1 **Inferring *Mycobacterium bovis* transmission between cattle and badgers using**

2 **isolates from the Randomised Badger Culling Trial**

3

4 Andries J. van Tonder[1,*], Mark Thornton[1], Andrew J.K. Conlan[1], Keith A. Jolley[2], Lee

5 Goolding[3], Andrew P. Mitchell[3], James Dale[3], Eleftheria Palkopoulou[3], Philip J.

6 Hogarth[3], R. Glyn Hewinson[4], James L.N. Wood[1], Julian Parkhill[1]

7

8 [1]Department of Veterinary Medicine, University of Cambridge, Cambridge, United

9 Kingdom

10 [2]Department of Zoology, University of Oxford, Oxford, United Kingdom

11 [3]Animal and Plant Health Agency, New Haw, United Kingdom

12 [4] IBERS, Aberystwyth University, Ceredigion, United Kingdom

13

14 *Corresponding author (ajv37@cam.ac.uk)

15

16 **Abstract**

17 *Mycobacterium bovis* (*M. bovis)* is a causative agent of bovine tuberculosis, a

18 significant source of morbidity and mortality in the global cattle industry. The

19 Randomised Badger Culling Trial was a field experiment carried out between 1998

20 and 2005 in the South West of England. As part of this trial, *M. bovis* isolates were

21 collected from contemporaneous and overlapping populations of badgers and cattle

22 within ten defined trial areas. We combined whole genome sequences from 1,442

23 isolates with location and cattle movement data, identifying transmission clusters and

24 inferred rates and routes of transmission of *M. bovis*. Most trial areas contained a

25    single transmission cluster that had been established shortly before sampling, often

26    contemporaneous with the expansion of bovine tuberculosis in the 1980s. The

27    estimated rate of transmission from badger to cattle was approximately two times

28    higher than from cattle to badger, and the rate of within-species transmission

29    considerably exceeded these for both species. We identified long distance

30    transmission events linked to cattle movement, recurrence of herd breakdown by

31    infection within the same transmission clusters and superspreader events driven by

32    cattle but not badgers. Overall, our data suggests that the transmission clusters in

33    different parts of South West England that are still evident today were established by

34    long-distance seeding events involving cattle movement, not by recrudescence from

35    a long-established wildlife reservoir. Clusters are maintained primarily by within-

36    species transmission, with less frequent spill-over both from badger to cattle and

37    cattle to badger.

38

39    **Introduction**

40    *Mycobacterium bovis* (*M. bovis*), a member of the *Mycobacterium tuberculosis*

41    complex (MTBC) and a pathogen with zoonotic potential [1], is the main causative

42    agent of bovine tuberculosis (bTB), a significant source of morbidity and mortality in

43    the global cattle industry.  In the United Kingdom (UK), the estimated annual cost of

44    managing this disease is £120 million [2].

45

46    *M. bovis* has a broad host range with different wildlife reservoirs depending on

47    geographic location: in Britain and Ireland the Eurasian badger is the predominant

48    wildlife host, in France wild boar and deer, and in New Zealand the introduced brush-

49 tail possum [3-5]. The presence of wildlife reservoirs makes the control and potential

50 elimination of bTB challenging even in countries such as the UK with extensive cattle

51 test and slaughter strategies, and movement restrictions imposed on herds with new

52 bTB incidents termed breakdowns [6].

53

54 The Randomised Badger Culling Trial (RBCT) was a large-scale ecological field

55 experiment carried out between 1998 and 2005 with the aim of quantifying the impact

56 of culling badgers on the incidence of bTB breakdowns in nearby cattle herds [7]. Ten

57 trial areas within the southwest of England and English Midlands, each of

58 approximately 100 km$^2$, were selected on the basis of high bTB incidence. Each trial

59 area was divided into triplets of randomly allocated interventions: proactive culling

60 (widespread and repeated culling across the trial areas), reactive culling (badgers

61 culled if breakdowns detected in nearby herds) and control or survey-only areas (no

62 badger culling). Approximately 9,000 badgers were culled and sampled in proactive

63 areas between 1998 and 2005 though culling was suspended between May 2001 and

64 January 2002 due to a national foot and mouth disease epidemic.

65

66 A number of previous studies have established epidemiological links between

67 badgers and nearby cattle although extent of transmission between the two host

68 species remains uncertain [8, 9]. More recent analyses making use of whole genome

69 sequencing (WGS), which offers much higher resolution for strain characterisation

70 and tracking transmission, have confirmed the close genetic relatedness of *M. bovis*

71 isolates from sympatric cattle and badger populations but, due to the low genomic

72 variability of the *M. bovis* genome and a lack of balanced sampling between the

73  different host species, have not been able to adequately address the direction of

74  transmission [10, 11].  The first direct estimate of the extent and directionality of

75  transmission between cattle and badgers suggested that transmission was up to ten

76  times higher from badgers to cattle than *vice versa* [9].  Subsequent studies have

77  estimated that cattle to badger transmission was at least three times or an order of

78  magnitude higher than badger to cattle transmission [12, 13].  However, these results

79  may not be applicable to the wider *M. bovis* population in different regions of the UK,

80  often being based on small, geographically narrow datasets chosen for the presence

81  of the same strain type (spoligotype SB0263). Those in Northern Ireland may

82  additionally reflect the lower density of badgers compared to Southern England and

83  the outbreak in Cumbria was an outbreak in a region with low incidence of bTB.

84

85  This Eradication of bovine tuberculosis (ERADbTB) project was set up with the aim of

86  using WGS data obtained from *M. bovis* isolates collected as part of the RBCT to

87  characterise the population structure of the bacterium within the trial areas, attempt

88  to quantify levels and directionality of *M. bovis* transmission between cattle and

89  badgers and track the longer-term persistence of genetic lineages of the bacterium.

90  Approximately 2,000 *M. bovis* isolates available from the RBCT were selected for

91  sequencing with the final dataset consisting of 1,442 genomes (690 from badgers

92  and 750 from cattle found to be infected in proactive cull trial areas respectively).

93

## Methods

*Sample selection, culturing and sequencing*

A total of 2,137 *M. bovis* isolates from cattle (n = 1,011) and badgers (n = 1,126) collected from proactive trial areas were selected for culturing, of which 1,838 isolates were located in the frozen archives maintained by the Animal and Plant Health Agency (APHA). Isolates were re-cultured and grown for up to six weeks or until sufficient growth was observed (n = 1,651). Isolates were heat killed in hot blocks at 80°C for 30 minutes. An adapted library construction protocol using an increased number of sixteen PCR cycles was used to generate Illumina libraries which were then sequenced at the Wellcome Sanger Institute using the Illumina HiSeq X10 platform to generate 2 x 150 bp paired-end reads. Metadata for the sequenced isolates is available on pubMLST (https://pubmlst.org/projects/mbovis-eradbtb) [14, 15]. A map of the geographical locations of isolate collection (latitude and longitude) was constructed using the R v 3.5.1 [16] library ggmap [17].


*Sequence QC*

FastQC v0.11.9 [18] was used to generate basic quality control metrics for the raw sequence data. Sequencing reads were prefiltered using Kraken v0.10.6 [19] against a database containing all RefSeq bacterial and archeal nucleotide sequences to identify reads with similarity to *Mycobacterium* species. Further sequence matching was done on the Kraken results using Bracken v1.0 [20]. Samples with < 70% reads mapping to a *Mycobacterium* species were excluded from further analyses (n = 183).

118 *In silico genotyping*

119 SpoTyping v2.0 [21] was used to extract the binary representation of spoligotype

120 patterns from the sequence reads and the *M. bovis* spoligotype database

121 (https://www.mbovis.org/database.php) was used to assign SB numbers. Novel

122 spoligotype patterns were submitted to the database to generate new SB numbers.

123 Clonal complexes were assigned to samples using RD-analyzer v1.0 [22] with

124 samples not identified as belonging to previously described clonal complexes (Eu1,

125 Eu2, Af1, Af2) designated as "Other" [23-26]. Further assignment of isolates marked

126 as "Other" to clonal complex was based on the phylogenetic lineages recently

127 identified by Loiseau *et al*. [27].

128

129 *Mapping and phylogenetics*

130 Sequence reads were mapped to the *Mycobacterium bovis* AF2122/97 reference

131 genome (NC0002945) using BWA mem v0.7.17 (minimum and maximum insert sizes

132 of 50 and 1000 respectively) [28]. Single nucleotide polymorphisms (SNPs) were

133 called using SAMtools v1.2 mpileup and BCFtools v1.2 (minimum base call quality of

134 50 and minimum root squared mapping quality of 30) as previously described [29].

135 Samples with reads mapping to less than 90% of the AF2122/97 reference were

136 excluded (n = 26). Genomic regions consisting of GC-rich sequences such as PPE

137 proteins and repeats were masked in the resulting alignment using previously

138 published coordinates [30] and variant sites in the subsequent masked alignment

139 were extracted using snp-sites v 2.5.1 [31]. Maximum likelihood phylogenetic trees

140 were constructed using IQ-tree v1.6.5 accounting for constant sites (-fconst;

141 determined using snp-sites -C) with the built-in model testing (-m MFP) to determine

142  the best phylogenetic model (GTR+F+R2) and 1000 ultrafast bootstraps (-bb 1000)

143  [32].  Pairwise SNP distances were calculated for all pairs of isolates from the SNP

144  alignment using pairsnp v1.0 (https://github.com/gtonkinhill/pairsnp).

145

146  To provide a global context for the isolates sequenced in this study, a published

147  clonal complex Eu1 dataset (n = 2,842; Supplementary File 2) spanning fourteen

148  countries was assembled [9, 10, 30, 33-47].   Sequence data was downloaded from

149  the European Nucleotide Archive (ENA) and trimmed using Trimmomatic v0.33 [48].

150  Sample QC, spoligotype assignment, mapping and phylogenetic tree construction

151  were performed as above.  The tree was rooted with a *Mycobacterium caprae* isolate

152  (SRR7617662).

153

154  *Transmission Clusters*

155  The R library iGRAPH [49] was used to define putative transmission clusters using a

156  pairwise SNP distance between any two samples of 15 as the threshold. This

157  threshold was chosen as it would allow for the possible identification of older

158  transmission events but also allow for any variance in the rates of mutation amongst

159  the sampled isolates, and has been previously used in a similar analysis of a human

160  *Mycobacterium tuberculosis* dataset [50]. Large clusters were manually divided

161  further on the basis of clear divisions within these clusters observed in the

162  phylogenetic tree (Clusters 5/6 and Clusters 8-12).  Transmission clusters with fewer

163  than 50 isolates were not analysed further leaving twelve transmission clusters for

164  further analyses.  New alignments were generated for each cluster as described

165  above.

166  The presence of a temporal signal in each transmission cluster was investigated by

167  plotting the root to tip distance for each isolate, calculated using the R library phytools

168  [51], against its sampling date (Supplementary Figure 1). The slope, x-intercept (most

169  recent common ancestor; MRCA), correlation coefficient and $R^2$ value were

170  calculated for each dataset in R. BEAST v1.8.4 [52] was run on each SNP alignment,

171  using tip sampling dates for calibration. Three runs of $10^8$ Markov chain Monte Carlo

172  (MCMC) iterations were performed using a HKY substitution model, strict or constant

173  molecular clock and constant or exponential population size and growth (12 separate

174  runs) for each transmission cluster. The performance of each model was assessed

175  through the comparison of posterior marginal likelihood estimates [53, 54] and the

176  model with the highest Bayes factor [55] (strict clock/constant population size) was

177  selected for each transmission cluster (Supplementary Table 1). The three selected

178  MCMC runs were combined using LogCombiner v1.8.4 (10% burnin) and

179  convergence was assessed (posterior effective sample size (ESS) > 200 for each

180  parameter). A maximum clade creditability tree summarizing the posterior sample of

181  trees in the combined MCMC runs was produced using TreeAnnotator v1.8.4. To

182  confirm the temporal signal in each tree generated, the R library TIPDATINGBEAST

183  [56] was used to resample tip dates from each alignment to generate 20 new datasets

184  with randomly assigned dates. BEAST was then run on each new dataset using the

185  same strict clock priors (Supplementary Figure 2). If the estimated substitution rates

186  in the observed data did not overlap with the estimated substitution rates in the

187  randomized data then the temporal signal observed in the observed data was

188  considered not to be obtained by chance.

189

190   Transmission reconstruction was performed on each cluster using the R library

191   TransPhylo [57] which allows for unsampled cases and within-host diversity.  The

192   same parameters (gamma shape = 1.6; scale = 3.5) were used for the infection and

193   generation time prior distribution.  The TransPhylo algorithm was run three times for

194   $10^7$ MCMC iterations sampling every 200,000 states and a burnin of 10% on each

195   cluster using the MCC trees generated previously.  The R library coda [58] was used

196   to assess convergence (Gelman and Rubin's Convergence Diagnostic < 1.05) and

197   ESS values > 100 for within-host diversity, reproductive rate and sampling proportion

198   (Supplementary Table 2). Post processing of each TransPhylo run was performed in

199   R.

200

201   The BEAST2 [59] package BASTA (Bayesian Structured coalescent Approximation)

202   [60] was used to estimate transmission rates between badgers and cattle, defined as

203   demes, in each transmission cluster.  A strict clock/equal population size model was

204   used and the BASTA analysis was repeated three times and run for $3 \times 10^8$ MCMC

205   iterations with 10% burnin.  Convergence was assessed as above.  Post processing

206   of the BASTA analysis was performed in R.

207

208   SNP-scaled phylogenetic trees were calculated for each transmission cluster using

209   pyjar (https://github.com/simonrharris/pyjar) [61] and plotted using the R libraries

210   treeio and ggtree [62, 63].

