 44 *M. bovis* isolates used in this work were recovered from cattle (*n*=16), red deer (*n*=16) and wild

108    boar (*n*=12) from TB hotspot areas in Portugal that encompass the districts (administrative level

109    sample unit) of Castelo Branco, Portalegre and Beja, which are located in inner central and south of

110    the mainland territory (Supplementary Fig. 1 and Supplementary Table 1). This dataset was selected

111    for WGS from a wider *M. bovis* dataset recovered in Portugal (*n*=487) that was previously submitted to

112    molecular characterization by classical genotyping techniques, namely spoligotyping and 8 *loci* MIRU-

113    VNTR (Supplementary Fig. 1) (Reis et al., 2020). To represent the population genetic diversity, two

114    selection criteria were applied: first, isolates should represent major spoligotyping-MIRU type groups

115    and adjacent variants; and second, they should cover different host species, geographic regions and

116    temporal/epidemiological contexts (using year of *M. bovis* isolation as proxy) (Supplementary Fig. 1).

117    **2.2. Ethical Approval**

118    The *M. bovis* dataset analysed here was selected for WGS from a wider *M. bovis* dataset recovered in

119    Portugal (Reis et al. 2020) in the scope of official control plans for animal TB. No animals were

120    sacrificed for the purposes of this study. Isolates were obtained in the national reference laboratory of

121    animal tuberculosis (INIAV, IP) from animal samples either presenting TB-compatible lesions during

122    official inspection and/or animal samples from reactor cattle submitted to official standard screening

123    test for TB [the single intradermal comparative cervical tuberculin (SICCT) test]. None of the authors

124    were responsible for the death of any animals nor were any samples used in the study collected by the

125    authors. All applicable institutional and/or national/international guidelines for the care and use of

126    animals have been followed.

127    **2.3. DNA extraction**

128    Bacteriological culture was performed as described by Reis et al. (2020). Frozen culture stocks of 44

129    *M. bovis* were successfully re-cultured on Middlebrook 7H9 (Difco) medium supplemented with 5%

130    sodium pyruvate and 10% ADS enrichment (50 g albumin, 20 g glucose, 8.5 g sodium chloride in 1 L

131    water), at 37 °C, on a level 3 biosecurity facility.

132    After four weeks' growth, the culture medium was renewed and the cultures were monitored regularly

133    until growth was observed. Cells were harvested, centrifuged and the culture pellet was re-suspended

5

134    in 500 µL PBS and inactivated by heating at 99ºC for 30 minutes. After centrifugation, the

135    supernatants were stored at -20°C until WGS.

136    **2.4. Whole-genome sequencing and SNP analysis**

137    The genomic DNA was sequenced using the Illumina Genome Analyser, according to the

138    manufacturer's specifications with the paired-end module attachment. Forty-two samples were

139    sequenced by MiSeq technology (2x250 pb) at the United States Department of Agriculture (USDA,

140    USA) and the remaining two by HiSeq (2x150 pb) (Eurofins, Germany).

141    The vSNP pipeline, currently available at https://github.com/USDA-VS/vSNP, was used to process the

142    FASTQ files obtained from Illumina sequencing. Briefly, reads were aligned to the *M. bovis* AF2122/97

143    reference genome (NCBI accession number NC_002945.4), using BWA and Samtools (Li & Durbin,

144    2009; Li et al., 2009). Base quality score recalibration, SNP and indel (insertion or deletion) discovery

145    were applied across all isolates using standard filtering parameters or variant quality score

146    recalibration according to Genome Analysis Toolkit (GATK)'s Best Practices recommendations

147    (Depristo et al., 2011; Mckenna et al., 2010; Van der Auwera et al., 2014). Results were filtered using

148    a minimum SAMtools quality score of 150 and AC = 2.

149    Integrated Genomics Viewer (IGV) (version 2.4.19) (Thorvaldsdóttir, Robinson, & Mesirov, 2012) was

150    used to visually validate SNPs and positions with mapping issues or alignment problems. SNPs that

151    fell within Proline-Glutamate (PE) and Proline-Proline Glutamate (PPE) genes were filtered from the

152    analysis, as well as indels.

153    The raw data is deposited in a public domain server at the National Centre for Biotechnology

154    Information (NCBI) SRA database, under BioProject accession number PRJNA682618.

155    **2.5. Phylogenetic analysis**

156    Validated and polymorphic SNPs were concatenated, resulting into a single 1842-nt sequence. MEGA

157    (Molecular Evolutionary Genetics Analysis, version 7.0) (Kumar, Stecher, & Tamura, 2016) was used

158    to conduct phylogenetic analysis, using the maximum likelihood method with 1,000 bootstrap

159    inferences.

160    Distribution of pairwise SNP distance were obtained by applying the Hamming distance, using the

161    library *ape*, and the corresponding heatmap was obtained by library *gplots*, both in R statistical

162    package (Paradis, Claude, & Strimmer, 2004).

163    **2.6. Transmission mapping**

164 The *SeqTrack* library in R statistical package (Jombart, Eggo, Dodd, & Balloux, 2011) was used to

165 infer total and local transmission networks. This library builds a network, minimizing the genomic

166 distance between links and keeping the disease onset dates coherent, reconstructing the most

167 plausible genealogy for a determined group of isolates. Matrices with information concerning host

168 species and geographic region were added to determine the transmission pathway (Jombart et al.,

169 2011).

170 The geographic coordinates of *M. bovis* isolates ($n$=36, for which georeferenced information was

171 complete) allowed the definition of two hotspot areas that were further analysed in local networks

172 (Supplementary Fig. 2). Each *M. bovis* isolate from cattle was assigned to the geographic coordinates

173 of the corresponding livestock herd; isolates from hunter-harvested red deer and wild boar were

174 associated to the centroids of officially-delimited hunting areas, and no substantial change in

175 geographic coordinates was assumed over the study period. Therefore, a total network with all *M.*

176 *bovis* and two local networks, one for Castelo Branco district ($n$=16) and another one for Portalegre

177 ($n$=12) were performed.

178 Putative recent transmission events were established based on a maximum inter-isolate pairwise SNP

179 distance of three SNPs within 5-year range (Crispell et al., 2020).

180 The number of secondary cases was inferred, based on the number of established links between

181 isolates and thereafter grouped by host species and geographic region, being displayed in a violin

182 density plot. The statistical significance between the inferred number of secondary cases and host

183 species or geographic region were tested using a one-way ANOVA test. Results were considered to

184 be significant if $p < 0.05$. All the statistical tests were performed in R software environment (version

185 3.4.6).

186 **2.7. Phylogeographic analysis**

187 A phylogeographic analysis was performed with a combined approach including Gephi version 0.9.2

188 (Bastian & Heymann, 2009) and QGIS (Quantum GIS development Team 2018, version 3.10).

189 A network analysis to explore the relationships established between *M. bovis* isolates, using as node

190 each *M. bovis* and as connection lines the number of shared SNPs, was performed in Gephi for four

191 cumulative temporal periods, and plotted in a map with QGIS. Four networks were generated with a

192 progressive temporal interval of three years between each other, resulting in the definition of period 1

193 (2003-2006), period 2 (2003-2009), period 3 (2003-2012) and period 4 (2003-2015). These temporal

194   windows follow the variation of animal TB epidemiological indicators in Portugal, considering herd

195   prevalence values in cattle and the surveillance measures implemented in wildlife (Supplementary Fig.

