

1 **Deciphering the tangible spatio-temporal spread of a 25 years tuberculosis**  
2 **outbreak boosted by social determinants**

3 Mariana G. López<sup>1</sup>, M<sup>a</sup> Isolina Campos-Herrero<sup>2</sup>, Manuela Torres-Puente<sup>1</sup>, Fernando Cañas<sup>3</sup>,  
4 Jessica Comín<sup>4</sup>, Rodolfo Copado<sup>5</sup>, Penelope Wintringer<sup>6</sup>, Zamin Iqbal<sup>6</sup>, Eduardo Lagarejos<sup>2</sup>,  
5 Miguel Moreno-Molina<sup>1</sup>, Laura Pérez-Lago<sup>7</sup>, Berta Pino<sup>8</sup>, Laura Sante<sup>9</sup>, Darío García de  
6 Viedma<sup>7,10,†</sup>, Sofía Samper<sup>4,10,†,\*</sup>, Iñaki Comas<sup>1,11\*</sup>

7 <sup>1</sup>Tuberculosis Genomics Unit, Instituto de Biomedicina de Valencia (IBV-CSIC), Valencia 46010,  
8 Spain

9 <sup>2</sup>Servicio de Microbiología, Hospital Universitario de Gran Canaria Dr. Negrín, Las Palmas de  
10 Gran Canaria 35010, Spain

11 <sup>3</sup>Hospital Universitario Insular de Gran Canaria, Las Palmas de Gran Canaria 35016, Spain

12 <sup>4</sup>Instituto Aragonés de Ciencias de la Salud, Fundación IIS Aragón, Zaragoza 50009, Spain

13 <sup>5</sup>Hospital José Molina Orosa, Las Palmas de Gran Canaria 35500, Spain

14 <sup>6</sup>European Molecular Biology Laboratory – European Bioinformatics Institute, Hinxton CB10  
15 1SD, UK

16 <sup>7</sup>Servicio Microbiología Clínica y Enfermedades Infecciosas, Hospital General Universitario  
17 Gregorio Marañón, Instituto de Investigación Sanitaria Gregorio Marañón, Madrid 28007, Spain

18 <sup>8</sup>Hospital Nuestra Señora de la Candelaria, Santa Cruz de Tenerife 38010, Spain

19 <sup>9</sup>Hospital Universitario de Canarias, Santa Cruz de Tenerife 38320, Spain

20 <sup>10</sup>CIBER Enfermedades Respiratorias, Instituto de Salud Carlos III, Madrid 28029, Spain

21 <sup>11</sup>CIBER Epidemiología y Salud Pública, Instituto de Salud Carlos III, Madrid 28029, Spain

22 † these authors equally contributed

23 \* corresponding authors

## 24 **Abstract**

25 **Background.** Outbreak strains are good candidates to look for intrinsic transmissibility as they  
26 are responsible for a large number of cases with sustained transmission. However, assessment  
27 of the success of long-lived outbreak strains has been flawed by the use of low-resolution typing  
28 methods and restricted geographical investigations. We now have the potential to address the  
29 nature of outbreak strains by combining large genomic datasets and phylodynamic approaches.

30 **Methods.** We retrospectively sequenced the whole genome of representative samples assigned  
31 to an outbreak circulating in the Canary Islands (GC) since 1993; accounting for ~20% of local  
32 TB cases. We selected a panel of specific SNP markers to in-silico search for additional  
33 outbreak related sequences within publicly-available TB genomic data. Using this information  
34 we inferred the origin, spread and epidemiological parameters of the GC-outbreak.

35 **Findings.** Our approach allowed us to accurately trace both the historical and recent dispersion  
36 of the strain. We evidenced its high success within the Canarian archipelago but found a limited  
37 expansion abroad. Estimation of epidemiological parameters from genomic data contradicts a  
38 distinct biology of the GC-strain.

39 **Interpretation.** With the increasing availability of genomic data allowing for an accurate  
40 inference of strain spread and key epidemiological parameters, we can now revisit the link  
41 between *Mycobacterium tuberculosis* genotypes and transmission, as routinely done for SARS-  
42 CoV-2 variants of concern. We show that the success of the GC-strain is better explained by  
43 social determinants rather than intrinsically higher bacterial transmissibility. Our approach can  
44 be used to trace and characterize strains of interest worldwide.