211

212

213

214 *Cattle Movements*

215 bTB metadata was extracted from APHA's Sam database which records all statutory

216 bTB testing information. Cattle movement metadata was extracted from APHA's

217 copy of the Department of Environment, Food and Rural Affairs' (DEFRA)'s Cattle

218 Tracing System (CTS). Movement data were extracted for 727/752 cattle where the

219 ear tag could be matched to the Sam database (it only became a legal requirement

220 to record cattle movement in the CTS after January 2001 so movement data may be

221 missing for the early part of the RBCT). Movements of TB test reactor cattle that were

222 not subjected to laboratory culture and/or sequencing of *M. bovis*, but may have

223 contributed to the spread of infection, were extracted from the CTS using the

224 following criteria: the animals passed through the same location as an animal with a

225 sequenced isolate, the animals were born before 2009 and the animals were

226 classified as "reactors". Animals were classified as reactors if they had a positive

227 tuberculin test result, had an inconclusive test result but were slaughtered and culture

228 positive for *M. bovis,* were culture positive for *M. bovis* following detection by routine

229 meat inspection at a slaughterhouse, or were culture-negative reactors that led to a

230 breakdown with other tuberculin test positive animals. UK grid coordinates were

231 extracted from the Sam database by matching to location IDs. Past breakdown

232 history was extracted by matching herds using county-parish-holding (CPH)

233 numbers. Where multiple herds had the same CPH number, the active dates of the

234 herds were checked and the individual animal test records were used to identify the

235 correct entries. Short stay locations and locations with missing coordinates were

236 excluded by creating animal records that removed missing locations or stays of fewer

237 than eight days. Where subsequent movements occurred, these were connected to

238    the previous movements to create a continuous record. Where the cattle ear tag IDs

239    of sequenced isolates could not be matched to the database, the CPH was used to

240    identify the final location and coordinates for plotting. The final herd of the animals

241    with sequenced isolates was determined as the location closest to death where the

242    length of stay was at least seven days. The data was queried and extracted from the

243    CTS using PostgreSQL. Pairwise geographic distances between each isolate in

244    kilometres were calculated using the distHaversine function from the R library

245    geosphere [64]. Herd and badger locations were randomly shifted by up to 1 km in

246    the horizontal and vertical planes for plotting using the R libraries maps and mapdata

247    [65].

248

249    **Results**

250    *Population structure*

251    A total of 1,442 *M. bovis* isolates from badgers (n = 690) and cattle (n = 752) were

252    sequenced and passed QC; the sites of collection for all 1,442 isolates are shown in

253    Figure 1A. The average number of sequenced isolates per trial area was 144 (range:

254    81-233) and the ratio of cattle to badger isolates varied from 0.22 (trial area D3) to

255    4.38 (trial area B2; Table 1). All sequenced isolates were collected between 1999 and

256    2010. The majority (1437/1442; 99.7%) of the isolates were clonal complex Eu1

257    whilst the remaining five isolates (all SB0134) belonged to an as yet undefined clonal

258    complex (labelled Unknown7 in Loiseau et al. [27]; Figure 1B).

259

260    Over 60 unique spoligotypes were identified with the most prevalent being SB0140

261    (n = 531), SB0263 (n = 491), SB0129 (n = 147), SB0274 (n = 85), SB0957 (n = 34) and
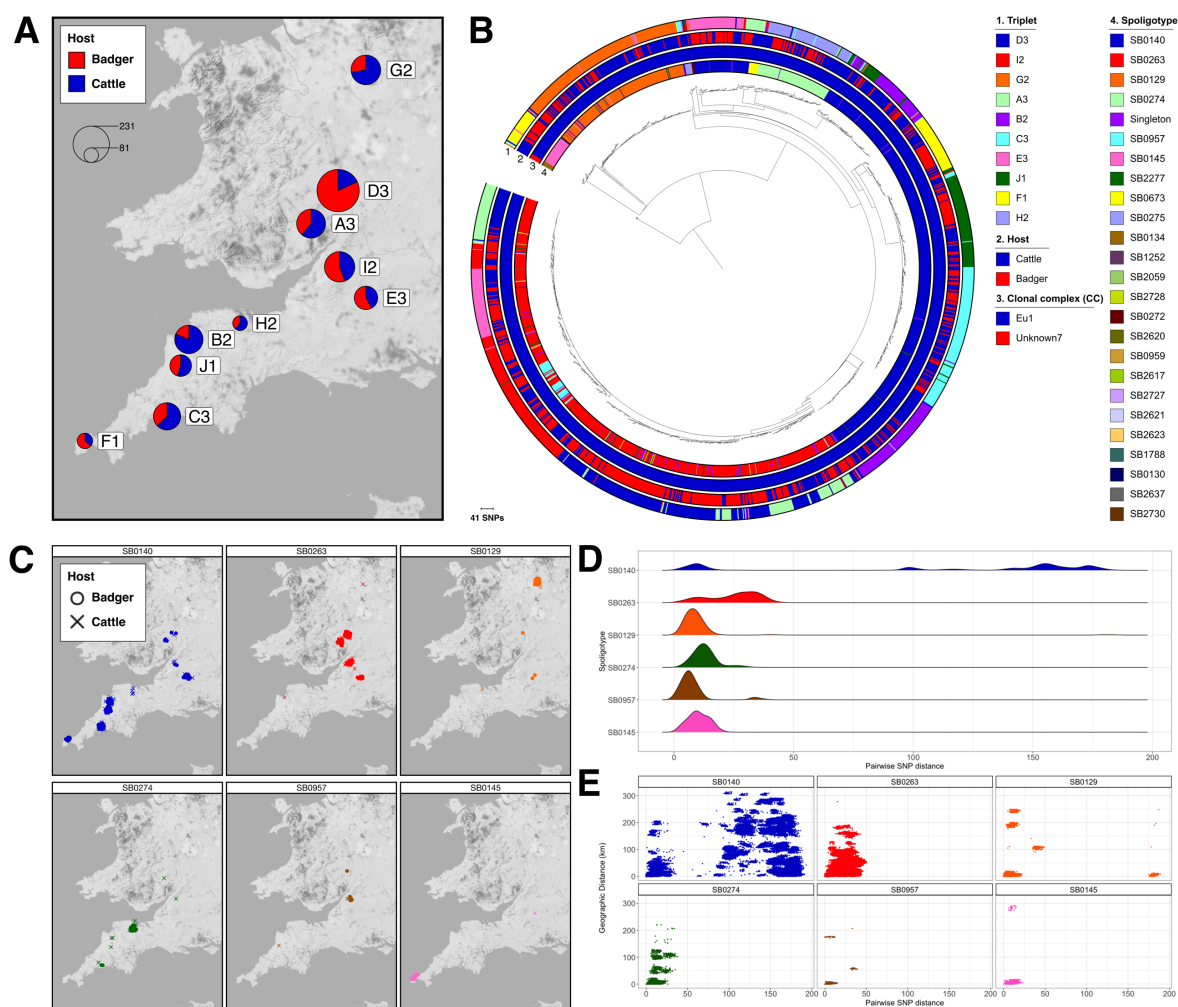
262    Table 1: Breakdown of the 1,442 sequenced *Mycobacterium bovis* isolates by trial area and

263    host. The table is ordered by total number of isolates from high to low.

| Trial area | Date range | Cattle (n) | Badgers (n) | Total |
|---|---|---|---|---|
| D3 | 2002-2010 | 42 | 191 | 233 |
| I2 | 2002-2007 | 74 | 93 | 167 |
| G2 | 2000-2006 | 118 | 45 | 163 |
| A3 | 2000-2007 | 97 | 62 | 159 |
| B2 | 1999-2007 | 127 | 29 | 156 |
| C3 | 1999-2006 | 94 | 56 | 150 |
| E3 | 2000-2007 | 54 | 75 | 129 |
| J1 | 2002-2008 | 65 | 54 | 119 |
| F1 | 2000-2005 | 31 | 54 | 85 |
| H2 | 2000-2006 | 50 | 31 | 81 |
| **Total** | **1999-2010** | **752** | **690** | **1442** |

264

265    SB0145 (n = 32).  With the exception of SB0140 and SB0263, which were found in

266    multiple trial areas, the geographical distributions of the most prevalent spoligotypes

267    were largely confined to a single trial area (Figure 1C).  Examination of the pairwise

268    SNP distances of isolates within the above spoligotypes showed that there were

269    considerable differences in diversity amongst the spoligotypes (Figure 1D).  High

270    levels of diversity were observed in spoligotypes SB0140 and SB0129 reflecting the

271    phylogenetic structure of the isolates with these spoligotypes. Figure 1E shows

272    pairwise SNP distances for all isolates plotted against geographic distance for each

273    of the most prevalent spoligotypes.

274

275

**Figure 1. Genomic epidemiology of Randomised Badger Culling Trial (RBCT) dataset:**

A) Map showing location of isolation for 1,442 sequenced *Mycobacterium bovis* isolates. Isolates collected from badgers and cattle are shown in red and blue respectively. The proportion of samples from each host is shown in the pie charts and the pie charts are scaled according to number of isolates. The RBCT triplet where each of the isolates were collected is labelled; B) Maximum likelihood phylogenetic tree of 1,442 *M. bovis* isolates rooted with isolates from the Unknown7 clonal complex. Trial area, host, clonal complex and spoligotype are shown as datastrips around the outside of the phylogenetic tree; C) Geographical distributions of the six most prevalent spoligotypes in the dataset. The host of each isolate is represented by a different shape: circle for badger and cross for cattle; D) Frequency distributions of pairwise SNP distances between all isolates belonging to the six most

287   prevalent spoligotypes; E) Scatterplots of pairwise SNP distance against geographic distance

288   in kilometres for all pairs of isolates belonging to the six most prevalent spoligotypes.

289

290   *Transmission*

291   *Transmission clusters*

292   A total of twelve putative transmission clusters, containing 1224/1442 (84.9%) of the

293   isolates, were defined using a conservative threshold of 15 SNPs.  The clusters varied

294   in size between 54 (Cluster 2) and 193 (Cluster 9) isolates (Table 2).  The ratio of cattle

295   to badger isolates in each transmission cluster varied from 0.15 (Cluster 9) to 5.44

296   (Cluster 12; Table 2).  The phylogenetic tree of all 1,442 isolates with the transmission

297   clusters overlaid on it is shown in Figure 2A and the geographical distribution of each

298   transmission cluster is shown in Figure 2B.  The geographical distribution of the

299   transmission clusters was strongly associated with trial area, with the majority of

300   isolates from a transmission cluster found in the same trial area (Figure 2B).

301

302   A multimodal distribution was observed for pairwise SNP distances of isolates

303   assigned to each of the transmission clusters (Supplementary Figure 3).  The first

304   mode comprised pairwise differences of 400 - 500 SNPs and was made up of

305   comparisons of isolates from Eu1 clades deeper in the phylogeny, and the second

306   and third modes between 100-200 SNPs were comprised of isolates from more

307   closely related clades.  The final modes between 0 and 50 SNPs were made up of

308   comparisons of isolates from the same clade and here the within and between

309   transmission cluster comparisons overlapped, although there was a clear peak below

310   15 SNPs representing the transmission clusters themselves.   There were no

311   observable differences between the distributions when the host of each isolate in a

312     pairwise comparison was considered i.e. there were no host-specific patterns of

313     genetic relatedness.

314

315     Table 2: Breakdown of the 12 putative transmission clusters by host. The table is ordered by

316     total number of isolates from high to low.

| Cluster | Badgers (n) | Cattle (n) | Total |
|---|---|---|---|
| **Cluster 9** | 168 | 25 | 193 |
| **Cluster 1** | 46 | 115 | 161 |
| **Cluster 6** | 53 | 86 | 139 |
| **Cluster 8** | 61 | 49 | 110 |
| **Cluster 5** | 50 | 47 | 97 |
| **Cluster 7** | 16 | 76 | 92 |
| **Cluster 10** | 41 | 46 | 87 |
| **Cluster 4** | 19 | 67 | 86 |
| **Cluster 3** | 34 | 49 | 83 |
| **Cluster 11** | 30 | 34 | 64 |
| **Cluster 12** | 9 | 49 | 58 |
| **Cluster 2** | 38 | 16 | 54 |
| **Total** | **565** | **659** | **1224** |

317

318

319     *Temporal analyses of transmission clusters*

320     To describe the temporal dynamics of the transmission clusters, each cluster was

321     independently tested for evidence of temporal signal.  Comparison of root to tip

322     distances with sampling dates did not find significant correlations for any of the

323     transmission clusters (Supplementary Figure 2).  However, dated tip randomisation

324     (DTR) analyses, where evidence of a temporal signal is shown by a lack of overlap

325     between the estimated substitution rates of the observed data and the randomised

326 datasets, showed that there was no overlap between the highest posterior densities

327 (HPD) of the real and randomised datasets for 5/12 of the transmission clusters

328 (Supplementary Figure 3). In a further 5/12 there were overlaps between the HPDs

329 but not medians of the real dataset and one or more of the randomised datasets. For

330 the final two clusters (Cluster 2 and Cluster 4), the median substitution rates of one

331 and five randomised datasets respectively overlapped that of the real datasets
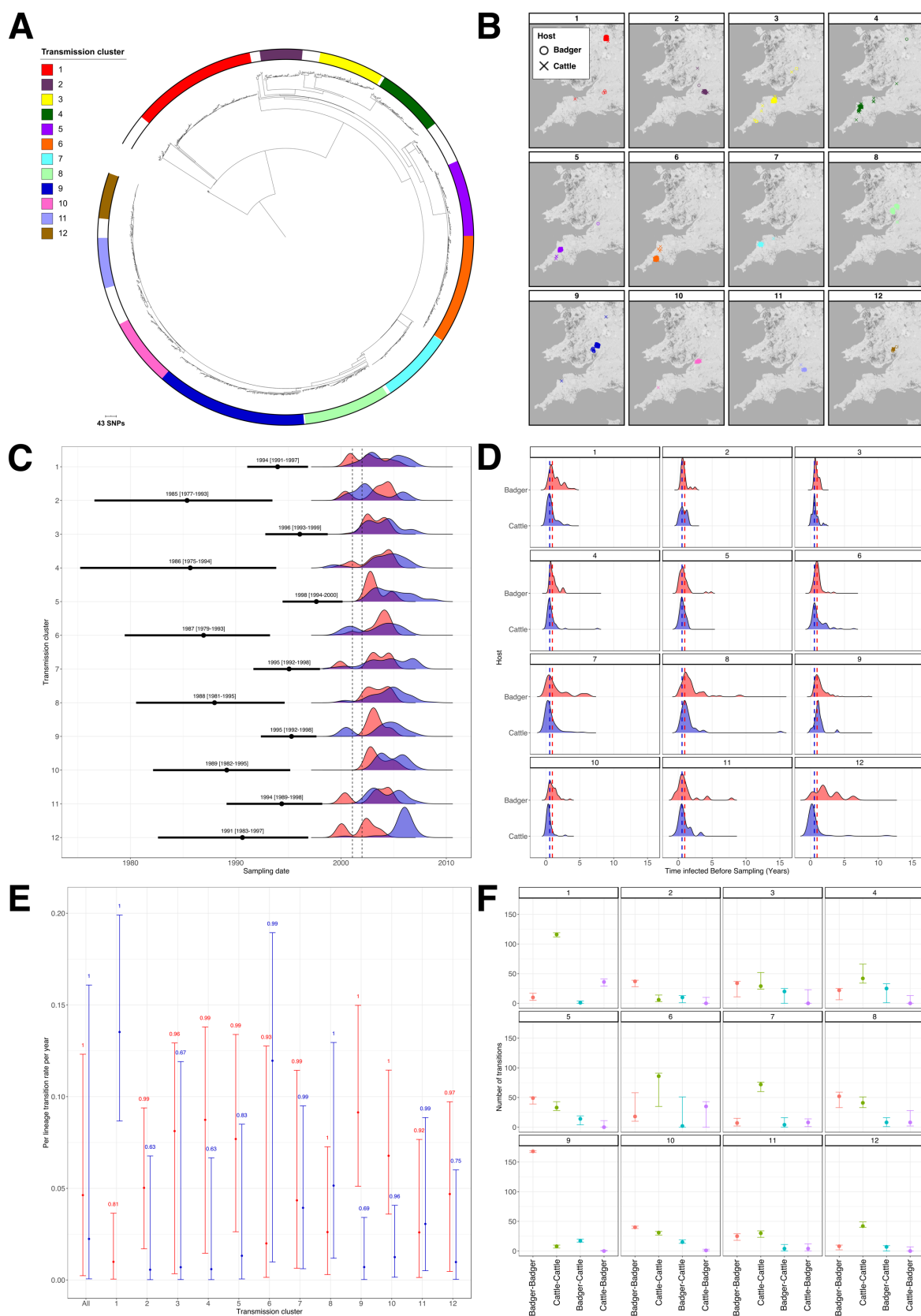
332 (Supplementary Figure 3).