196   3). Thus, in periods 1 and 2, the values of herd prevalence steadily decreased; in period 3, an

197   increase in herd prevalence was registered and carcass examination of hunted big game species

198   became mandatory in the epidemiological risk area that is under analysis in this study; and, finally, in

199   the last years added to the timeline, a decrease in herd prevalence was observed (Supplementary Fig.

200   3). The connection lines were established based on absolute values of SNPs; no statistical

201   transformation was performed.

202   A total of 36 *M. bovis* had information concerning geographic coordinates of the corresponding

203   livestock herd or officially-delimited hunting area. For eight *M. bovis* strains [cattle ($n$=5), wild boar

204   ($n$=2), red deer ($n$=1)], the specific geographic coordinates were absent, so the coordinates were

205   assumed as the centroid of the geographic region, so that all *M. bovis* could be plotted in the map. No

206   substantial change in geographic coordinates was assumed over the study period.

207   **2.8. Isolation by distance analysis**

208   A Mantel test was conducted in the R software environment package *ade4* (Dray & Dufour, 2007) to

209   assess correlation between spatial and SNP distances, using 10,000 permutations to assess

210   significance. Only the isolates with geographical coordinates were included in this analysis. The

211   Mantel test was applied to the total dataset ($n$=36) and to local networks of Castelo Branco ($n$=16) and

212   Portalegre ($n$=12).

213   **2.9. Molecular dating**

214   The temporal structure of the sequences was explored with *Tempest* (TEMPoral Exploration of

215   Sequences and Trees, version 1.5.1) (Rambaut, Lam, Carvalho, & Pybus, 2016) and with a date

216   randomization test (DRT) with LSD (Least-Squares Dating, version 0.3-beta) (To, Jung, Lycett, &

217   Gascuel, 2015). In the approach with *Tempest*, a phylogenetic tree performed with the *M. bovis*

218   AF2122/97 reference genome as outgroup was used as input.

219   The least square method implemented in LSD v0.3-beta (To et al., 2015) was applied to estimate the

220   molecular clock rate in the observed data and to perform a DRT with 20 randomized datasets. The

221   QPD (quadratic programming dating) algorithm was used, allowing to estimate the position of the root

222   (option -r as) and calculating the confidence intervals (options -f 100 and -s) (To et al., 2015).

223

224

225 **3. Results**

226 **3.1. SNP-based genotyping and phylogenetic analyses**

227 The sequence reads of 44 *M. bovis* whole genomes representing the genetic diversity of strains

228 circulating in TB hotspots in Portugal were mapped to the assembled reference genome of *M. bovis*

229 AF2122/97 (NC_002945.4) (Supplementary Table 1). The average depth of coverage and genome

230 coverage was 93.6 and 99.57%, respectively (Supplementary Table 2). The SNP alignment had a total

231 of 1842 polymorphic positions, being the majority (86.5%) located in coding regions.

232 The phylogenetic distribution of SNPs grouped *M. bovis* into five related clades, each one with more

233 than 100 clade-defining SNP sites, i.e. polymorphic positions specifically found within each clade

234 member (Fig. 1a and Table 1). Clades were named from A to E, being clade A the largest, counting

235 with 14 *M. bovis* genomes, while clade C is the smallest, encompassing only three (Fig. 1a).

236 The topology of the SNP-based phylogenetic tree and the one based on classical genotyping

237 techniques (combination of spoligotyping and MIRU-VNTR), agree in the large branch division

238 between clades A to D and clade E, and in the clustering of clade D and E members (Fig. 1a and 1b).

239 The majority of *M. bovis* isolates (*n*=34, 77%, clades A to D) cluster within the highly structured clonal

240 complex European 2 (Eu2) (Rodriguez-Campos et al., 2012) that is widely distributed in the Iberian

241 Peninsula. All members from clades A to D possessed the *guaA* gene (G→A) synonymous mutation,

242 which is the hallmark of this clonal complex (Rodriguez-Campos et al., 2012).

243 The clades in the upper phylogenetic branch (clades A to D) registered between 108 to 217 clade-

244 defining SNP sites, while the lower phylogenetic branch (clade E) presented a total of 360 SNPs

245 (Table 1). Moreover, for clades C, D and E it was also possible to identify clade-monomorphic SNP

246 sites (i.e. polymorphic positions present only in clade members and common to them all): 49 SNPs in

247 clade C, 82 in clade D and 352 in clade E (Table 1).  When accounting for the total SNP sites

248 registered per clade, intra-clade homogeneity (i.e. the proportion of monomorphic SNP sites within

249 each clade) ranged from 0% to 82%, in decreasing order: clade E (82%), clade D (22%), clade C

250 (15%) and clades A and B (0%), pointing clade E as the most homogeneous.

251 The differences between phylogenetic branches are also clearly expressed in a heatmap based on the

252 absolute SNP differences between strains, which supports a very clear separation between members

253 of clades A to D and clade E, with a high average diversity (mean SNP distance value of 384 SNPs)

254 (Fig. 2). Grouping by host species revealed that the mean genetic difference within each group is

255    similar (382 SNPs within cattle; 397 SNPs within red deer; and 398 within wild boar), while

256    dissimilarities strike out when comparing the three geographic regions covered by the analyses (376

257    SNPs within Castelo Branco; 405 within Portalegre; and 244 within Beja) (data not shown).

258    Clades A and B were the most widely distributed, being found across the three TB hotspots under

259    study. Clade C was exclusive of Castelo Branco, clade D was found to be absent in Portalegre, clade

260    E was absent in Beja (Fig. 3). In contrast, similar clade distributions were found at the host species

261    level, with strains from the three hosts clustering in the five clades (Fig. 3). When considering

262    geographic region or host species as grouping criteria, group-defining SNP sites could also be

263    identified for all categories under analysis, with Castelo Branco (396 SNPs) and cattle (379 SNPs)

264    yielding the highest number of polymorphic SNPs. However, no monomorphic SNP sites were

265    identified per host species or geographic region.

266    The temporal evolution of the established SNP network was assessed to get insights into the strength

267    of relationships established through time by *M. bovis* strains, considering four progressive temporal

268    windows and using as connection links the number of common SNPs (Fig. 4). Based on this SNP

269    dataset, each clade could not be distinguished from the others based solely on the sampling time of

270    the isolates (Supplementary Fig. 4). Furthermore, this analysis highlighted strong global and local

271    networks in Portugal, with a particular focus to the strength of connections established within

272    geographic regions (maximum shared SNPs between strains in Beja=265, Castelo Branco=370 and

273    Portalegre=365) and between Castelo Branco and Portalegre (*n*=365) (Fig. 4).

274    **3.2. Transmission network**

275    The global transmission network for Portuguese isolates based on the SNP alignment with 1842

276    positions and available metadata was constructed. The connections between *M. bovis* represent the

277    number of SNP differences.

278    In the total network, 43% of individuals are linked as likely sources of infection to at least one another

279    individual (Fig. 5). This analysis has shown that individuals could be sources of infection to up nine

280    others, with seven (16%) linked to one host, three (7%) linked to two hosts, 22 (7%) linked to three

281    hosts and 6 (14%) linked to four or more hosts. Links with zero SNP differences were identified within

282    the same host species (cattle-cattle and wild boar-wild boar) and between different host species (red

283    deer-wild boar), suggesting intra- and inter-specific transmission events (Fig. 5).