45 **Funding.** European Research Council (101001038-TB-RECONNECT), the Ministerio de  
46 Economía, Industria y Competitividad (PID2019-104477RB-I00), Instituto de Salud Carlos III  
47 (FIS18/0336), European Commission –NextGenerationEU (Regulation EU 2020/2094), through  
48 CSIC's Global Health Platform (PTI Salud Global) to IC. Gobierno de Aragón/Fondo Social  
49 Europeo "Construyendo Europa desde Aragón" to SS

50 **Keywords:** Tuberculosis, outbreak, whole-genome sequencing, phylodynamics, genomic  
51 epidemiology

## 52 **Research in context**

### 53 **Evidence before this study**

54 Identification of intrinsically highly transmissible strains of *Mycobacterium tuberculosis* remains  
55 elusive. Among candidates are those strains that have been thriving in a community for decades  
56 representing a significant contribution to the long-term local TB burden. These long-lived  
57 outbreak strains have been identified in different parts of the world and the speculation is that  
58 their success is linked to higher transmissibility. Several studies have attempted to analyze the  
59 epidemiological characteristics of these strains as well as their genomic composition to look for  
60 potential transmission determinants. However those studies are usually circumscribe to their  
61 original geographic boundaries. By contrast, this transmissibility should be replicated in different  
62 parts of the world, a lesson learnt from SARS-CoV-2 variants of concern. Previous attempts  
63 failed to examine the success of these outbreak strains at a global scale. Thus, it is unknown  
64 whether the long-lived outbreak strains had a similar or different trajectory in other countries,  
65 casting doubts about their transmissibility potential.

### 66 **Added value of this study**

67 Here we analyzed a strain causing a long-lived outbreak in the Canary Islands since 1993 using  
68 whole genome sequencing. As in previous studies with other similar outbreak strains, we  
69 analyzed the diversity and phylodynamics of the outbreak in the area where it was originally  
70 described. However, thanks to the possibility of interrogating the entire European Nucleotide  
71 Archive, we had the unique chance to look at the spread of the strains beyond its original  
72 geographic boundaries. This approach allowed us to comprehensively trace the real spatio-  
73 temporal spread of the outbreak from the emergence of its ancestor about 700 years ago to its  
74 recent transmission outside the Canary Islands. As a result, there is limited evidence for similar  
75 success of the strains outside Canary Islands. Furthermore, we complemented the analysis with  
76 epidemiological data of the early cases and with phylodynamic analysis to estimate key  
77 epidemiological parameters linked to the strain spread. All evidence strongly suggests that  
78 factors related to the host, instead of the bacteria, are behind the persistence and expansion of  
79 the outbreak strain.

## 80 **Implications of all the available evidence**

81 Infectious disease outbreaks are a major problem for public health. Tracing outbreak expansion  
82 and knowing the main factors behind their emergence and persistence are key to an effective  
83 disease control. Our study allows researchers and public health authorities to use WGS-based  
84 methods to trace outbreaks, and include available epidemiological information to evaluate the  
85 factors underpinning outbreak persistence. Taking advantage of all the information freely  
86 available in public repositories, researchers can accurately establish the expansion of the  
87 outbreak behind its original boundaries; and they can determine the potential risk of the strain to  
88 inform health authorities which, in turn, can define target strategies to mitigate its expansion and  
89 persistence. Finally, we show the need to evaluate strain transmissibility in different geographic  
90 contexts to unequivocally associate its spread to local or pathogen factors, a major lesson taken  
91 from SARS-CoV-2 genomic surveillance.

## 92 **Introduction**

93 Tuberculosis (TB) has been the first cause of death by an infectious disease for years  
94 surpassing HIV according to the World Health Organization (WHO). In 2019 were reported 10  
95 million new TB cases and 1.4 million deaths [1], with these numbers likely to increase due to  
96 the COVID-19 pandemic [2]. Outbreaks, defined as the concentration of an abnormal number of  
97 disease cases in space and time, are very common in TB. While outbreaks have generally been  
98 assumed to be short-lived, genotyping has been able to identify outbreak strains that became  
99 highly prevalent in specific regions and end up transmitting over decades. These outbreak  
100 strains can be found worldwide and across MTBC lineages and are usually drivers of local TB  
101 burden [3–8]. However, success of outbreak strains has been largely evaluated only in their  
102 original location, raising the question of whether those strains have any intrinsic transmissibility  
103 advantage or if their success is associated with local population processes like founder effects  
104 or ecological drivers of transmission [3–8].