333

334 The median substitution rate of each transmission cluster varied between 0.51

335 (Cluster 12) and 6.0 (Cluster 3) substitutions per genome per year (Supplementary

336 Table 3). Phylogenetic dating analysis using BEAST showed that the estimated date

337 of the MRCA of each transmission cluster varied between 1985 (Cluster 2; 95%

338 confidence interval [CI]: 1977 to 1993) and 1997 (Cluster 5; 95% CI: 1994 to 2000;

339 Figure 2C). The median difference between the MRCA and the date of collection of

340 the first sample was 8.1 (range 4.0 – 15.0) years.

341

342 *Transmission analysis with TransPhylo*

343 TransPhylo was used to estimate the number of unsampled cases. The median

344 number of sampled cases per transmission cluster was 89.5 (range: 54 - 193)

345 compared to a median of 130.1 (range: 41 - 203) inferred unsampled cases

346 suggesting a median case finding rate of 43.2% (range: 28.3% - 74.2%;

347 Supplementary Figure 4). The median inferred time from infection to sampling across

348 all transmission clusters was 0.56 years (95% CI: 0.07 – 2.26 years) for cattle and

349 0.96 years (95% CI: 0.25 – 3.36 years) for badgers (Figure 2D).

350

**Figure 2: Transmission in the Randomised Badger Culling Trial (RBCT) dataset:** A)

Maximum likelihood phylogenetic tree of 1,442 isolates with the twelve putative transmission

353  clusters annotated; B) Geographical distributions of the twelve putative transmission clusters.

354  The host of each isolate is represented by a different shape: circle for badger and cross for

355  cattle; C) Molecular dating of transmission clusters. The inferred median and 95% confidence

356  intervals of the MRCA for each transmission cluster is shown in black. The dates of collection

357  of samples within each transmission cluster are shown as frequency distributions and

358  coloured according to host (red for badgers and blue for cattle). The time period of the

359  suspension of badger culling due to FMD is represented by dashed lines; D) Median length

360  of time of infection for all isolates before sampling per transmission cluster.  The medians of

361  all isolates are shown by red and blue dashed lines for cattle and badgers respectively; E)

362  Estimated inter-species transmission rates for each transmission cluster.  The vertical lines

363  show the lower and upper (2.5% and 97.5%) bounds of the transmission rate distribution for

364  each transmission cluster. The values above the vertical lines represent the posterior

365  probability of each rate and the distributions are coloured according to direction of

366  transmission (red for badger-to-cattle and blue for cattle-to-badger transmission); F) Number

367  of transmissions between known and estimated species counted on each phylogenetic tree

368  in the posterior distribution for each transmission cluster.  The vertical lines show the lower

369  and upper (2.5% and 97.5%) bounds of the distributions. The distributions are coloured

370  according to the type of transmission (red for badger-badger, green for cattle-cattle, blue for

371  badger-cattle and purple for cattle-badger).

372

373  A total of 84 highly supported transmission pairs (posterior probability of transmission

374  between isolate 1 and isolate 2 > 0.5) were identified within the twelve transmission

375  clusters (Table 3) using TransPhylo.  The majority of these transmissions (60/84) were

376  within-species whilst 24/84 were between species. No highly supported transmission

377  pairs were identified in Cluster 3.  The median pairwise SNP distances for the highly

378  supported transmission pairs across all transmission clusters were 1 (range: 0-8), 1

379  (range: 0-5) and 1 (range: 0-6) for cattle to cattle, badger to badger and between-

380  species transmission respectively.

381

382  *Directionality of transmission between host species*

383  BASTA was used to determine the dominant direction of transmission between host

384  species and found higher rates of transmission from badgers to cattle in 8/12

385  transmission clusters and from cattle to badgers in 4/12 clusters (the confidence

386  intervals for each direction overlapped in all clusters except clusters 1 and 9; Figure

387  2E).  For the networks with higher badger to cattle transmission, this direction of

388  infection occurred between 1.1 (Cluster 7) and 14.8 (Cluster 4) times more frequently

389  than in the opposite direction.  By comparison, in the four networks with higher cattle

390  to badger transmission, the frequency was between 1.2 (Cluster 11) and 13.6 (Cluster

391  1) times higher than in the opposite direction.  The overall median badger to cattle

392  transmission rate for all transmission clusters was 2.1 (95% CI: 0.8-3.8) times higher

393  than the cattle to badger transmission rate.

394

395  As BASTA does not directly calculate the transmission rate within-species, the lower

396  bound of the number of transmissions (the count of transitions in the posterior trees)

397  between different animals, regardless of host, was also calculated for each

398  transmission cluster (Figure 2F).  For each of the clusters, the estimated number of

399  transmission events is consistent with the estimated inter-species transmission rates.

400  Across the twelve clusters, the number of within-species transmission events was

401  higher than the between-species transmissions with the average number of cattle to

402    cattle transmission events 4.9 (range 0-31) times greater than the number of cattle to

403    badger transmission events and 17 (range 0.5-116) times greater than the number of

404

405    Table 3: Highly supported transmission pairs within each transmission cluster.  For each

406    transmission cluster, the number of intra- and inter-species transmission pairs with a

407    probability > 0.5 is listed.  The median SNP distance for each set of transmission pairs is

408    given in parentheses.

| Transmission cluster | Number highly supported transmission pairs | Number Cattle-Cattle transmission pairs (median SNP distance) | Number Badger-Badger transmission pairs (median SNP distance) | Number Between-species transmission pairs (median SNP distance) |
|---|---|---|---|---|
| Cluster 1 | 9 | 4 (2) | 1 (0) | 4 (1) |
| Cluster 2 | 1 | 0 | 0 | 1 (1) |
| Cluster 3 | 0 | 0 | 0 | 0 |
| Cluster 4 | 1 | 1 (1) | 0 | 0 |
| Cluster 5 | 14 | 3 (1) | 7 (1) | 4 (1) |
| Cluster 6 | 5 | 4 (0.5) | 1 (2) | 0 |
| Cluster 7 | 15 | 11 (1) | 0 | 4 (1) |
| Cluster 8 | 10 | 4 (1) | 3 (1) | 3 (1) |
| Cluster 9 | 17 | 0 | 12 (0) | 5 (3) |
| Cluster 10 | 2 | 1 (1) | 1 (0) | 0 |
| Cluster 11 | 5 | 3 (2) | 1 (2) | 1 (5) |
| Cluster 12 | 5 | 0 | 3 (7) | 2 (4) |
| Total | 84 | 31 | 29 | 24 |

409

410    badger to cattle transmission events. The number of badger to badger transmission

411    events was 4.7 (range 0.9-10) times higher than the number of badger to cattle

412    transmission events and 4.5 (range 0-40) times higher than the number of cattle to
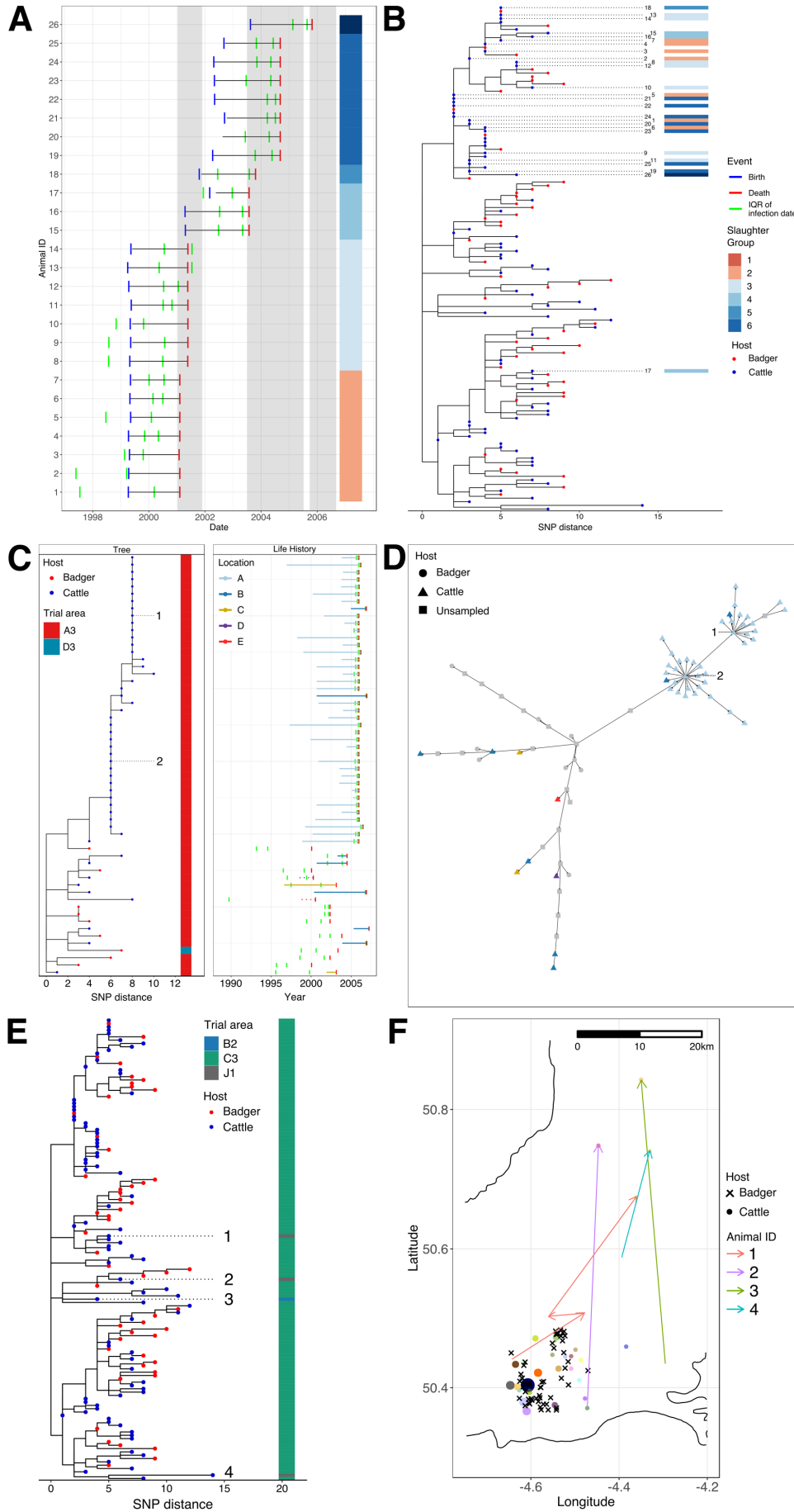
413    badger transmission events.

414

415    *Other Transmission dynamics*

416    Metadata from the SAM database was incorporated for each of the cattle within a

417    transmission cluster that could be matched, which allowed examples of recurrence

418    (detection of infections subsequent to previous outbreaks or breakdowns),

419    superspreading (individual hosts that have a disproportionate effect on the spread of

420    infection) and long-distance transmission (transmission between different trial areas)

421    to be characterised.

422

423    *Recurrence*

424    A total of 47 isolates formed a distinct clade within the Cluster 6 phylogeny (Figure

425    3B). Of these, 25 were isolates from cattle slaughtered between February 2001 and

426    October 2005 as part of three recorded breakdowns on the same farm comprising 35

427    confirmed cases (Figure 3A). Based on the date of slaughter, these isolates were

428    divided into six slaughter groups (Figure 3A). Two pairs of isolates from different

429    animals (1 and 20; 6 and 23) and different slaughter groups (1 and 5) were 0 SNPs

430    apart despite the subsequently infected animals (20 and 23) only moving to the farm

431    a year after the first animals were slaughtered (Figure 3A). The subsequently infected

432    animals were then slaughtered three years later as part of a later breakdown. The

433    majority of the rest of the cattle isolates in this clade were also very similar despite

434

435 **Figure 3: Integration of genomics and cattle movement data:** A) Life histories of animals

436 included in Cluster 6. Date of birth is shown in blue, date of death in red and the interquartile

437 range (IQR) for the estimated date of infection is shown in green. The grey shading represents

438 breakdowns and Slaughter Groups are labelled based on date of death; B) SNP-scaled

439 phylogenetic subtree for Cluster 6.  Tips are coloured according to host (red for badger and

440 blue for cattle), relevant isolates are labelled and Slaughter Group is shown as a data strip;

441 C) SNP-scaled phylogenetic tree for Cluster 12.  Tips are coloured according to host (red for

442 badger and blue for cattle), isolates inferred as superspreaders are labelled and Trial area is

443 shown as a data strip. Life histories are shown for each isolate. Date of birth is shown in blue,

444 date of death in red and the interquartile range for the estimated date of infection is shown

445 in green. The farm location identifier of each cattle isolate is coloured according to the legend

446 and the length of each bar reflects the length of time an animal spent at its final location; D)

447 Transmission network for Cluster 12.  The shape of each node is based on the host (circle

448 for badger or triangle for cattle) or else a square for inferred unsampled cases.  Nodes are

449 coloured based on their location (grey for unsampled and badger isolates which don't have

450 a farm location identifier).  Isolates inferred as superspreaders are labelled; E) SNP-scaled

451 phylogenetic tree for Cluster 6.  Tips are coloured according to host (red for badger and blue

452 for cattle), isolates highlighted as examples of long-distance transmission are labelled and

453 the Trial area for each isolate is shown as a data strip; F) Map showing the location of isolation

454 for each isolate in Cluster 6. The shape of each isolate is based on the host (cross for badger

455 or circle for cattle), the colour of the cattle isolates is based on the farm location and the size

456 of the cattle isolates reflects the number of animals in Cluster 6 at that location.  The

457 movements of the animals from which isolates highlighted as examples of long-distance

458 transmission were collected are shown as arrows and coloured according to animal identifier.