11

284    Three large branches could be defined in the transmission network, with strains recovered from each

285    one of the three TB hotspots being placed at the bottom of each branch, respectively (Fig. 5). The left

286    subdivision of the network was exclusively composed by clade E members, while the right subdivision

287    encompasses clades A and B affiliates. The middle branch is more diverse, with members from clades

288    B, C and D (Fig. 5).

289    The number of links established by each strain was used to infer the number of secondary infections,

290    and therefore infectivity. A stratified approach performed by host species and geographic region was

291    represented in comparative violin plots (Fig. 6a and 6b). The graphics reveal the dispersion of the

292    number of secondary infections, from minimum values of zero cases to a maximum of nine cases for

293    cattle and Portalegre, if considering host species or geographic region, respectively. Plus, the higher

294    density areas indicate that, in the majority of cases, the total number of secondary infections

295    generated is maintained between 0 to 3, independently of host species or geographic region, with the

296    exception of Beja (Fig. 6a e 6b). Moreover, it is also possible to conclude that cattle and Beja yield the

297    higher median values of secondary infection (Fig. 6a and 6b). No statistically significant differences

298    were verified when comparing secondary infections median values by host species or geographic

299    location.

300    A comparison of pairwise genetic distance (SNPs) and geographic distance was performed, only

301    including cases with accurate geographic coordinates ($n$=36). The Mantel test revealed no statistically

302    significant association (simulated $p$-value= 0.541).

303    **3.2.1 Local transmission networks**

304    Two hotspot areas in Castelo Branco ($n$=16) and Portalegre ($n$=12) were analysed and a cut-off value

305    of 3 SNP differences in the same temporal context was used to define recent transmission events.

306    Clustering by SNP clade is evident in both networks, with recent transmission events only involving *M.*

307    *bovis* included in the same SNP clade (Fig. 7a and 7b). In Castelo Branco network, two recent

308    transmission events were identified, with one case including individuals from different host species. In

309    both events the wildlife hosts belong to different officially defined hunting area (Fig. 7a). Considering

310    Portalegre's network, one event of recent transmission was identified (Fig. 7b), with the host

311    individuals being sampled from different herds.

312    The Mantel test applied to Castelo Branco and Portalegre datasets revealed that there was no

313    statistically significant association between genetic and spatial distance (simulated *p*-value=0.619 and

314    0.190, respectively), in agreement with findings for the global dataset from Portugal.

315    **3.3. Molecular clock analysis**

316    Before performing an evolutionary analysis, it is crucial to evaluate the clock-like structure of the entire

317    dataset, i.e. the relationship between observed genetic distances and time. Therefore, root-to-tip

318    regression with Tempest and DRT analysis were applied. Under a perfect clock-like behaviour, the

319    distance between the root of the phylogenetic tree and the tips is a linear function of the tip's sampling

320    year, being $R^2$ the degree of clock-like structure. This *M. bovis* dataset exhibits a positive correlation

321    between genetic divergence and sampling time, however the $R^2$ value obtained was low ($R^2$ =0.15),

322    suggesting that this dataset lacks a strong clock-like behaviour. Since the approach with Tempest can

323    only be used as an exploratory analysis of temporal structure, a DRT was performed. The year of

324    sampling was reshuffled 20 times and the clock rate of observed and randomized datasets was

325    estimated using LSD. The clock rate estimate obtained from the observed data was $1.30 \times 10^{-4}$ [95% CI

326    $(1.0 \times 10^{-10}$-$3.98 \times 10^{-4})$] substitutions per-site-per-year, a value that overlaps the range of estimates

327    obtained from the randomized sets, reinforcing that the observed data does not have a strong

328    temporal signal (Rieux & Balloux, 2016).Therefore, considering the weak clock-like structure of this

329    dataset, Bayesian molecular clock and evolutionary analyses were not performed.

330    **4. Discussion**

331    This work is the first study that applied WGS approaches to *M. bovis* from Portugal. The aim was to

332    explore the fine-scale genomic signatures of field isolates and to highlight the value of phylogenetic

333    inference and transmission networks approaches to understand the history of this pathogen at the

334    livestock-wildlife interface.

335    In the current research, the SNP-based phylogeny established five main clades, with clades C, D and

336    E presenting lower levels of intra-clade diversity and clade-specific monomorphic SNP sites that can

337    be explored as phylogenetic markers in future diagnostic and epidemiological studies. The topology of

338    the SNP-based phylogenetic tree agrees with the dendrogram obtained by combining classical

339    genotyping methods, being evident that the partition between clades A-D and clade E remains well-

340    structured and that the members of clades D and E remain together. Similar findings have been

341    reported by other works that register clustering agreement between spoligotyping and WGS (Hauer et

342  al., 2019; Lasserre et al., 2018). For five strains, the spoligotyping profile obtained by the reverse-

343  hybridization method in the wet lab is not the same as the profile determined *in silico* through the

344  vSNP pipeline, a finding that has also been reported by others (Hijikata et al., 2017; Maeda et al.,

345  2020). This mismatch might contribute to discrepancies in phylogenetic trees topology generated with

346  classical genotyping methods and the SNP data.

347  The SNP-based phylogeny revealed geographic clustering to some extent, with clade C confined to

348  Castelo Branco and clade E being only present in Castelo Branco and Portalegre. These results are in

349  agreement with previous work performed on this *M. bovis* population from the same geographical

350  settings, that was based on the molecular characterization of isolates by classical genotyping methods

351  (Reis et al., 2020). This previous study suggested the existence of five *M. bovis* ancestral populations

352  with geographic specificities. Furthermore, a closer analysis into the temporal evolution of the SNP

353  network evidenced strong local dynamics within *M. bovis* strains from the same geographic region,

354  highlighting higher mean values of SNPs in common comparing with the links established between *M.*

355  *bovis* strains from different geographic regions. As expected, the maximum values of common shared

356  SNPs occurred between *M. bovis* from the bordering regions of Castelo Branco and Portalegre.

357  The global transmission network across regions in Portugal supports the information provided by the

358  phylogenetic tree, with *M. bovis* isolates grouping in the same SNP clade contributing to the same

359  transmission chains. The latency periods and long generation times of *M. bovis* can contribute to the

360  non-agreement between phylogeny and transmission trees. Moreover, works already performed

361  mainly with *M. tuberculosis* pointed to the variability in the number of SNPs shared between strains

362  recovered from individuals with epidemiological links with the differential number of SNPs being used

363  as cut-off to consider a transmission event (J. A. Guerra-Assunção et al., 2015). In the global and local

364  transmission trees concerning this *M. bovis* population from Portugal, links with SNP differences

365  inferior to the established cut-off value of three were identified, from the same or different host species

366  circulating in the same spatiotemporal context, therefore suggesting recent chains of transmission.

367  The Castelo Branco and Portalegre networks allowed a fine-scale detail of these events with the

368  identification of three situations where the individuals were originated from different officially delimited

369  hunting areas and herds, shedding light on the importance of animal mobility to the spread and

370  transmission of *M. bovis* in this multi-host system. The TB eradication program implemented in

371  Portugal is based on active and passive surveillance, but exclusively applied to the cattle population.

372    In contrast, surveillance in wildlife is exclusively passive and dependent upon irregular sanitary

373    evaluation of hunter-harvested animals, meaning that *M. bovis* infection and underlying transmission

374    might be established for long periods of time before sample collection and molecular characterization

375    of isolates occurs, therefore limiting transmission reconstruction inferences.