459 Herd and badger locations were randomly shifted by up to 1 km in the horizontal and vertical

460 planes.

461

462    the long time periods from when these animals arrived on the farm and when they

463    were slaughtered.   Cattle were still present on the farm between the different

464    breakdowns suggesting that infection was being maintained locally, either in this or

465    a neighbouring herd or else within the local badger population.

466

467    *Superspreading*

468    The structure of the Cluster 12 phylogeny showed a number of cattle isolates

469    clustering together within a very flat tree structure i.e. the majority of these isolates

470    were 0 SNPs apart (Figure 3C).  Of these, 38 isolates were from animals slaughtered

471    at a single location (A) as part of a breakdown between April 2005 and July 2007 that

472    identified 59 reactors from which 39 had *M. bovis* cultured.   The resulting

473    transmission network inferred two distinct superspreading events of 13 and 21 cases

474    inferred to be centred around two animals (1 and 2; Figure 3D).  Two of the animals

475    in these superspreading events were from a different location (B) though this was only

476    0.83 km away, suggesting potential epidemiological links between locations A and B.

477

478    *Long-distance transmission*

479    Isolates from four cattle in Cluster 6 were not from the predominant trial area (C3) of

480    this cluster (Figure 3E).  For three of these isolates (1, 2 and 3), the inferred dates of

481    infection showed that the animals, which had moved between farms, were likely

482    infected at a location in or near trial area C3.  For the remaining animal, 4, the inferred

483    date of infection didn't overlap with a previous location in or near trial area C3 but the

484    location data for 58 days of this animal's life is missing from the database.   The

485    distances moved by the infected animals ranged between 23 and 46 km providing

486    evidence of movement mediated transmission (Figure 3F).

487

488    **Discussion**

489    The RBCT was set up to assess the impact of badger culling on the incidence of bTB

490    in nearby cattle herds.  In this study, the resulting badger and cattle isolates have

491    been sequenced, meaning that, for the first time, WGS of large numbers of *M. bovis*

492    isolates from co-located populations of both species collected contemporaneously

493    in well-defined geographical areas can be used to address key questions surrounding

494    bTB transmission in the high-risk regions of England.

495

496    A total of 60 unique spoligotype patterns were identified in the dataset though the

497    majority of the isolates (1320/1442; 91.5%) were made up of the six most prevalent

498    spoligotypes confirming that there is relatively little genetic heterogeneity at this level

499    amongst *M. bovis* in the high prevalence areas of England.  The observed prevalence

500    of spoligotypes in this study closely matched previously published data, generated

501    using traditional typing methods, from the same time period [66], showing that the

502    dataset accurately reflects the known population structure of *M. bovis* at that time.

503

504    Genotyping methods such as spoligotyping are used to infer close relationships

505    between isolates (low genetic diversity) and assume monophyly.  However, it is clear

506    from plotting spoligotypes on the phylogenetic tree in Figure 1B, that some of the

507    spoligotypes, in particular SB0140, are polyphyletic.  Whilst the majority of SB0140

508    isolates sit adjacent to each other in the phylogenetic tree, there is a single clade that

509     sits separately with the SB0274 and SB0673 clades falling between them. Another

510     example is the intermingling of SB0957 and SB0263 isolates in trial area I2. An

511     alternative way to view this data is to calculate and plot pairwise SNP distances for

512     each of the most prevalent spoligotypes (Figure 1D). From this, we see that while

513     there is a maximum of approximately 50 SNPs between members of four of the six

514     most prevalent spoligotypes, for SB0140 and SB0129 the maximum pairwise SNP

515     distance was between 150 and 200 SNPs, demonstrating higher levels of diversity

516     for the *M. bovis* population as evidenced from genome-wide data compared to

517     traditional typing data.

518

519     As we observed in our data, it has been recognised for a long time that spoligotypes

520     can be homoplasic and identical spoligotype patterns can be found in

521     phylogenetically unrelated strains [66]. Despite this, spoligotyping, along with

522     Mycobacterial interspersed repetitive unit-variable number tandem repeat (MIRU-

523     VNTR) analysis, has continued to be the most commonly applied method for

524     genotyping *M. bovis* as it is cheap and comparatively straightforward to implement

525     in the laboratory. However, given the much higher resolution offered by WGS and

526     the fact that national bodies such as APHA and the United States Department of

527     Agriculture (USDA) are now moving towards routinely sequencing all cases of

528     *M.bovis*, it may now be time to move towards a SNP-based method of typing *M.*

529     *bovis* isolates similar to the Coll method adopted for typing *M. tuberculosis sensu*

530     *stricto* [67].

531

532   The predominance of clonal complex Eu1 in our dataset was unsurprising as previous

533   work, including the study that first defined Eu1 [24], has shown that this clonal

534   complex is ubiquitous in Great Britain, Northern Ireland and the Republic of Ireland

535   whilst being uncommon in mainland Europe, where *M. bovis* genetic diversity is much

536   greater [42].  Previous work using both PCR and genomics to assign clonal complex

537   shows that Eu1 is likely the most predominant globally-circulating *M. bovis* lineage

538   and it has been found in many countries that have historically traded cattle with the

539   UK such as New Zealand, the USA, Mexico and Uruguay [33-35, 41].  The presence

540   of this clonal complex and its spoligotypes such as SB0140 and SB0263 in these

541   countries suggests that Eu1 has been present in the UK for at least 200 years or more.

542   This is supported by the high level of pairwise SNP diversity observed within SB0140.

543

544   The Unknown7 isolates in the dataset were found in trial areas A3, D3 and G2 and

545   were a maximum of 12 SNPs apart.  The small number of isolates suggest that it is

546   uncommon in the UK and the low level of genetic diversity suggests that it may be

547   part of a single, recent introduction.  This clonal complex has also been found in

548   France, Mali and the USA [27, 33, 42] and given the geographical proximity of France

549   and the ongoing trade of cattle between the two countries this may be the most likely

550   origin for this lineage.

551

552   The prevailing hypothesis for the population history of *M. bovis*, in particular Eu1, in

553   the UK is that there was a single introduction followed by long term endemicity and

554   a population bottleneck due to effective control measures beginning in the 1930s [66].

555   However, the observed phylogenetic structure of isolates when incorporated with a

556   global collection of Eu1 isolates, indicates that there have likely been multiple,

557   perhaps as many as four, introductions of Eu1 into England (Supplementary Figure

558   5).  Whilst we did not attempt to date these introductions in this study, the availability

559   of archived isolates from the 1980s along with contemporaneous isolates will be used

560   alongside the RBCT dataset and other published UK datasets in future work to

561   provide estimated dates for these introductions.

562

563   Due to the large size and clear phylogenetic and geographical structure of the dataset

564   as well as our study aims, we defined transmission clusters using a conservative

565   pairwise SNP threshold of 15 SNPs.  There is currently no consensus as to what is

566   the best threshold to apply to *Mycobacterium* genome datasets with previous studies

567   using thresholds between three and fifteen SNPs [40, 50].  We chose the 15 SNP

568   threshold as it would allow for the possible identification of older transmission events

569   but also allow for any variance in the rates of mutation amongst the sampled isolates.

570   We chose to use the software package TransPhylo as it integrates dates of isolate

571   collection and genetic relatedness and allows for within-host diversity and unsampled

572   cases.  The first step of the analysis was to generate molecular dated phylogenies for

573   each of the transmission clusters.  Assessing the presence of temporal signal in a

574   genome dataset is typically done in two ways: examining the linear relationship

575   between root to tip distance and sampling date (under a perfect clocklike behaviour,

576   then $R^2 = 1$ [68]), and dated tip randomization (DTR) analysis.  In DTR, the dates of

577   sampling are repeatedly shuffled amongst the taxa and the clock rates between the

578   observed and random data calculated and compared [69].  If there is no overlap

579   between the estimated substitution rates of the observed data and the randomized

580    datasets then we can conclude that the observed dataset has a stronger temporal

581    signal than expected by chance [69]. We obtained very low or negative values for $R^2$

582    for all our transmission clusters which is normally interpreted as evidence for a lack

583    of temporal signal or else overdispersion in lineage-specific clock rates [70]

584    (Supplementary Figure 2). The previously reported slow substitution rate of *M. bovis*

585    [10, 34] and the short window of sampling (twelve years) may explain the lack of

586    association between root to tip distances and sampling dates. As root to tip

587    regression is only a tool for exploratory analysis [70] we performed DTR on all of the

588    transmission clusters (Supplementary Figure 3). From this, we observed that there

589    was strong evidence of temporal signal in 5/12 of the transmission clusters, moderate

590    temporal signal in 5/12 transmission clusters and weak temporal signal in two

591    clusters, particularly Cluster 2 (Supplementary Figure 2). Despite the especially weak

592    evidence of temporal signal in Cluster 2, we decided to include it in our analyses as,

593    given the close relatedness of all of the clusters, it is highly likely that the mutational

594    process, and therefore the molecular clock, will be similar in each of them. As well

595    as being used as input for TransPhylo molecular dating of each of the transmission

596    clusters provided additional insights into the dataset.

597

598    Firstly, we were able to calculate substitution rates for each transmission cluster

599    which ranged between 0.5 and 6 SNPs per genome per year (Supplementary Table

600    3). Published estimates for the median substitution rate of *M. bovis* vary between

601    0.15 and 0.53 substitutions per genome per year [10, 34]. Previous work has shown

602    that there are lineage and study specific differences in the substitution rates within

603    the *Mycobacterium tuberculosis* complex (MTBC) [71]. This analysis showed that

604 higher substitution rates were found in smaller datasets with narrow sample date

605 ranges, which may explain the much higher substitution rates of up to six

606 substitutions per genome per year we observed in our analyses. Secondly, we were

607 able to infer the date of the MRCA of each transmission cluster (Figure 2C). The short

608 time period (4 – 15 years) between the inferred date of the MRCAs and the earliest

609 sample collection dates suggests that the transmission clusters are likely the result

610 of recent seeding events and not a consequence of endemic disease in the form of

611 long-term maintenance within herds or an endemic wildlife reservoir, in this case

612 badgers. The introduction of a compulsory test and slaughter scheme in the UK in

613 the 1950s saw a sustained decline in the annual number of infected animals removed

614 as TB test reactors and infected cattle herds with only a few hundred reactors being

615 detected annually in the early 1980s [66]. However, this decline was reversed from

616 the late 1980s, with the UK now having one of the highest incidences of bTB in

617 Europe. From our analyses, it is clear that the dates for the MRCA of each of the

618 transmission clusters overlap with this population expansion.

619

620 From the TransPhylo analysis we were able to estimate what proportion of infected

621 hosts we managed to sample for each transmission cluster (Supplementary Figure

622 4). Sampling of all hosts infected with a disease is never complete due to a range of

623 factors such as detection, failure to culture and in the case of genomics, issues

624 related to sequence quality. In this study we estimate that we managed to sequence

625 a median of 43.2% of cases across the transmission clusters though this varied from

626 less than 30% to as high as 75% depending on the cluster. Obviously, the success

627 of sampling has an impact on the types and quality of the inferences we are able to

628 make. For instance, we were able to confidently identify and confirm a

629 superspreading event in Cluster 12 due to having sequenced 38/39 of the confirmed

630 cases in a breakdown (see below).

631

632 The incubation period of TB in cattle is generally believed to be several months and

633 potentially years, although there is some evidence of much shorter incubation periods

634 in other mammalian species such as cats [72]. There is also typically a lag period

635 (occult period) between infection and detection where infections are undetectable to

636 the standard tuberculin test [73]. The organism may also persist for several years

637 within infected animals before they are detected (latency) and reactivation has not

638 been demonstrated in cattle. To date, there are no firm estimates for either the

639 duration of the occult period or of epidemiological latency, which is problematic for

640 fitting transmission models [74] and predicting the impact of control polices [75].

641 Based on our analysis using TransPhylo, we can provide estimates for how long both

642 badgers and cattle were infected before sampling. The analysis showed that, on

643 average, badgers were infected for twice as long as cattle before sampling. The

644 median period of infection for cattle of 0.56 years is consistent with the annual testing

645 schedule imposed on cattle during the RBCT. Whilst there was a wide range of

646 estimates for the length of infection the 95% confidence intervals for both badgers

647 and cattle were within the normal lifespans of both species.

648

649 We were able to identify a small number of highly supported direct transmission

650 events, defined as transmission pairs that had a posterior probability greater than 0.5

651 (Table 3). Although the majority (60/84) of these transmission events occurred

652   between the same species, there were also 24 interspecies transmission pairs across

653   the transmission clusters with pairwise SNP distances varying between 1 and 5

654   SNPs. To date, there is limited evidence of badgers and cattle directly interacting

655   and the majority of transmission is considered to be indirect i.e. through the

656   environment [76]. Given the inferred number of unsampled cases and small number

657   of highly supported transmission pairs, more intensive sampling would need to be

658   performed to better establish transmission dynamics between the different bTB host

659   species. Despite the logistical challenges around detecting and culturing *M. bovis* in

660   environmental samples, the inclusion of samples from faeces and feed troughs and

661   other potential hosts such as rodents and cervids should be an integral part of any

662   future work.

663

664   One of the aims of this study was to assess and quantify the directionality of

665   transmission between cattle and badgers. For this we used a Bayesian evolutionary

666   tool, BASTA (Bayesian Structured coalescent Approximation), to estimate the

667   interspecies transmission rates in each of the transmission clusters. BASTA was

668   designed to estimate evolutionary dynamics in structured populations and account

669   for sampling biases. For the majority of our transmission clusters, badger to cattle

670   transmission occurred more frequently even in clusters with approximately equal

671   numbers of cattle and badgers (Figure 2E). It is worth noting that the estimated

672   transmission rates were very low with the median number of badger to cattle

673   transmissions across all transmission clusters estimated as 0.05 transmissions per

674   lineage per year and the median number of cattle to badger transmissions estimated

675   as 0.02 transmissions per lineage per year (Figure 2E). Whilst BASTA does not

676    directly estimate intra-species transmission rates we could calculate the number of

677    transmission events between each host species from the posterior log and tree files.

678    These are conservative counts of the minimum number of transitions between

679    sampled animals and their ancestors but do allow us to compare the number of inter-

680    and intra-species transmissions. From this we were able to demonstrate that inter-

681    species transmission occurs much less frequently than intra-species transmission in

682    our transmission clusters and cattle to cattle transmission is more common than

683    badger to badger transmission (Figure 2F).   Three previous studies, each on small

684    geographically localized populations, have used BASTA to estimate rates of

685    transmission between badgers and cattle; the first estimated that badger to cattle

686    transmission was 10.4 times more frequent than cattle to badger transmission, the

687    second estimated that cattle to badger transmission was at least an order of

688    magnitude higher than badger to cattle transmission and the third estimated that

689    cattle to badger transmission was at least three times higher than badger to cattle

690    transmission (a similar result was obtained using a similar transmission analysis

691    package, MASCOT) [9, 12, 13]. These results, along with those described in this

692    study, suggest that the directionality of transmission may vary between sampling area

693    although badger to cattle transmission does appear to be more frequent. What is

694    consistent across all the studies, however, is that intra-species transmission occurs

695    much more frequently than inter-species transmission.

696

697    Beyond the original aims of the project such as characterising the population

698    structure of *M. bovis* isolates collected as part of the RBCT and investigating

699    interspecies transmission, the utility of WGS was also shown through its application

700  to other important aspects of bTB transmission in the UK. The combination of

701  genomics and the extensive cattle tracing database allowed us to characterise

702  examples of recurrence, superspreading and long-distance transmission within the

703  dataset. Previous work has shown that prior history of bTB within a herd is an

704  important predictor of breakdown: 38% of herds that clear movement restrictions

705  experience another breakdown within 24 months [77]. This suggests that infection is

706  being maintained within herds despite repeated testing and it is estimated that

707  between 24% and 50% of recurrent breakdowns are due to persistence within the

708  herd [74]. We were able to use pairwise genome comparisons to identify near

709  identical isolates that were collected up to four years apart and which were part of

710  confirmed herd breakdowns. Examination of the cattle movements confirmed that

711  some of these isolates were collected from animals that arrived subsequent to the

712  dates of slaughter of infected animals as part of previous breakdowns. The similarity

713  of these more recent isolates to the earlier isolates would suggest that the animals

714  were infected after their arrival in the new location and that control measures following

715  the prior breakdowns were insufficiently effective.

716

717  TransPhylo allowed us to generate plausible transmission networks where star like

718  nodes representative of potential superspreaders (individual hosts that have a

719  disproportionate effect on the spread of infection) could be identified (Figure 3C/3D).

720  We were then able to incorporate data from the CTS to identify the cattle likely acting

721  as the source of the infections. Whilst previous work using modelling or network

722  analysis has highlighted the importance of small numbers of farms or herds as hubs

723  of transmission which act as superspreaders of infection [78, 79], we provide the first

724   evidence, based on genomics and cattle movement data, that particular animals

725   within herds may also act as superspreaders potentially contributing to increased

726   transmission between different locations if these animals are not identified before

727   being moved.  We were unable to identify any superspreaders amongst any of the

728   sampled badgers.

729

730   From the temporal analysis of the transmission clusters we showed that these

731   clusters are comparatively young and likely recently seeded.  The most likely

732   mechanism for this is the movement of infected cattle into a location followed by

733   subsequent onward transmission within the herd and into the local badger

734   populations.  Given the median estimate of the MRCA of the transmission clusters

735   was eight years before sampling began, this precluded any possibility of us sampling

736   the index case for any of the transmission clusters.  However, by incorporating cattle

737   movement information with our transmission clusters, we were able to identify cattle

738   infected with a particular lineage in one trial area moving to a trial area further away,

739   highlighting the potential for long distance transmission events to seed new

740   transmission clusters (Figure 3E/3F).  This was also recently demonstrated by Rossi

741   et al. who identified an imported infected animal or animals as being responsible for

742   a bTB outbreak in a region of England with no previously known wildlife infections

743   [12].  This has important implications for infection control; even with the limited

744   sampling we conducted, the combination of genomics and cattle movements still

745   allowed us to identify these potential seeding events.  More targeted testing and

746   sequencing before animals are moved, particularly to lower incidence areas, would

747  potentially identify these likely sources of infection before they are able to become

748  established in other locations.

749

750  Potential limitations of our analysis were the choice and number of samples included

751  in the study and known issues surrounding the lack of a strong temporal signal in *M.*

752  *bovis* that may affect the results of any analyses based on molecular dating.  Any

753  sampling strategy we selected would not have been perfect; ideally, we would have

754  tried to sequence all samples collected as part of the RBCT; however, this was not

755  possible due to cost and manpower constraints so we chose to sequence only the

756  badger and cattle isolates collected from proactive triplets excluding isolates from

757  infected badgers culled in reactive triplets and infected cattle culled as part of

758  contemporaneous breakdowns.  From our TransPhylo analysis we estimated that we

759  managed to sample approximately 40% of infected cases across our transmission

760  clusters.  Despite this, the size of the dataset was still large enough to generate

761  several large transmission clusters that allowed us to draw robust conclusions about

762  transmission, notably directionality of transmission between badgers and cattle.

763  Comparison of the spoligotype distribution in our study to earlier work confirmed that

764  our dataset was representative of the known population structure during the RBCT.

765

766  We know from previous work that the lack of a strong temporal signal is a potential

767  issue when attempting to accurately date the origin of particular lineages [71].  The

768  results of the dated tip randomization analysis indicated that there was moderate or

769  strong temporal signal in nearly all of our transmission clusters; however, two of our

770  transmission clusters notably Cluster 2 had a weak temporal signal.  The range of

771   substitution rates we estimated for some of our transmission clusters was also higher

772   than previously observed which may have affected the estimated dates of those

773   transmission cluster's MRCAs.  Overall, however, even if individual networks such as

774   Cluster 2 with little or no temporal signal or Cluster 3 with a high substitution rate are

775   of concern, the conclusions we have drawn are based on considering the results from

776   twelve different transmission clusters composed of over 1,200 genomes and thus can

777   be considered robust.

778

779   Multiple previous studies have shown that bTB transmission is complicated, unlikely

780   to be driven by a single mechanism and is strongly associated with the setting and

781   host dynamics of the system being studied.  Here we used the largest single country

782   genome dataset alongside the national cattle movement database to attempt to

783   address key questions around bTB transmission in a multi-host, intensive setting.

784   Whilst both the TransPhylo and BASTA results support inter-species transmission

785   with some evidence that there is broadly more badger to cattle transmission than in

786   the opposite direction, it is clear that the majority of ongoing transmission is occurring

787   within cattle herds and within the badger populations.  Spillover in either direction

788   could then be considered to be occurring at a low level and, based on the dates of

789   their MRCAs, the transmission clusters we defined are likely to have been the result

790   of recent seeding events and are primarily being maintained by within-species

791   transmission.  We have also provided the first genomics-based estimates for the

792   length of time that badgers and cattle are infected with bTB before sampling.  Finally,

793   we were able to characterise recurrence, superspreading and long-distance

794   transmission within our transmission clusters.

**Data availability**

Raw sequencing reads were deposited at the European Nucleotide Archive (https://www.ebi.ac.uk/ena/browser/home) under project PRJEB19799; all accessions used in this project are listed in Supplementary File 1. Metadata for the sequenced isolates is available on pubMLST (https://pubmlst.org/projects/mbovis-eradbtb).

**Code availability**

The R code used to perform data analyses in this study is available in GitHub (https://github.com/avantonder/RBCT).

**Author contributions**

A.J.K.C., R.G.H., J.L.N.W. and J.P. conceived the study. J.D. cultured, heat inactivated and submitted the *M. bovis* isolates for sequencing. L.G. and A.P.M.

819    extracted metadata for the study isolates from the CTS and Sam databases and

820    uploaded this metadata to a BIGSdb database created by K.A.J. A.J.v.T. and M.T.

821    performed the data analysis. A.J.v.T. coordinated the study and wrote the initial draft

822    of the manuscript.  A.J.v.T, A.J.K.C., E.P. P.J.H., J.L.N.W. and J.P. contributed to the

823    final version of the manuscript. All authors read and approved the manuscript.

824

825    **Competing interests**

826    The authors declare no competing interests.

827

828    **References**

829    1.    Olea-Popelka F, Muwonge A, Perera A, Dean AS, Mumford E, Erlacher-Vindel

830          E, Forcella S, Silk BJ, Ditiu L, El Idrissi A, et al: **Zoonotic tuberculosis in**

831          **human beings caused by *Mycobacterium bovis*-a call for action.** *Lancet*

832          *Infectious Diseases* 2017, **17:**E21-E25.

833    2.    Godfray HCJ, Donnelly C, Hewinson G, Winter M, Wood J: **Bovine TB**

834          **strategy review.** 2018.

835    3.    Muirhead RH, Gallagher J: **Tuberculosis in Wild Badgers in Gloucestershire**

836          **- Epidemiology.** *Veterinary Record* 1974, **95:**552-555.

837    4.    Hauer A, De Cruz K, Cochard T, Godreuil S, Karoui C, Henault S, Bulach T,

838          Banuls AL, Biet F, Boschiroli ML: **Genetic Evolution of *Mycobacterium bovis***

839          **Causing Tuberculosis in Livestock and Wildlife in France since 1978.** *Plos*

840          *One* 2015, **10**.

841    5.    Morris RS, Pfeiffer DU: **Directions and issues in bovine tuberculosis**

842          **epidemiology and control in New Zealand.** *New Zealand Veterinary Journal*

843          1995, **43:**256-265.

844    6.    Godfray HCJ, Donnelly CA, Kao RR, Macdonald W, McDonald RA,

845          Petrokofsky G, Wood JLN, Woodroffe R, Young DB, McLean AR: **A**

846          **restatement of the natural science evidence base relevant to the control**

847        **of bovine tuberculosis in Great Britain.** *Proceedings of the Royal Society B-*
848        *Biological Sciences* 2013, **280**.

849   7.    Bourne J: *Bovine TB: the scientific evidence: a science base for a sustainable*
850        *policy to control TB in cattle: an epidemiological investigation into bovine*
851        *tuberculosis.* Department for Environment, Food and Rural Affairs; 2007.

852   8.    Woodroffe R, Donnelly CA, Johnston WT, Bourne FJ, Cheeseman CL, Clifton-
853        Hadley RS, Cox DR, Gettinby G, Hewinson RG, Le Fevre AM, et al: **Spatial**
854        **association of *Mycobacterium bovis* infection in cattle and badgers Meles**
855        **meles.** *Journal of Applied Ecology* 2005, **42:**852-862.

856   9.    Crispell J, Benton CH, Balaz D, De Maio N, Ahkmetova A, Allen A, Biek R,
857        Presho EL, Dale J, Hewinson G, et al: **Combining genomics and**
858        **epidemiology to analyse bi-directional transmission of *Mycobacterium***
859        ***bovis* in a multi-host system.** *Elife* 2019, **8**.

860  10.   Biek R, O'Hare A, Wright D, Mallon T, McCormick C, Orton RJ, McDowell S,
861        Trewby H, Skuce RA, Kao RR: **Whole Genome Sequencing Reveals Local**
862        **Transmission Patterns of *Mycobacterium bovis* in Sympatric Cattle and**
863        **Badger Populations.** *Plos Pathogens* 2012, **8**.

864  11.   Ahlstrom C, Barkema HW, Stevenson K, Zadoks RN, Biek R, Kao R, Trewby
865        H, Haupstein D, Kelton DF, Fecteau G, et al: **Limitations of variable number**
866        **of tandem repeat typing identified through whole genome sequencing of**
867        ***Mycobacterium avium* subsp. paratuberculosis on a national and herd**
868        **level.** *BMC Genomics* 2015, **16:**161.

869  12.   Rossi G, Crispell J, Brough T, Lycett SJ, White PCL, Allen A, Ellis RJ, Gordon
870        SV, Harwood R, Palkopoulou E, et al: **Phylodynamic analysis of an emergent**
871        ***Mycobacterium bovis* outbreak in an area with no previously known**
872        **wildlife infections.** *bioRxiv* 2020**:**2020.2011.2012.379297.

873  13.   Akhmetova A, Guerrero J, McAdam P, Salvador LCM, Crispell J, Lavery J,
874        Presho E, Kao RR, Biek R, Menzies F, et al: **Genomic epidemiology of**
875        ***Mycobacterium bovis* infection in sympatric badger and cattle**
876        **populations in Northern Ireland.** *bioRxiv* 2021**:**2021.2003.2012.435101.

877 14. Jolley KA, Bray JE, Maiden MC: **Open-access bacterial population genomics: BIGSdb software, the PubMLST. org website and their applications.** *Wellcome open research* 2018, **3**.

880 15. Jolley KA, Maiden MCJ: **BIGSdb: Scalable analysis of bacterial genome variation at the population level.** *Bmc Bioinformatics* 2010, **11**.

882 16. Team RC: **R: A language and environment for statistical computing.** Vienna, Austria; 2013.

884 17. Kahle D, Wickham H: **ggmap: Spatial Visualization with ggplot2.** *R Journal* 2013, **5:**144-161.

886 18. Andrews S: **FastQC: a quality control tool for high throughput sequence data.** Babraham Bioinformatics, Babraham Institute, Cambridge, United Kingdom; 2010.