376    In the global transmission tree, there are three large transmission chains with origin in three strains

377    recovered from each one of the three TB hotspots, giving support to geographic clustering; plus,

378    different degrees of inferred infectivity could be associated to each region. Beja is the most

379    homogeneous one, presenting a lower mean value of SNP differences (244 SNPs) and a regular

380    inferred infectivity plot; while Castelo Branco and Portalegre presented more diverse *M. bovis*

381    populations, with higher mean SNP differences between strains ($n$= 376 for Castelo Branco and $n$=

382    405 for Portalegre) and with the inferred infectivity plots revealing high maximum values (six

383    secondary cases for Castelo Branco and nine for Portalegre).

384    Inferences from a previous molecular approach based on two types of genomic regions - direct

385    repeats and tandem repeats - (Reis et al., 2020) pointed the *M. bovis* population from Beja as the

386    most ancestral, when comparing with Castelo Branco and Portalegre and populations from the latter

387    as expanding populations. The high values of secondary cases reported in Castelo Branco and

388    Portalegre, although in a smaller proportion, indicate higher dissemination of the pathogen. Despite

389    the fact that this analysis only reflects the reality of these 44 cases, it mirrors the differences between

390    expanding and stable populations.

391    The lack of temporal signal indicated poor correlation between genetic divergence and sampling time,

392    stopping us to advance into evolutionary analyses. This limitation is common given the very slow

393    mutation rate of *M. bovis* and a larger and over-sampled isolate dataset would be required to make

394    any firm inferences regarding evolution. The  large genetic distance among some isolates in the global

395    and local transmission trees more likely reflect frequent introductions of new, more distantly related

396    lineages into the same geographic regions.

397    When considering host species, the transmission network registered recent transmission events

398    (including zero SNP differences) involving the same wildlife species or two wildlife host species,

399    suggesting the occurrence of intra- and inter-specific transmission events. The results from inferred

400    infectivity by host support the importance of both wildlife species in *M. bovis* transmission, with the red

401    deer dataset presenting a higher maximum value of inferred secondary cases than wild boar,

402    suggesting that red deer could exert a prominent role in this multi-host system. This finding is in line

403    with previous observational studies in Portugal (Vieira-Pinto et al., 2011) and with the characteristics

404    of the pathophysiology of TB in cervids that develop more open lesions and thus are more propense to

405    excrete the bacilli (Cunha et al., 2012). Consequently, focus on these two wild hosts should be

406    considered when designing new interventions aiming to improve control of animal TB.

407    **5. Concluding remarks**

408    There is clearly a need to characterise *M. bovis* lineages using clade-defining SNPs. Although limited

409    to a modest dataset, this work confirms and reinforces the value of WGS application to the study of *M.*

410    *bovis* transmission and persistence at the livestock-wildlife interface. The combination of WGS data

411    and epidemiological information provided insights into *M. bovis* demographic history in a multi-host

412    system, enabling the analyses of recent transmission events in several situations.

413    The knowledge of disease status and routes of transmission in wildlife are crucial to design and

414    implement effective control measures. The findings reported in this work contribute to support the idea

415    that eradication actions in the wildlife population are increasingly necessary, particularly when one

416    considers infectivity by host species. Altogether, we provide quantitative evidence that future control

417    measures in livestock production systems must not ignore wildlife related parameters, such as

418    abundance, behaviour and interaction with livestock, with the possibility for a differential approach

419    regarding red deer and wild boar. Furthermore, our fine scale molecular analyses suggest that Castelo

420    Branco and Portalegre are to be considered priority areas of research and intervention and that

421    adjacent livestock populations are to be tested more frequently.

422    Future WGS studies with a larger dataset from Castelo Branco and Portalegre, and from a broader

423    time period, could help resolve the transmission networks and potentially provide stronger temporal

424    signal, enabling evolutionary analyses, with estimation of evolutionary parameters, such as

425    substitution rates, the probability of host species transition, time to the most recent common ancestor

426    and timescales for clade divergence. Whole genome sequencing of *M. bovis* from across the world is

427    utterly needed, not only to enlighten epidemiological scenarios, but also to build experience and tools

428    to deal with the characteristic lack of temporal signal when slow evolving, latent microorganisms are

429    involved. Most algorithms and simulation tools were originally developed for the evolutionary analyses

430    of fast-evolving microorganisms (e.g. RNA viruses) and often are inadequate to deal with molecular

431    data from highly clonal, monomorphic organisms. Future studies with a larger dataset could give

432 further support to the use of SNP monomorphic sites as phylogenetic markers in settings where WGS

433 may not be easy to be implemented and also contribute to understand pathogen adaptation processes

434 to hosts.

435 **6. Acknowledgements**

448

449 **7.    Author Contributions**

450 MVC conceived the study. AB and TA provided the *M. bovis* isolates from the NRL. ACR performed

451 the experimental work. SRA contributed with isolate sequencing and quality control of genome

452 sequences. ACR, LS, and MVC thoroughly analysed the data. ACR and MVC wrote the manuscript.

453 LS and RT contributed with critical discussion. All authors approved the final manuscript.

454

455 **8.    Data Availability Statement**

456 Data sharing will be granted by the corresponding author upon reasonable request.

457

458 **9.    Conflict of Interest Statement**

459 The authors declare that no competing interests exist.

460

461 **10. References**

17

462    Barasona, J. A., Torres, M. J., Aznar, J., Gortázar, C., & Vicente, J. (2017). DNA Detection Reveals
463        Mycobacterium tuberculosis Complex Shedding Routes in Its Wildlife Reservoir the Eurasian
464        Wild    Boar.    *Transboundary    and    Emerging    Diseases*,    *64*,    906–915.
465        https://doi.org/10.1111/tbed.12458

466    Bastian, M., & Heymann, S. (2009). Gephi: An Open Source Software for Exploring and Manipulating
467        Networks. In *International AAAI Conference on Weblogs and Social Media*.

468    Biek, R., O'Hare, A., Wright, D., Mallon, T., Mccormick, C., Orton, R. J., … Kao, R. R. (2012). Whole
469        Genome Sequencing Reveals Local Transmission Patterns of Mycobacterium bovis in Sympatric
470        Cattle    and    Badger    Populations.    *PLoS    Pathogens*,    *8*,    1003008.
471        https://doi.org/10.1371/journal.ppat.1003008

472    Brites, D., Loiseau, C., Menardo, F., Borrell, S., Boniotti, M. B., Warren, R., … Gagneux, S. (2018). A
473        New Phylogenetic Framework for the Animal-Adapted Mycobacterium tuberculosis Complex.
474        *Frontiers in Microbiology*, *9*, 2820. https://doi.org/10.3389/fmicb.2018.02820

475    Corner, L. A. L., Murphy, D., & Gormley, E. (2011). Mycobacterium bovis Infection in the Eurasian
476        Badger (Meles meles): the Disease, Pathogenesis, Epidemiology and Control. *Journal of
477        Comparative Pathology*, *144*, 1–24. https://doi.org/10.1016/j.jcpa.2010.10.003

478    Crispell, J., Benton, C. H., Balaz, D., De Maio, N., Ahkmetova, A., Allen, A., … Kao, R. R. (2019).
479        Combining genomics and epidemiology to analyse bi-directional transmission of mycobacterium
480        bovis in a multi-host system. *ELife*, *8*, 45833. https://doi.org/10.7554/eLife.45833

481    Crispell, J., Cassidy, S., Kenny, K., McGrath, G., Warde, S., Cameron, H., … Gordon, S. V. (2020).
482        Mycobacterium bovis genomics reveals transmission of infection between cattle and deer in
483        Ireland. *Microbial Genomics*, *6*(8), 1–8. https://doi.org/10.1099/mgen.0.000388

484    Crispell, J., Zadoks, R. N., Harris, S. R., Paterson, B., Collins, D. M., De-Lisle, G. W., … Price-carter,
485        M. (2017). Using whole genome sequencing to investigate transmission in a multi-host system:
486        bovine tuberculosis in New Zealand. *BMC Genomics*, *18*, 180. https://doi.org/10.1186/s12864-
487        017-3569-x

488    Cunha, M. V., Matos, F., Canto, A., Albuquerque, T., Alberto, J. R., Aranha, J. M., … Botelho, A.
489        (2012). Implications and challenges of tuberculosis in wildlife ungulates in Portugal: A molecular
490        epidemiology    perspective.    *Research    in    Veterinary    Science*,    *92*,    225–235.
491        https://doi.org/10.1016/j.rvsc.2011.03.009

492    Depristo, M., Banks, E., Poplin, R., Garimella, K. V, Maguire, J., Hartl, C., … Daly, M. (2011). A
493        framework for variation discovery and genotyping using next-generation DNA sequencing data.
494        *Nature Genetics*, *43*(5), 491–498. https://doi.org/10.1038/ng.806.A

495    Dray, S., & Dufour, A.-B. (2007). The ade4 Package: Implementing the Duality Diagram for Ecologists.
496        *Journal of Statistical Software*, *22*(4), 1–20. https://doi.org/10.18637/jss.v022.i04

497   Duarte, E. L., Domingos, M., Amado, A., Cunha, M. V., & Botelho, A. (2010). MIRU-VNTR typing adds
498       discriminatory value to groups of Mycobacterium bovis and Mycobacterium caprae strains
499       defined    by    spoligotyping.    *Veterinary    Microbiology*,    *143*,    299–306.
500       https://doi.org/10.1016/j.vetmic.2009.11.027

501   Fitzgerald, S. D., & Kaneene, J. B. (2012). Wildlife Reservoirs of Bovine Tuberculosis Worldwide:
502       Hosts, Pathology, Surveillance, and Control. *Veterinary Pathology*, *50*, 488–499.
503       https://doi.org/10.1177/0300985812467472

504   Gagneux, S. (2018, April 1). Ecology and evolution of Mycobacterium tuberculosis. *Nature Reviews*
505       *Microbiology*. Nature Publishing Group. https://doi.org/10.1038/nrmicro.2018.8

506   Glaser, L., Carstensen, M., Shaw, S., Robbe-Austerman, S., Wunschmann, A., Grear, D., …
507       Thomsen, B. (2016). Descriptive Epidemiology and Whole Genome Sequencing Analysis for an
508       Outbreak of Bovine Tuberculosis in Beef Cattle and White-Tailed Deer in Northwestern
509       Minnesota. *PLoS ONE*, *11*, e0145735. https://doi.org/10.1371/journal.pone.0145735

510   Gortázar, C., Torres, M. J., Vicente, J., Acevedo, P., Reglero, M., de la Fuente, J., … Aznar-Martín, J.
511       (2008). Bovine tuberculosis in Doñana Biosphere Reserve: The role of wild ungulates as disease
512       reservoirs   in   the   last   Iberian   lynx   strongholds.   *PLoS   ONE*,   *3*,   e2776.
513       https://doi.org/10.1371/journal.pone.0002776

514   Guerra-Assunção, J. A., Houben, R. M. G. J., Crampin, A. C., Mzembe, T., Mallard, K., Coll, F., …
515       Glynn, J. R. (2015). Recurrence due to Relapse or Reinfection With Mycobacterium tuberculosis:
516       A Whole-Genome Sequencing Approach in a Large, Population- Based Cohort With a High HIV
517       Infection Prevalence and Active Follow-up. *Journal of Infectious Diseases*, *211*, 1154–1163.
518       https://doi.org/10.1093/infdis/jiu574

519   Guerra-Assunção, J., Crampin, A., Houben, R., Mzembe, T., Mallard, K., Coll, F., … Glynn, J. (2015).
520       Large-scale whole genome sequencing of M . tuberculosis provides insights into transmission in
521       a high prevalence area. *ELife*, *4*, e05166. https://doi.org/10.7554/eLife.05166

522   Hauer, A., Michelet, L., Cochard, T., Branger, M., Nunez, J., Boschiroli, M. L., & Biet, F. (2019).
523       Accurate Phylogenetic Relationships Among Mycobacterium bovis Strains Circulating in France
524       Based on Whole Genome Sequencing and Single Nucleotide Polymorphism Analysis. *Frontiers*
525       *in Microbiology*, *10*, 955. https://doi.org/10.3389/fmicb.2019.00955

526   Hijikata, M., Keicho, N., Duc, L. Van, Maeda, S., Hang, N. T. Le, Matsushita, I., & Kato, S. (2017).
527       Spoligotyping and whole-genome sequencing analysis of lineage 1 strains of Mycobacterium
528       tuberculosis   in   Da   Nang,   Vietnam.   *PLoS   ONE*,   *12*,   e0186800.
529       https://doi.org/10.1371/journal.pone.0186800

530   Jombart, T., Eggo, R., Dodd, P., & Balloux, F. (2011). Reconstructing disease outbreaks from genetic
531       data : a graph approach. *Heredity*, *106*(2), 383–390. https://doi.org/10.1038/hdy.2010.78

532   Kumar, S., Stecher, G., & Tamura, K. (2016). MEGA7 : Molecular Evolutionary Genetics Analysis

Version 7.0 for Bigger Datasets. *Molecular Biology and Evolution*, *33*, 1870–1874. https://doi.org/10.1093/molbev/msw054

Lasserre, M., Fresia, P., Greif, G., Iraola, G., Castro-Ramos, M., Juambeltz, A., … Berná, L. (2018). Whole genome sequencing of the monomorphic pathogen Mycobacterium bovis reveals local differentiation of cattle clinical isolates. *BMC Genomics*, *19*(2), 1–14. https://doi.org/10.1186/s12864-017-4249-6

Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows – Wheeler transform. *Bioinformatics*, *25*, 1754–1760. https://doi.org/10.1093/bioinformatics/btp324

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … Durbin, R. (2009). The Sequence Alignment / Map format and SAMtools. *Bioinformatics*, *25*(16), 2078–2079. https://doi.org/10.1093/bioinformatics/btp352

Maeda, S., Hijikata, M., Hang, N. T. Le, Thoung, P. H., Huan, H. Van, Hoang, N. P., … Keicho, N. (2020). Genotyping of Mycobacterium tuberculosis spreading in Hanoi, Vietnam using conventional and whole genome sequencing methods. *Infection, Genetics and Evolution*, *78*, 104107. https://doi.org/10.1016/j.meegid.2019.104107

Mckenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., … Depristo, M. A. (2010). The Genome Analysis Toolkit□: A MapReduce framework for analyzing next-generation DNA sequencing data sequencing data. *Genome Research*, *20*, 1297–1303. https://doi.org/10.1101/gr.107524.110.20