889 19. Wood DE, Lu J, Langmead B: **Improved metagenomic analysis with Kraken 2.** *Genome Biology* 2019, **20**.

891 20. Lu J, Breitwieser FP, Thielen P, Salzberg SL: **Bracken: estimating species abundance in metagenomics data.** *Peerj Computer Science* 2017.

893 21. Xia E, Teo YY, Ong RTH: **SpoTyping: fast and accurate in silico Mycobacterium spoligotyping from sequence reads.** *Genome Medicine* 2016, **8**.

896 22. Faksri K, Xia E, Tan JH, Teo YY, Ong RTH: **In silico region of difference (RD) analysis of *Mycobacterium tuberculosis* complex from sequence reads using RD-Analyzer.** *Bmc Genomics* 2016, **17**.

899 23. Rodriguez-Campos S, Schurch AC, Dale J, Lohan AJ, Cunha MV, Botelho A, De Cruz K, Boschiroli ML, Boniotti MB, Pacciarini M, et al: **European 2-A clonal complex of *Mycobacterium bovis* dominant in the Iberian Peninsula.** *Infection Genetics and Evolution* 2012, **12:**866-872.

903 24. Smith NH, Berg S, Dale J, Allen A, Rodriguez S, Romero B, Matos F, Ghebremichael S, Karoui C, Donati C, et al: **European 1: A globally important clonal complex of *Mycobacterium bovis*.** *Infection Genetics and Evolution* 2011, **11:**1340-1351.

907 25. Muller B, Hilty M, Berg S, Garcia-Pelayo MC, Dale J, Boschiroli ML, Cadmus S, Ngandolo BNR, Godreuil S, Diguimbaye-Djaibe C, et al: **African 1, an**

**Epidemiologically Important Clonal Complex of *Mycobacterium bovis* Dominant in Mali, Nigeria, Cameroon, and Chad.** *Journal of Bacteriology* 2009, **191:**1951-1960.

26.    Berg S, Garcia-Pelayo MC, Muller B, Hailu E, Asiimwe B, Kremer K, Dale J, Boniotti MB, Rodriguez S, Hilty M, et al: **African 2, a Clonal Complex of *Mycobacterium bovis* Epidemiologically Important in East Africa.** *Journal of Bacteriology* 2011, **193:**670-678.

27.    Loiseau C, Menardo F, Aseffa A, Hailu E, Gumi B, Ameni G, Berg S, Rigouts L, Robbe-Austerman S, Zinsstag J, et al: **An African origin for *Mycobacterium bovis.*** *Evolution Medicine and Public Health* 2020**:**49-59.

28.    Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform.** *Bioinformatics* 2009, **25:**1754-1760.

29.    Harris SR, Feil EJ, Holden MT, Quail MA, Nickerson EK, Chantratita N, Gardete S, Tavares A, Day N, Lindsay JA, et al: **Evolution of MRSA during hospital transmission and intercontinental spread.** *Science* 2010, **327:**469-474.

30.    Price-Carter M, Brauning R, de Lisle GW, Livingstone P, Neill M, Sinclair J, Paterson B, Atkinson G, Knowles G, Crews K, et al: **Whole Genome Sequencing for Determining the Source of *Mycobacterium bovis* Infections in Livestock Herds and Wildlife in New Zealand.** *Frontiers in Veterinary Science* 2018, **5.**

31.    Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR: **SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments.** *Microb Genom* 2016, **2:**e000056.

32.    Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ: **IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies.** *Mol Biol Evol* 2015, **32:**268-274.

33.    Orloski K, Robbe-Austerman S, Stuber T, Hench B, Schoenbaum M: **Whole Genome Sequencing of *Mycobacterium bovis* Isolated From Livestock in the United States, 1989-2018.** *Frontiers in Veterinary Science* 2018, **5.**

34.    Crispell J, Zadoks RN, Harris SR, Paterson B, Collins DM, de-Lisle GW, Livingstone P, Neill MA, Biek R, Lycett SJ, et al: **Using whole genome**

940       **sequencing to investigate transmission in a multi-host system: bovine**
941       **tuberculosis in New Zealand.** *Bmc Genomics* 2017, **18**.

942   35.   Perea Razo CA, Rodriguez Hernandez E, Ponce SIR, Milian Suazo F, Robbe-
943       Austerman S, Stuber T, Canto Alarcon GJ: **Molecular epidemiology of cattle**
944       **tuberculosis in Mexico through whole-genome sequencing and**
945       **spoligotyping.** *PLoS One* 2018, **13:**e0201981.

946   36.   Sandoval-Azuara SE, Muniz-Salazar R, Perea-Jacobo R, Robbe-Austerman S,
947       Perera-Ortiz A, Lopez-Valencia G, Bravo DM, Sanchez-Flores A, Miranda-
948       Guzman D, Flores-Lopez CA, et al: **Whole genome sequencing of**
949       *Mycobacterium bovis* **to obtain molecular fingerprints in human and cattle**
950       **isolates from Baja California, Mexico.** *Int J Infect Dis* 2017, **63:**48-56.

951   37.   Salvador LCM, O'Brien DJ, Cosgrove MK, Stuber TP, Schooley AM, Crispell
952       J, Church SV, Grohn YT, Robbe-Austerman S, Kao RR: **Disease management**
953       **at the wildlife-livestock interface: Using whole-genome sequencing to**
954       **study the role of elk in** *Mycobacterium bovis* **transmission in Michigan,**
955       **USA.** *Mol Ecol* 2019, **28:**2192-2205.

956   38.   Trewby H, Wright D, Breadon EL, Lycett SJ, Mallon TR, McCormick C,
957       Johnson P, Orton RJ, Allen AR, Galbraith J, et al: **Use of bacterial whole-**
958       **genome sequencing to investigate local persistence and spread in bovine**
959       **tuberculosis.** *Epidemics* 2016, **14:**26-35.

960   39.   Glaser L, Carstensen M, Shaw S, Robbe-Austerman S, Wunschmann A, Grear
961       D, Stuber T, Thomsen B: **Descriptive Epidemiology and Whole Genome**
962       **Sequencing Analysis for an Outbreak of Bovine Tuberculosis in Beef**
963       **Cattle and White-Tailed Deer in Northwestern Minnesota.** *PLoS One* 2016,
964       **11:**e0145735.

965   40.   Crispell J, Cassidy S, Kenny K, McGrath G, Warde S, Cameron H, Rossi G,
966       MacWhite T, White PCL, Lycett S, et al: *Mycobacterium bovis* **genomics**
967       **reveals transmission of infection between cattle and deer in Ireland.**
968       *Microb Genom* 2020, **6**.

969   41.   Lasserre M, Fresia P, Greif G, Iraola G, Castro-Ramos M, Juambeltz A, Nunez
970       A, Naya H, Robello C, Berna L: **Whole genome sequencing of the**

971  monomorphic pathogen *Mycobacterium bovis* reveals local
972  differentiation of cattle clinical isolates. *BMC Genomics* 2018, **19:**2.

973  42. Hauer A, Michelet L, Cochard T, Branger M, Nunez J, Boschirolil ML, Biet F:
974  **Accurate Phylogenetic Relationships Among *Mycobacterium bovis***
975  **Strains Circulating in France Based on Whole Genome Sequencing and**
976  **Single Nucleotide Polymorphism Analysis.** *Frontiers in Microbiology* 2019,
977  **10**.

978  43. Andrievskaia O, Duceppe MO, Lloyd D: **Genome Sequences of Five**
979  ***Mycobacterium bovis* Strains Isolated from Farmed Animals and Wildlife**
980  **in Canada.** *Genome Announc* 2018, **6**.

981  44. Dippenaar A, Parsons SDC, Miller MA, Hlokwe T, Gey van Pittius NC, Adroub
982  SA, Abdallah AM, Pain A, Warren RM, Michel AL, van Helden PD: **Progenitor**
983  **strain introduction of *Mycobacterium bovis* at the wildlife-livestock**
984  **interface can lead to clonal expansion of the disease in a single**
985  **ecosystem.** *Infect Genet Evol* 2017, **51:**235-238.

986  45. Acosta F, Chernyaeva E, Mendoza L, Sambrano D, Correa R, Rotkevich M,
987  Tarte M, Hernandez H, Velazco B, de Escobar C, et al: ***Mycobacterium bovis***
988  **in Panama, 2013.** *Emerg Infect Dis* 2015, **21**.

989  46. Meiring C, Higgitt R, Dippenaar A, Roos E, Buss P, Hewlett J, Cooper D,
990  Rogers P, de Klerk-Lorist LM, van Schalkwyk L, et al: **Characterizing**
991  **epidemiological and genotypic features of *Mycobacterium bovis* infection**
992  **in wild dogs (Lycaon pictus).** *Transboundary and Emerging Diseases* 2020.

993  47. Kohl TA, Kranzer K, Andres S, Wirth T, Niemann S, Moser I: **Population**
994  **Structure of *Mycobacterium bovis* in Germany: a Long-Term Study Using**
995  **Whole-Genome Sequencing Combined with Conventional Molecular**
996  **Typing Methods.** *J Clin Microbiol* 2020, **58**.

997  48. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**
998  **sequence data.** *Bioinformatics* 2014, **30:**2114-2120.

999  49. Csardi G, Nepusz, T: **The igraph software package for complex network**
1000  **research.** *InterJournal, Complex Systems* 2006.

1001  50. Xu Y, Cancino-Munoz I, Torres-Puente M, Villamayor LM, Borras R, Borras-
1002  Manez M, Bosque M, Camarena JJ, Colomer-Roig E, Colomina J, et al: **High-**

1003        **resolution mapping of tuberculosis transmission: Whole genome**
1004        **sequencing and phylogenetic modelling of a cohort from Valencia Region,**
1005        **Spain.** *PLoS Med* 2019, **16:**e1002961.

1006   51.   Revell LJ: **phytools: an R package for phylogenetic comparative biology**
1007        **(and other things).** *Methods in Ecology and Evolution* 2012, **3:**217-223.

1008   52.   Drummond AJ, Suchard MA, Xie D, Rambaut A: **Bayesian Phylogenetics with**
1009        **BEAUti and the BEAST 1.7.** *Molecular Biology and Evolution* 2012, **29:**1969-
1010        1973.

1011   53.   Baele G, Lemey P, Bedford T, Rambaut A, Suchard MA, Alekseyenko AV:
1012        **Improving the Accuracy of Demographic and Molecular Clock Model**
1013        **Comparison While Accommodating Phylogenetic Uncertainty.** *Molecular*
1014        *Biology and Evolution* 2012, **29:**2157-2167.

1015   54.   Baele G, Li WLS, Drummond AJ, Suchard MA, Lemey P: **Accurate Model**
1016        **Selection of Relaxed Molecular Clocks in Bayesian Phylogenetics.**
1017        *Molecular Biology and Evolution* 2013, **30:**239-243.

1018   55.   Kass RE, Raftery AE: **Bayes Factors.** *Journal of the American Statistical*
1019        *Association* 1995, **90:**773-795.

1020   56.   Rieux A, Khatchikian CE: **TIPDATINGBEAST: an R package to assist the**
1021        **implementation of phylogenetic tip-dating tests using BEAST.** *Molecular*
1022        *Ecology Resources* 2017, **17:**608-613.

1023   57.   Didelot X, Fraser C, Gardy J, Colijn C: **Genomic Infectious Disease**
1024        **Epidemiology in Partially Sampled and Ongoing Outbreaks.** *Molecular*
1025        *Biology and Evolution* 2017, **34:**997-1007.

1026   58.   Plummer M, Best N, Cowles K, Vines K: **CODA: convergence diagnosis and**
1027        **output analysis for MCMC.** *R news* 2006, **6:**7-11.

1028   59.   Bouckaert R, Vaughan TG, Barido-Sottani J, Duchene S, Fourment M,
1029        Gavryushkina A, Heled J, Jones G, Kuhnert D, De Maio N, et al: **BEAST 2.5:**
1030        **An advanced software platform for Bayesian evolutionary analysis.** *Plos*
1031        *Computational Biology* 2019, **15**.

1032   60.   De Maio N, Wu CH, O'Reilly KM, Wilson D: **New Routes to Phylogeography:**
1033        **A Bayesian Structured Coalescent Approximation.** *Plos Genetics* 2015, **11**.

1034 61.  Pupko T, Pe'er I, Shamir R, Graur D: **A fast algorithm for joint reconstruction**
1035      **of ancestral amino acid sequences.** *Mol Biol Evol* 2000, **17:**890-896.

1036 62.  Yu GC, Smith DK, Zhu HC, Guan Y, Lam TTY: **GGTREE: an R package for**
1037      **visualization and annotation of phylogenetic trees with their covariates**
1038      **and other associated data.** *Methods in Ecology and Evolution* 2017, **8:**28-36.

1039 63.  Wang LG, Lam TT, Xu S, Dai Z, Zhou L, Feng T, Guo P, Dunn CW, Jones BR,
1040      Bradley T, et al: **Treeio: An R Package for Phylogenetic Tree Input and**
1041      **Output with Richly Annotated and Associated Data.** *Mol Biol Evol* 2020,
1042      **37:**599-603.

1043 64.  **Geosphere: spherical trigonometry**

1044 65.  Becker R, Wilks, A.R.: **Maps in S.** *AT&T Bell Laboratories Statistics Research*
1045      *Report* 1993, **93**.

1046 66.  Smith NH, Gordon SV, de la Rua-Domenech R, Clifton-Hadley RS, Hewinson
1047      RG: **Bottlenecks and broomsticks: the molecular evolution of**
1048      *Mycobacterium bovis*. *Nature Reviews Microbiology* 2006, **4:**670-681.

1049 67.  Coll F, McNerney R, Guerra-Assuncao JA, Glynn JR, Perdigao J, Viveiros M,
1050      Portugal I, Pain A, Martin N, Clark TG: **A robust SNP barcode for typing**
1051      *Mycobacterium tuberculosis* **complex strains.** *Nature Communications*
1052      2014, **5**.

1053 68.  Korber B, Muldoon M, Theiler J, Gao F, Gupta R, Lapedes A, Hahn BH,
1054      Wolinsky S, Bhattacharya T: **Timing the ancestor of the HIV-1 pandemic**
1055      **strains.** *Science* 2000, **288:**1789-1796.

1056 69.  Ramsden C, Melo FL, Figueiredo LM, Holmes EC, Zanotto PM, Consortium V:
1057      **High rates of molecular evolution in hantaviruses.** *Mol Biol Evol* 2008,
1058      **25:**1488-1492.