Naranjo, V., Gortázar, C., Vicentea, J., & de la Fuente, J. (2008). Evidence of the role of European wild boar as a reservoir of Mycobacterium tuberculosis complex. *Veterinary Microbiology*, *127*, 1–9. https://doi.org/10.1016/j.vetmic.2007.10.002

Palmer, M. V. (2007). Tuberculosis: A reemerging disease at the interface of domestic animals and wildlife. *Current Topics in Microbiology and Immunology*, *315*, 195–215. https://doi.org/10.1007/978-3-540-70962-6_9

Palmer, Mitchell V., Thacker, T. C., Waters, W. R., Gortázar, C., & Corner, L. a L. (2012). Mycobacterium bovis: A model pathogen at the interface of livestock, wildlife, and humans. *Veterinary Medicine International*, *2012*, 236205. https://doi.org/10.1155/2012/236205

Paradis, E., Claude, J., & Strimmer, K. (2004). APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics*, *20*, 289–290. https://doi.org/10.1093/bioinformatics/btg412

Price-Carter, M., Brauning, R., de Lisle, G. W., Livingstone, P., Neill, M., Sinclair, J., … Collins, D. M. (2018). Whole Genome Sequencing for Determining the Source of Mycobacterium bovis Infections in Livestock Herds and Wildlife in New Zealand. *Frontiers in Veterinary Science*, *5*, 272. https://doi.org/10.3389/fvets.2018.00272

Rambaut, A., Lam, T. T., Carvalho, L. M., & Pybus, O. G. (2016). Exploring the temporal structure of

568       heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evolution*, *2*, vew007.

569       https://doi.org/10.1093/ve/vew007

570 Reis, A. C., Tenreiro, R., Albuquerque, T., Botelho, A., & Cunha, M. V. (2020). Long-term molecular

571       surveillance provides clues on a cattle origin for Mycobacterium bovis in Portugal. *Scientific*

572       *Reports*, *10*(1), 1–18. https://doi.org/10.1038/s41598-020-77713-8

573 Rieux, A., & Balloux, F. (2016). Inferences from tip-calibrated phylogenies : a review and a practical

574       guide. *Molecular Ecology*, *25*, 1911–1924. https://doi.org/10.1111/mec.13586

575 Rodriguez-Campos, S., Schürch, A. C., Dale, J., Lohan, A. J., Cunha, M. V., Botelho, A., … Smith, N.

576       H. (2012). European 2 – A clonal complex of Mycobacterium bovis dominant in the Iberian

577       Peninsula. *Infection, Genetics and Evolution*, *12*(4), 866–872.

578       https://doi.org/10.1016/j.meegid.2011.09.004

579 Salvador, L., O'Brien, D., Cosgrove, M., Stuber, T., Schooley, A., Crispell, J., … Kao, R. (2019).

580       Disease management at the wildlife-livestock interface: using whle-genome sequencing to study

581       the role of elk in Mycobacterium bovis transmission in Michigan, USA. *Molecular Ecology*, *28*,

582       2192–2205. https://doi.org/10.1111/mec.15061

583 Santos, N., Correia-Neves, M., Ghebremichael, S., Källenius, G., Svenson, S. B., & Almeida, V.

584       (2009). Epidemiology of Mycobacterium bovis infection in wild boar (Sus scrofa) from Portugal.

585       *Journal of Wildlife Diseases*, *45*, 1048–1061. https://doi.org/10.7589/0090-3558-45.4.1048

586 Thorvaldsdóttir, H., Robinson, J. T., & Mesirov, J. P. (2012). Integrative Genomics Viewer (IGV): high-

587       performance genomics data visualization and exploration. *Briefings in Bioinformatics*, *14*, 178–

588       192. https://doi.org/10.1093/bib/bbs017

589 To, T.-H., Jung, M., Lycett, S., & Gascuel, O. (2015). Fast Dating Using Least-Squares Criteria and

590       Algorithms. *Systematics Biology*, *65*, 82–97. https://doi.org/10.1093/sysbio/syv068

591 Trewby, H., Wright, D., Breadon, E. L., Lycett, S. J., Mallon, T. R., Mccormick, C., … Kao, R. R.

592       (2016). Use of bacterial whole-genome sequencing to investigate local persistence and spread in

593       bovine tuberculosis. *Epidemics*, *14*, 26–35.

594       https://doi.org/dx.doi.org/10.1016/j.epidem.2015.08.003

595 Van der Auwera, G., Carneiro, M. O., Hartl, C., Poplin, R., del Angel, G., Levy-moonshine, A., …

596       Depristo, M. A. (2014). From FastQ data to high confidence variant calls: the Genome Analysis

597       Toolkit best practices pipeline. *Current Protocols Bioinformatics*, *43*, 11.10.1-11.10.33.

598       https://doi.org/10.1002/0471250953.bi1110s43.From

599 Vieira-Pinto, M., Alberto, J., Aranha, J., Serejo, J., Canto, A., Cunha, M. V., & Botelho, A. (2011).

600       Combined evaluation of bovine tuberculosis in wild boar (Sus scrofa) and red deer (Cervus

601       elaphus) from Central-East Portugal. *European Journal of Wildlife Research*, *57*, 1189–1201.

602       https://doi.org/10.1007/s10344-011-0532-z

603

604    **Table 1.** Number of *M. bovis* strains within each SNP clade (A to E) and identification of clade-

605    defining and clade-monomorphic SNP sites.

| SNP clade | Total SNP sites | Clade-defining SNP sites[a] | Clade-monomorphic SNP sites[b] |
|---|---|---|---|
| A (*n*=14) | 622 | 108 | - |
| B (*n*=11) | 611 | 133 | - |
| C (*n*=3) | 320 | 184 | 49 |
| D (*n*=6) | 372 | 217 | 82 |
| E (*n*=10) | 431 | 360 | 352 |
| A to D (*n*=34) | 1419 | 1411 | 106 |

606    (a) Polymorphic positions present only in the clade-members.
607    (b) Polymorphic positions present only in the clade-members and common to all members.

608    **Figures**

609



610

611    **Figure 1.** Phylogeny of *M. bovis* field strains from TB hotspots in Portugal. (A) Maximum Likelihood

612    Tree using GTR model with input taken as an alignment file containing only informative and validated

613    SNPs. The tree is drawn to scale, with branch lengths measured in the number of substitutions per

614    site. (B) UPGMA tree, applying categorical option as a similarity coefficient, with input taken as the

615    combined dataset based on spoligotyping and 8-*loci* MIRU-VNTR data. Colours identify the different