1059 70.  Rambaut A, Lam TT, Max Carvalho L, Pybus OG: **Exploring the temporal**
1060      **structure of heterochronous sequences using TempEst (formerly Path-O-**
1061      **Gen).** *Virus Evol* 2016, **2:**vew007.

1062 71.  Menardo F, Duchene S, Brites D, Gagneux S: **The molecular clock of**
1063      *Mycobacterium tuberculosis*. *PLoS Pathog* 2019, **15:**e1008067.

1064 72.  **Zoonotic Tuberculosis in Mammals , including Bovine and Caprine**
1065      **Tuberculosis Infections.** In.; 2019

1066 73.  Williams JE: **The tuberculin test in cattle.** *J Lancet* 1959, **79:**212-213.

1067 74.  Conlan AJ, McKinley TJ, Karolemeas K, Pollock EB, Goodchild AV, Mitchell
1068     AP, Birch CP, Clifton-Hadley RS, Wood JL: **Estimating the hidden burden of**
1069     **bovine tuberculosis in Great Britain.** *PLoS Comput Biol* 2012, **8:**e1002730.

1070 75.  Conlan AJK, Pollock EB, McKinley TJ, Mitchell AP, Jones GJ, Vordermeier M,
1071     Wood JLN: **Potential Benefits of Cattle Vaccination as a Supplementary**
1072     **Control for Bovine Tuberculosis.** *Plos Computational Biology* 2015, **11**.

1073 76.  Woodroffe R, Donnelly CA, Ham C, Jackson SYB, Moyes K, Chapman K,
1074     Stratton NG, Cartwright SJ: **Badgers prefer cattle pasture but avoid cattle:**
1075     **implications for bovine tuberculosis control.** *Ecology Letters* 2016,
1076     **19:**1201-1208.

1077 77.  Karolemeas K, McKinley TJ, Clifton-Hadley RS, Goodchild AV, Mitchell A,
1078     Johnston WT, Conlan AJ, Donnelly CA, Wood JL: **Recurrence of bovine**
1079     **tuberculosis breakdowns in Great Britain: risk factors and prediction.** *Prev*
1080     *Vet Med* 2011, **102:**22-29.

1081 78.  Brooks-Pollock E, Roberts GO, Keeling MJ: **A dynamic model of bovine**
1082     **tuberculosis spread and control in Great Britain.** *Nature* 2014, **511:**228-231.

1083 79.  Mekonnen GA, Ameni G, Wood JLN, Berg S, Conlan AJK, Aseffa A, Mihret A,
1084     Tessema B, Belachew B, Fekadu EW, et al: **Network analysis of dairy cattle**
1085     **movement and associations with bovine tuberculosis spread and control**
1086     **in emerging dairy belts of Ethiopia.** *Bmc Veterinary Research* 2019, **15**.

1087

1088

1089

1090

1091

1092

1093

1094

1095 **Supplementary Tables and Figures**

1096

1097 Supplementary Table 1: Model performance based on Maximum Likelihood Estimates (MLE)

1098 and Bayes Factors for all transmission clusters

| Transmission cluster | Model | log MLE | log Bayes Factor | Strength of Evidence (Kass & Raftery, 1995) |
|---|---|---|---|---|
| Cluster 1 | Relaxed exponential | -5536634 | 0.056 | Not worth more than a bare mention |
|  | Relaxed constant | -5536655 | 21.289 | Very strong |
|  | Strict exponential | -5536634 | - | - |
|  | **Strict constant** | **-5536659** | **25.277** | **Very strong** |
| Cluster 2 | Relaxed exponential | -5534141 | - | - |
|  | Relaxed constant | -5534156 | 14.906 | Very strong |
|  | Strict exponential | -5534146 | 4.800 | Strong |
|  | **Strict constant** | **-5534161** | **19.609** | **Very strong** |
| Cluster 3 | Relaxed exponential | -5535466 | - | - |
|  | Relaxed constant | -5535498 | 31.000 | Very strong |

|  | Strict exponential | -5535468 | 2.014 | Positive |
|---|---|---|---|---|
|  | **Strict constant** | **-5535505** | **39.717** | **Very strong** |
| **Cluster 4** | Relaxed exponential | -5534382 | 0.101 | Not worth more than a bare mention |
|  | Relaxed constant | -5534394 | 12.002 | Very strong |
|  | Strict exponential | -5534382 | - | - |
|  | **Strict constant** | **-5534396** | **13.793** | **Very strong** |
| **Cluster 5** | Relaxed exponential | -5534469 | - | - |
|  | Relaxed constant | -5534474 | 4.936 | Strong |
|  | Strict exponential | -5534470 | 1.661 | Positive |
|  | **Strict constant** | **-5534476** | **7.052** | **Very strong** |
| **Cluster 6** | Relaxed exponential | -5535426 | - | - |
|  | Relaxed constant | -5535452 | 25.594 | Very strong |
|  | Strict exponential | -5535428 | 2.099 | Positive |
|  | **Strict constant** | **-5535453** | **27.358** | **Very strong** |
| **Cluster 7** | Relaxed exponential | -5533415 | - | - |

|  | Relaxed constant | -5533417 | 2.257 | Positive |
|---|---|---|---|---|
|  | Strict exponential | -5533421 | 5.893 | Very strong |
|  | **Strict constant** | **-5533421** | **6.083** | **Very strong** |
| **Cluster 8** | Relaxed exponential | -5535839 | - | - |
|  | Relaxed constant | -5535863 | 23.890 | Very strong |
|  | Strict exponential | -5535841 | 2.074 | Positive |
|  | **Strict constant** | **-5535866** | **26.912** | **Very strong** |
| **Cluster 9** | Relaxed exponential | -5538088 | - | - |
|  | Relaxed constant | -5538147 | 59.427 | Very strong |
|  | Strict exponential | -5538088 | 0.389 | Not worth more than a bare mention |
|  | **Strict constant** | **-5538157** | **69.010** | **Very strong** |
| **Cluster 10** | Relaxed exponential | -5534887 | - | - |
|  | Relaxed constant | -5534891 | 3.526 | Strong |
|  | Strict exponential | -5534888 | 0.799 | Not worth more than a bare mention |

|  |  |  |  |  |
|---|---|---|---|---|
|  | **Strict constant** | **-5534894** | **7.076** | **Very strong** |
| **Cluster 11** | Relaxed exponential | -5533358 | 1.383 | Positive |
|  | Relaxed constant | -5533358 | 0.961 | Not worth more than a bare mention |
|  | Strict exponential | -5533357 | - | - |
|  | **Strict constant** | **-5533361** | **3.610** | **Strong** |
| **Cluster 12** | Relaxed exponential | -5533082 | 1.546 | Positive |
|  | Relaxed constant | -5533081 | - | - |
|  | **Strict exponential** | **-5533085** | **4.393** | **Strong** |
|  | Strict constant | -5533083 | 2.369 | Positive |

1099

1100

1101    Supplementary Table 2: Effective Sample Size (ESS) for TransPhylo parameters.

| **Transmission cluster** | **Sampling proportion pi** | **Within-host coalescent rate Ne*** | **Basic reproduction R** |
|---|---|---|---|
| **Cluster 1** | 1721 | 367 | 2407 |
| **Cluster 2** | 6027 | 1119 | 3400 |
| **Cluster 3** | 1890 | 391 | 3196 |
| **Cluster 4** | 4761 | 581 | 2649 |
| **Cluster 5** | 108783 | 933 | 10035 |
| **Cluster 6** | 2477 | 322 | 4910 |

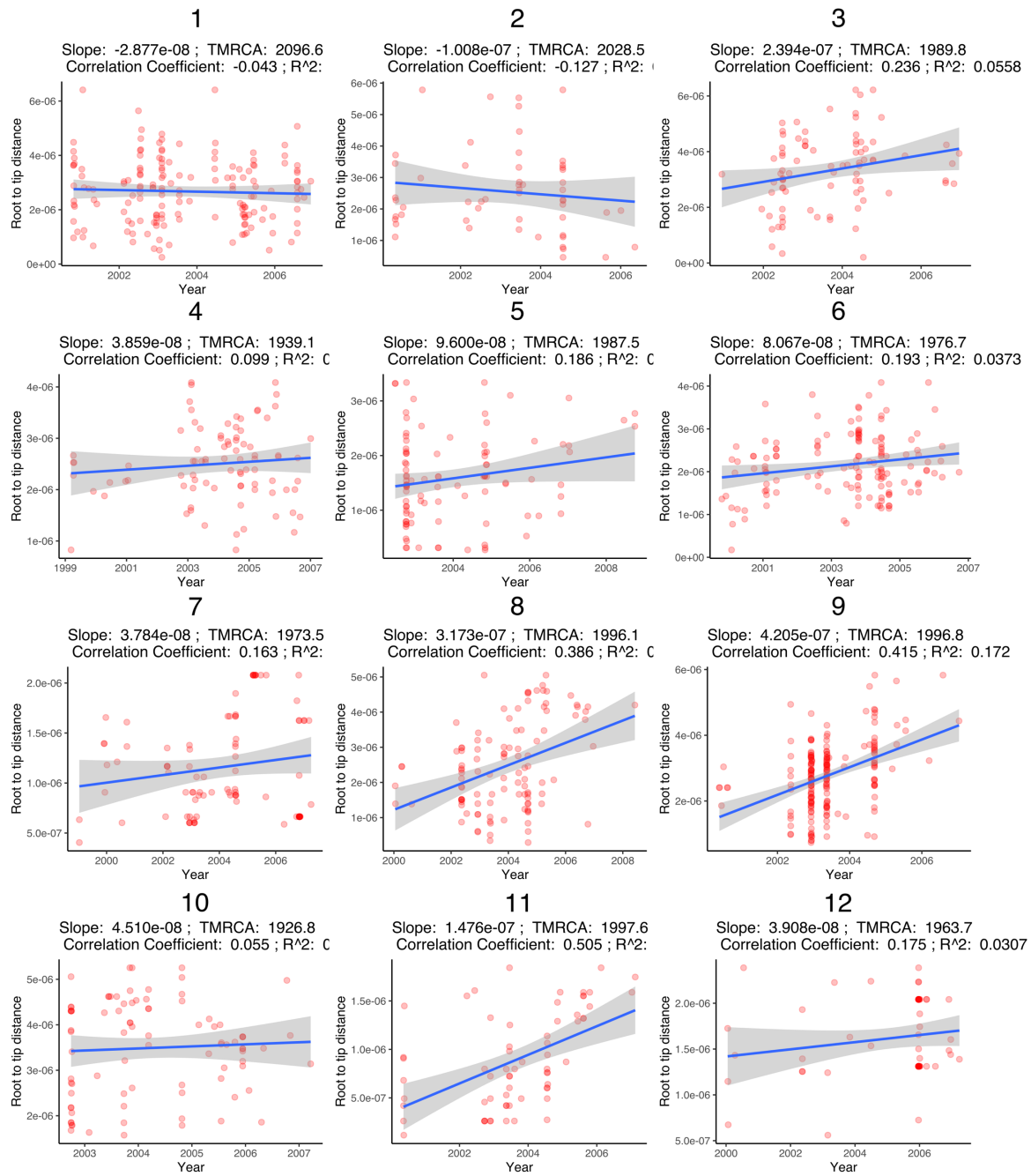| | | | |
|---|---|---|---|
| **Cluster 7** | 86600 | 568 | 3311 |
| **Cluster 8** | 3393 | 720 | 13012 |
| **Cluster 9** | 10640 | 202 | 17844 |
| **Cluster 10** | 2620 | 150 | 3533 |
| **Cluster 11** | 22204 | 1579 | 5094 |
| **Cluster 12** | 11576 | 8014 | 3128 |

1102

1103    Supplementary Table 3: Substitution rates for each transmission cluster

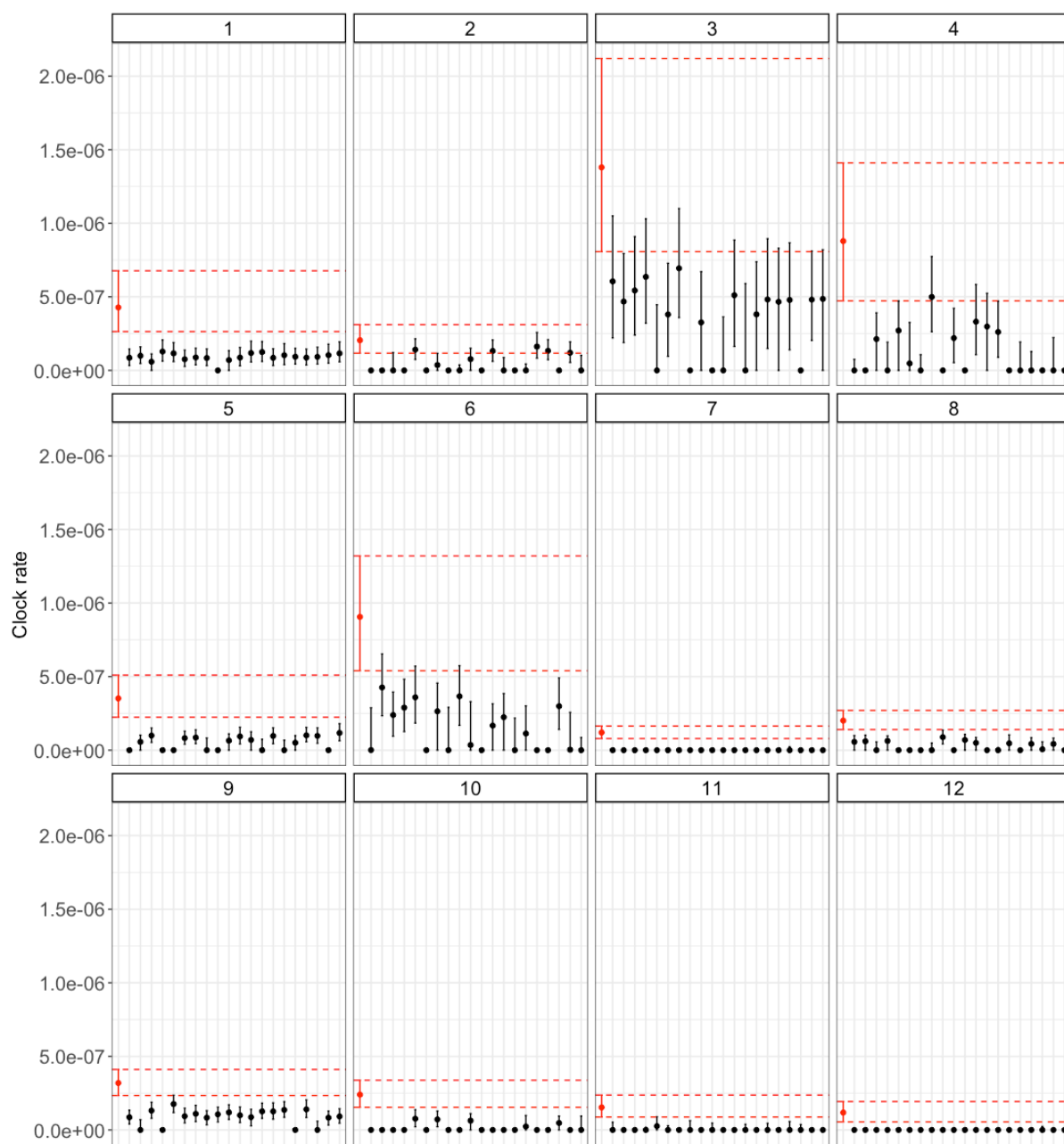| Transmission cluster | Median substitution rate (substitutions/site/year) | Median substitution rate (substitutions/genome/year) |
|---|---|---|
| **Cluster 1** | $4.3 \times 10^{-7}$ | 1.86 |
| **Cluster 2** | $2.1 \times 10^{-7}$ | 0.89 |
| **Cluster 3** | $1.4 \times 10^{-6}$ | 6.00 |
| **Cluster 4** | $8.8 \times 10^{-7}$ | 3.82 |
| **Cluster 5** | $3.5 \times 10^{-7}$ | 1.53 |
| **Cluster 6** | $9.1 \times 10^{-7}$ | 3.94 |
| **Cluster 7** | $1.2 \times 10^{-7}$ | 0.52 |
| **Cluster 8** | $2.0 \times 10^{-7}$ | 0.87 |
| **Cluster 9** | $3.2 \times 10^{-7}$ | 1.39 |
| **Cluster 10** | $2.4 \times 10^{-7}$ | 1.04 |
| **Cluster 11** | $1.5 \times 10^{-7}$ | 0.66 |
| **Cluster 12** | $1.2 \times 10^{-7}$ | 0.51 |

1104
1105

Supplementary Figure 1: Root to tip distances plotted against sampling dates for all isolates in each transmission cluster.
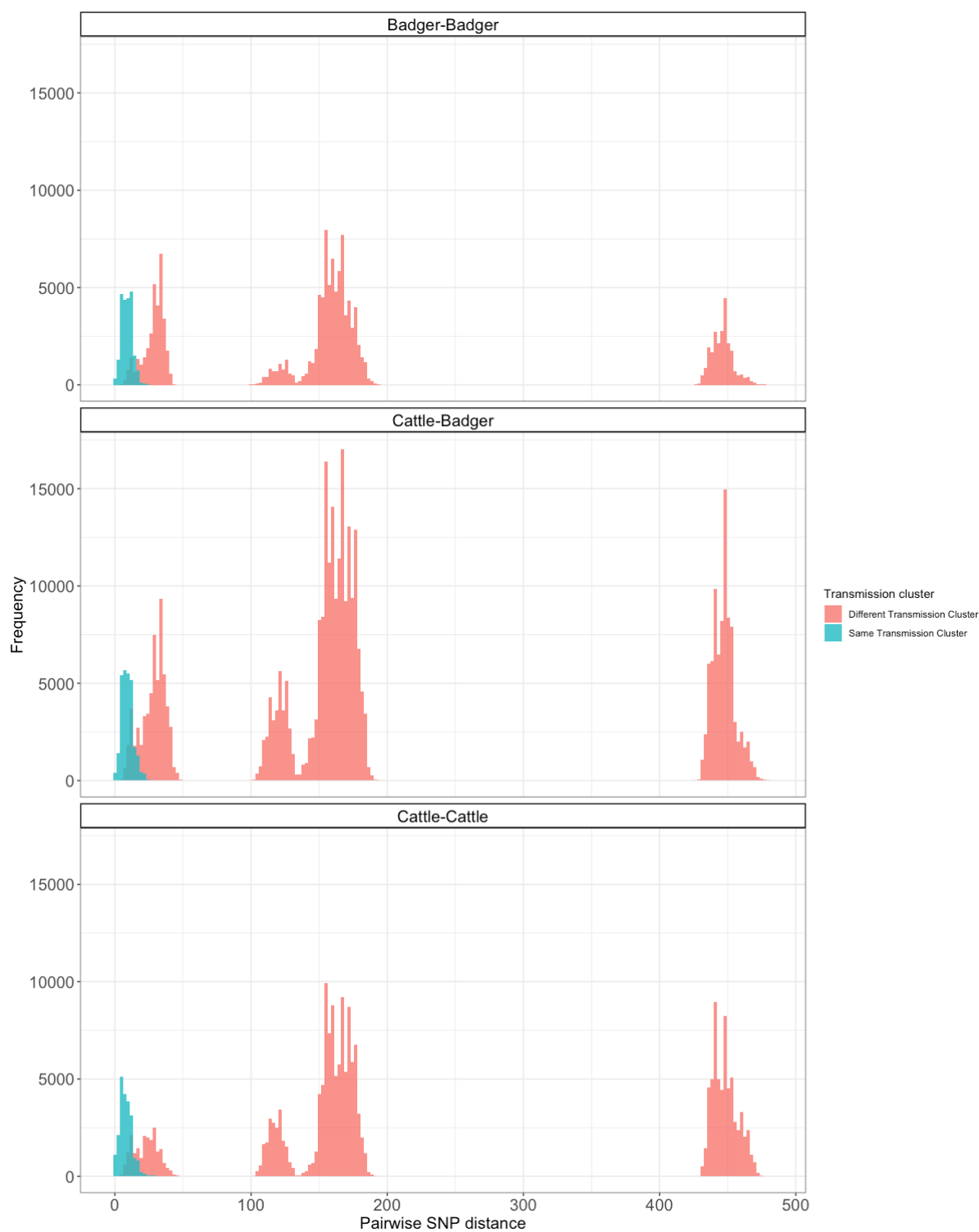
Supplementary Figure 2: Date randomization (DTR) analysis in BEAST for each transmission cluster. Estimated substitution rates (mean and highest posterior density) shown in red for the observed dataset and black for the randomized datasets.
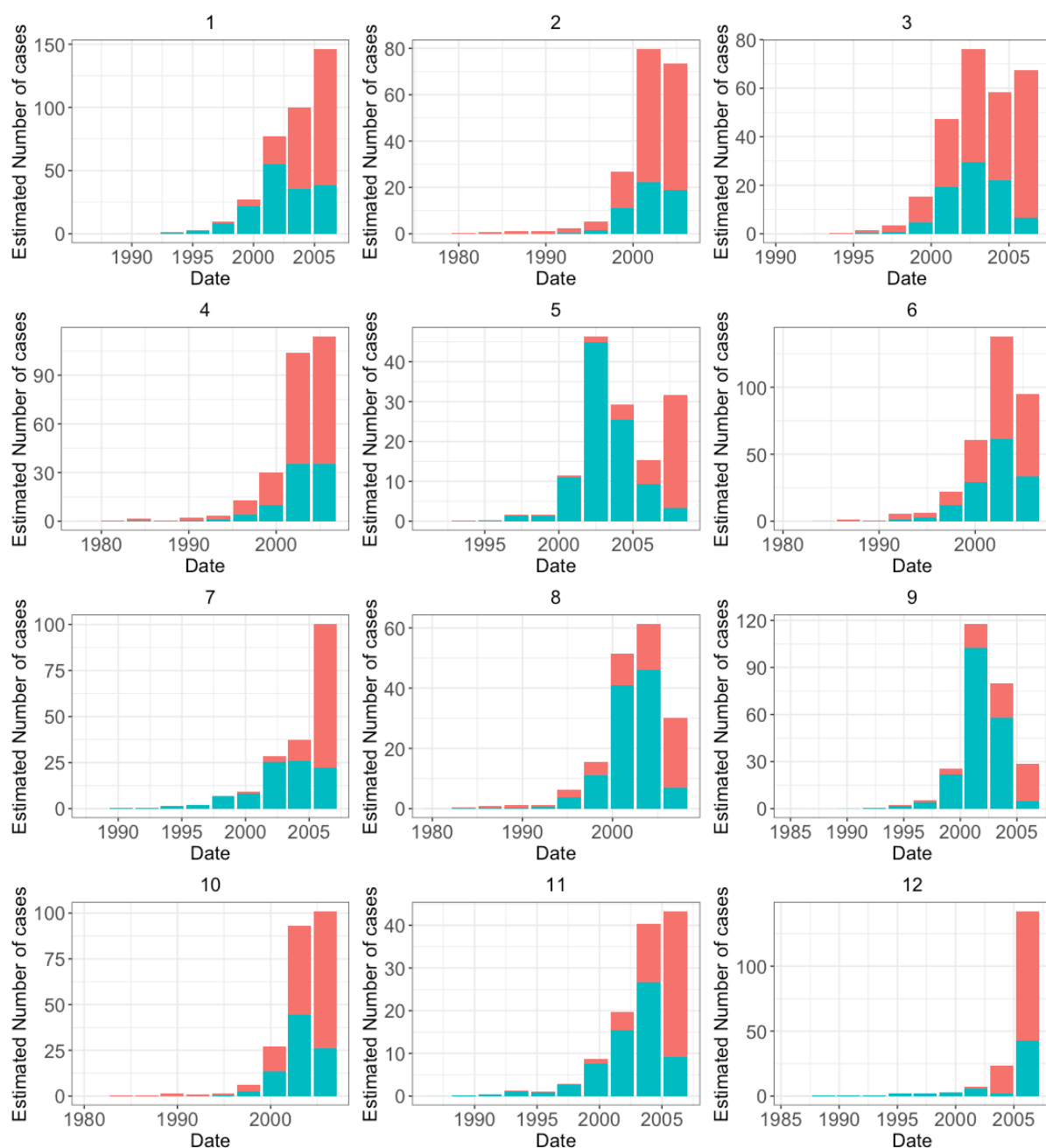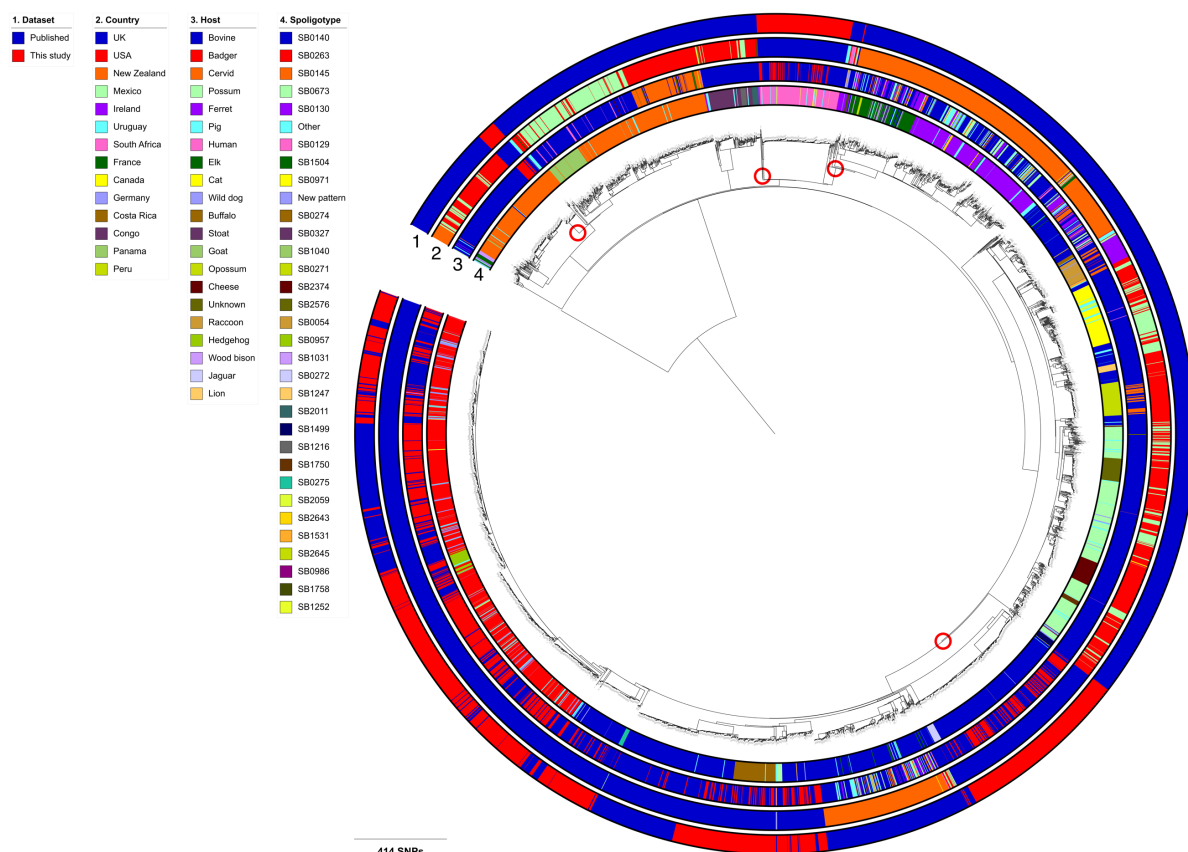
1117
1118    Supplementary Figure 3: Pairwise distance histograms for all samples, coloured by

1119    between/within transmission cluster and separated by host pair.

1120

1121
1122
1123    Supplementary Figure 4: Proportion of sampled and estimated unsampled cases for each

1124    transmission cluster. Sampled and unsampled cases are shown in red and blue respectively

1125
1126
1127    Supplementary Figure 5: Maximum likelihood phylogenetic tree of 4,281 *Mycobacterium*

1128    *bovis* Eu1 isolates rooted with a *M. caprae* isolate as the outgroup. Dataset, country, host

1129    and spoligotype are shown as datastrips around the outside of the phylogenetic tree.

1130    Potential introductions of Eu1 into England are highlighted with red circles.

1131