616    SNP clades (A – dark cyan, B – orange, C – purple, D – pink and E – yellow).
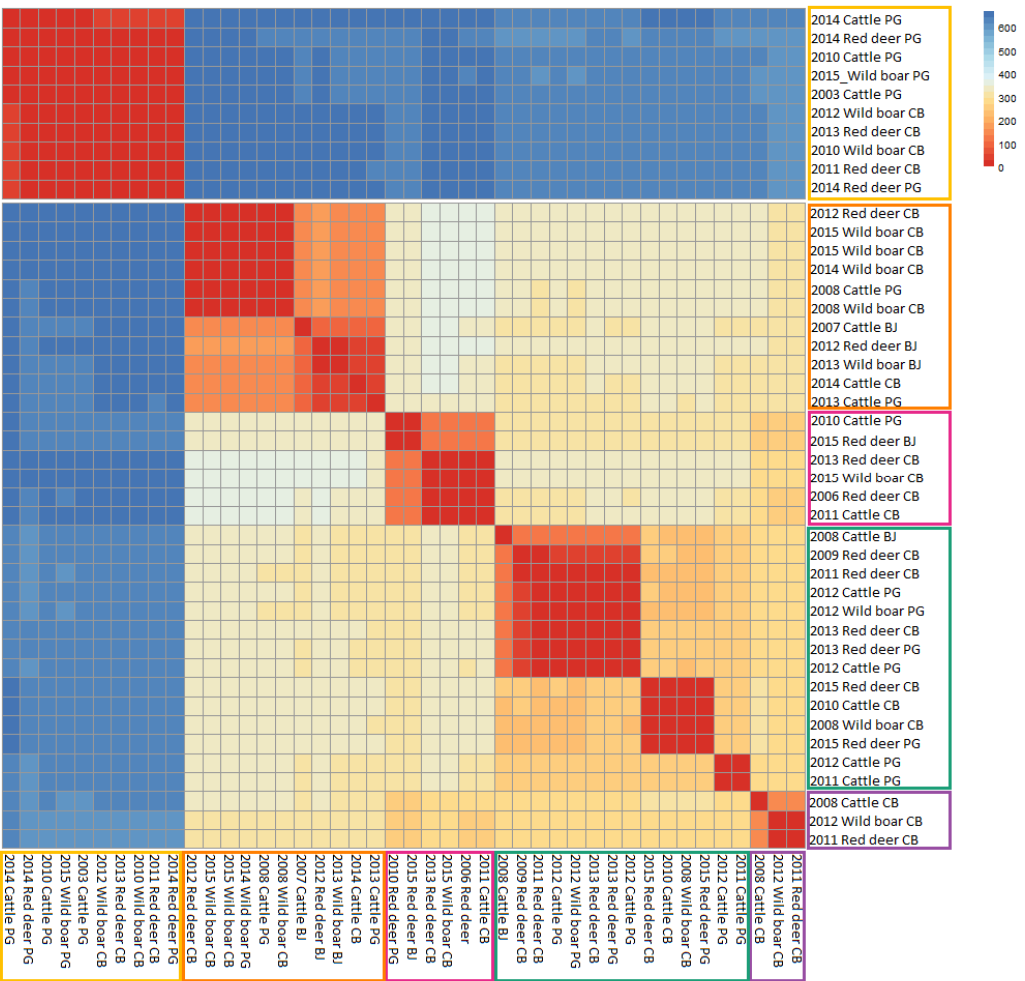
617

618

619

620

**Figure 2.** Heatmap of pairwise SNP distances based on the absolute differences of SNPs. The colours of the boxes identify the different SNP clades (A –dark cyan, B – orange, C – purple, D – pink and E – yellow).
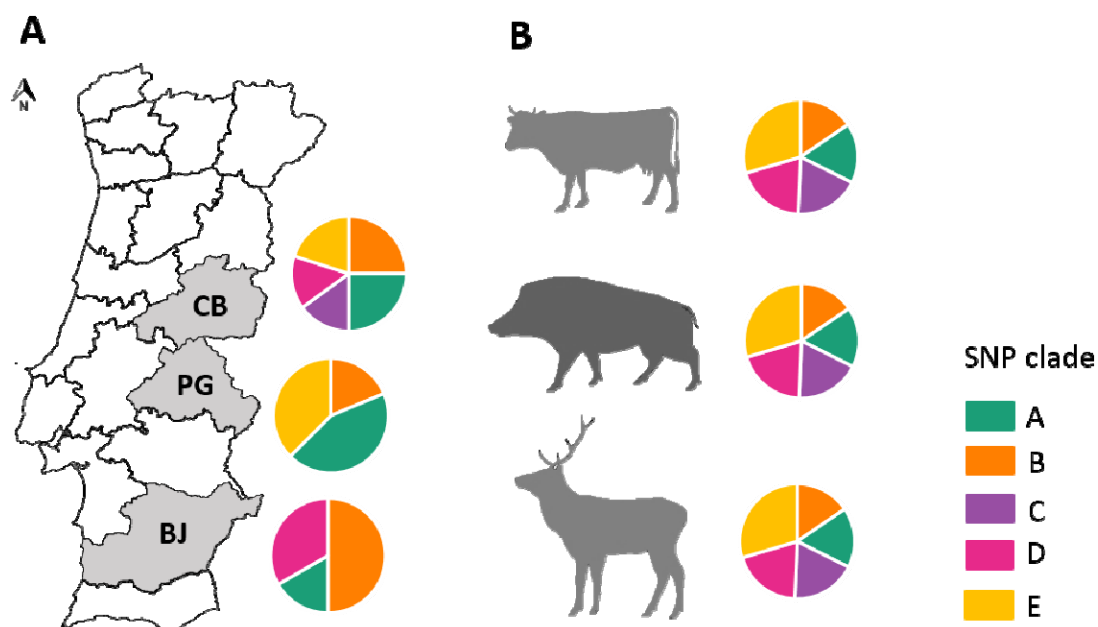
25

**Figure 3.** Distribution of SNP clades by geographic region (A) [Beja (BJ), Castelo Branco (CB) and Portalegre (PG)] and host species (B). The colours of the pie charts identify the different SNP clades (A –dark cyan, B – orange, C – purple, D – pink and E – yellow).
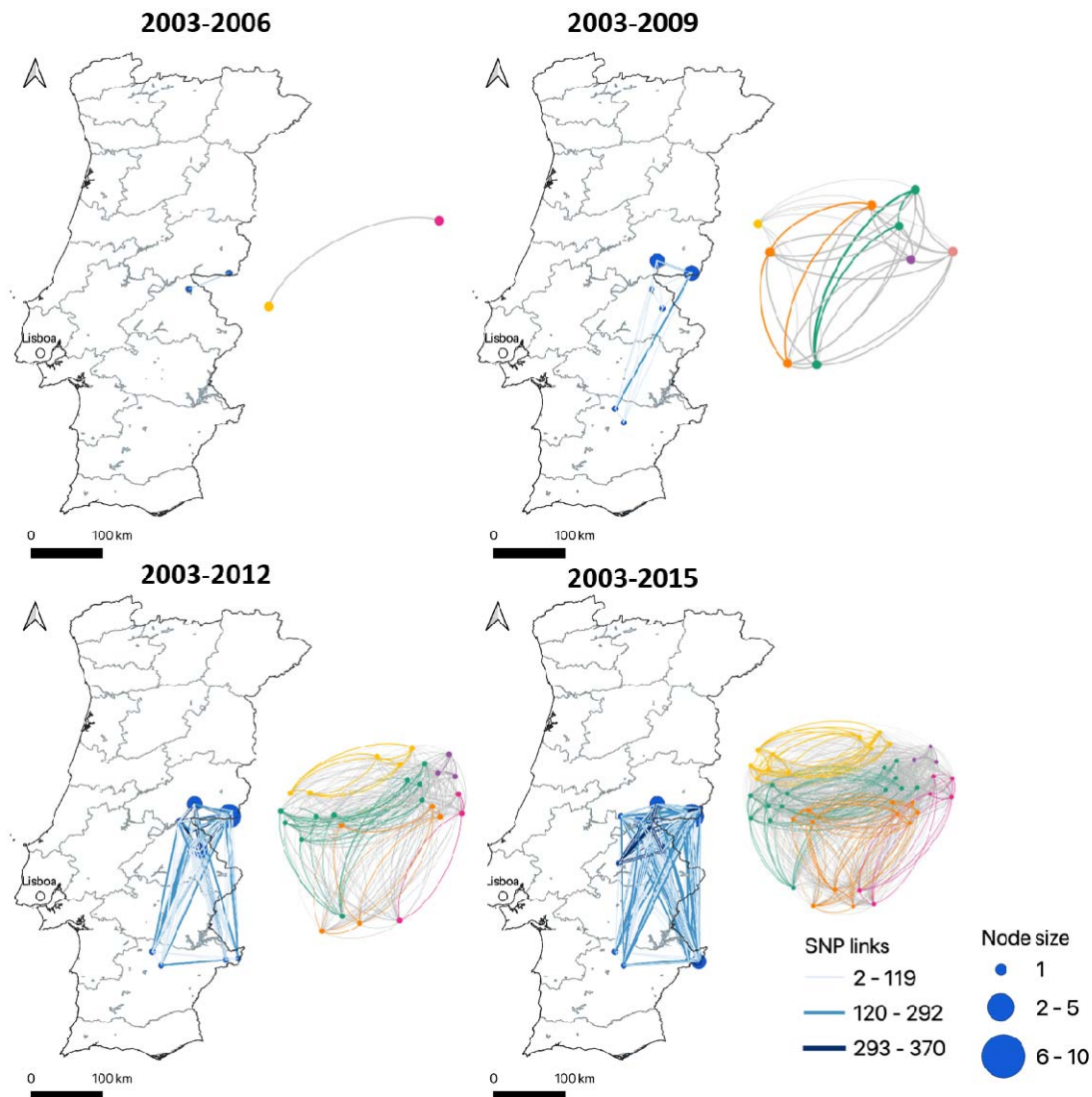
**Figure 4.** Temporal evolution of the SNP network established for the 44 *M. bovis* isolates recovered between 2003 and 2015. Four progressive time periods were considered according with the epidemiological scenario in Portugal. In the map of Portugal, the nodes represent *M. bovis* strains and the complexity of connections is based on the number of shared SNPs between strains. The side network evidences the relation established between *M. bovis* strains grouped within the same clade – each node represents one *M. bovis* strain and the colours identify the connections according with SNP clade (A – dark cyan, B – orange, C – purple, D – pink and E – yellow and grey for connections between different clades).
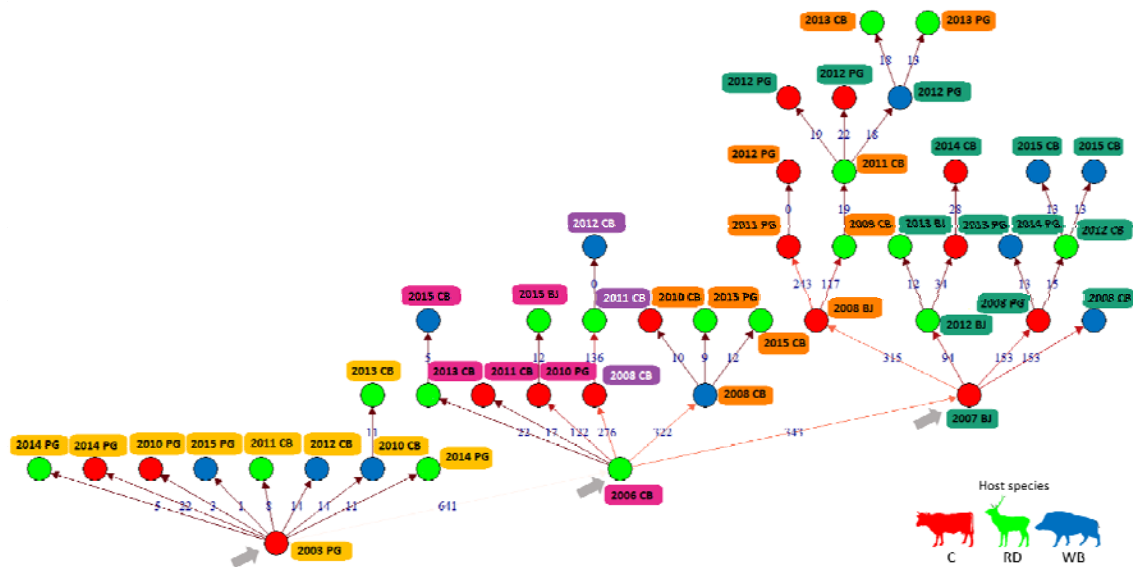
**Figure 5.** Transmission tree: reconstruction of transmission chains between cattle, red deer and wild boar, based on the SNP analyses, the year of isolation, host species and geographic region [Beja (BJ), Castelo Branco (CB) and Portalegre (PG)] of TB positive animals. The nodes are coloured by host species (cattle – red, red deer – green and wild boar – blue) and the established connections are based in SNP differences between *M. bovis* strains. Grey arrow point to the beginning of the three major branches.
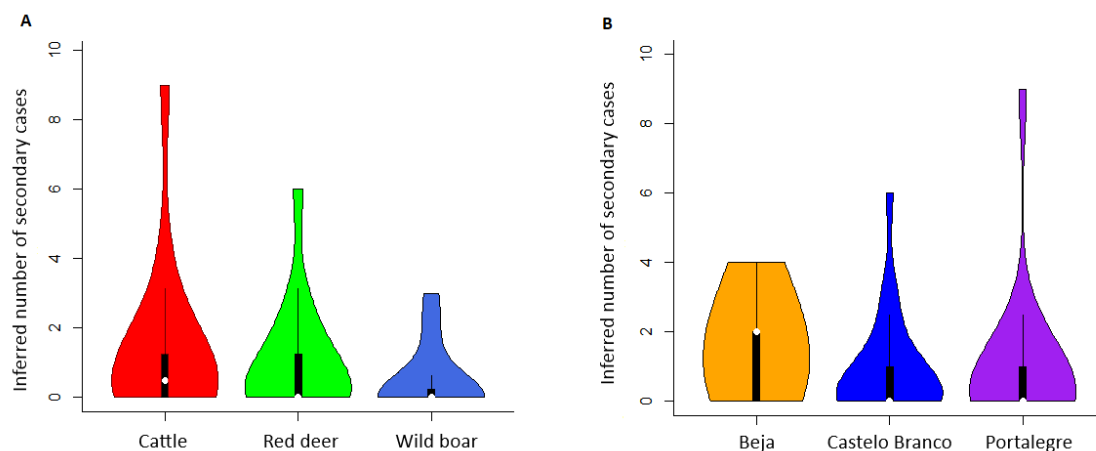
656

657



658
659     **Figure 6.** Violin plots of inferred number of secondary cases grouped by host species (A) and

660     geographic region (B).

661

662

**Figure 7**. Transmission tree: reconstruction of transmission chains between cattle, red deer and wild boar, based on the SNP analyses, the year of isolation and host species in Castelo Branco (A) and Portalegre (B) of TB positive animals. The nodes are coloured by host species (cattle – red, red deer – green and wild boar – blue) and the established connections are based in SNP differences between *M. bovis*. The apostrophe numbers identify the recent transmission events (in transmission tree A, the event number "2" identifies a zero SNP link).