105 In this work we investigate in detail one of those outbreak strains and compare them to others,  
106 previously published, to learn about the origin and epidemiology of long-lived strains. In 2001 a  
107 TB outbreak was identified in the Canary Islands. 651 strains were retrospectively analyzed  
108 between 1991-1996 in Gran Canaria (GC), the most populated island of the archipelago, using  
109 RFLP IS6110. A big cluster of 75 isolates was recognized with the first 10 cases diagnosed in

110 1993. The likely index case was a Liberian refugee who arrived on the island 6 months before  
111 diagnosis, in July of 1993 [9]. The strain belongs to the Beijing genotype, corroborated by  
112 spoligotyping, and was named GC1237. More recently, three GC-outbreak single nucleotide  
113 polymorphisms (SNPs) were selected and used to design a specific PCR to rapidly identify  
114 secondary cases belonging to GC-outbreak [10]. The fast spread of the GC strain to the other  
115 islands of the Canarian archipelago was evidenced by using different molecular typing methods,  
116 and also new cases outside the archipelago were identified up to 2014 [10–12]. Until now, GC-  
117 outbreak has never been studied by exploring the information obtained from whole genome  
118 sequencing (WGS).

119 In this work, we track the outbreak strain in its original geographic boundaries but also  
120 elsewhere by querying the entire European Nucleotide Archive (ENA) database. Then, we  
121 combine epidemiological and sequencing data, and apply a phylodynamic analysis to trace the  
122 origin of the GC outbreak strain, to track its spread, to understand its dynamics and to define the  
123 factors that underpin its expansion.

## 124 **Materials and method**

### 125 *Study population*

126 *M. tuberculosis* DNA samples from 80 different patients, 3 additional ones from the index case,  
127 all from Canary Islands, and 1 from Zaragoza were supplied by the Instituto de Investigación  
128 Sanitaria de Aragón. In addition, DNA samples from 5 different patients from the peninsula  
129 (Madrid) were supplied by the Hospital Universitario Gregorio Marañón, 2 of which belonged to  
130 recent immigrants from Guinea (Table S1). A representative number of samples from the  
131 Canary Islands, ranging the defined period of the outbreak (1993-2014), was selected. All  
132 samples were previously identified as part of the Gran Canaria (GC) outbreak by genotyping  
133 methods.

134 Epidemiological information was obtained from Servicio de Microbiología (Hospital Universitario  
135 de Gran Canaria Dr. Negrín) under the Ethics Committee Application CEIm H.U.G.C. Dr.  
136 Negrín 2019-502-1.

### 137 *Study design*

138 The whole genome of *M. tuberculosis* from the first set of 89 samples (86 patients) was  
139 sequenced. Additional samples were looked for by querying our local database of sequences  
140 [13–15], with the three GC outbreak specific SNPs [10] previously selected; Rv2524 (C1398T,

141 SNP2847935), Rv3869 (G1347C, SNP4346385) and Rv0926c (G162A, SNP1033625).  
142 Redefinition of specific SNPs was conducted based on phylogeny and distance analysis. New  
143 samples were identified by inspecting the whole ENA repository and samples from ongoing  
144 projects. Phylogeographic and phylodynamics analysis were performed in order to evaluate the  
145 outbreak. Comparison with other outbreaks was performed; sequences corresponding to  
146 Denmark (PRJEB20214)[4]; Thailand (PRJN244659)[3]; Bern (PRJEB5925)[6]; and Buenos  
147 Aires (BA, PRJEB7669)[5] outbreaks were downloaded from ENA and analyzed with the  
148 bioinformatics pipeline detailed below.

#### 149 *Whole-Genome Sequencing and bioinformatics analysis*

150 DNA samples were used to prepare sequencing libraries with a Nextera XT DNA library  
151 preparation kit (Illumina), following the manufacturer's instructions. Sequencing was performed  
152 on an Illumina MiSeq instrument, applying a 2 x 300bp paired-end chemistry. General  
153 bioinformatics analysis is described in  
154 <https://gitlab.com/tbgenomicsunit/ThePipeline/-/tree/master/>. Briefly, read files were trimmed  
155 and filtered with fastp [16]. Kraken software was used to remove nonMTBC (non-  
156 Mycobacterium tuberculosis complex) reads [17]. Sequences were then mapped to an inferred  
157 MTBC common ancestor genome (<https://doi.org/10.5281/zenodo.3497110>) using bwa [18].  
158 SNPs were called with SAMtools [19] and VarScan2 [20]. GATK HaplotypeCaller [21] was used  
159 for InDels calling. SNPs with a minimum of 20 reads (20X) in both strands and quality 20 were  
160 selected. InDels with less than 20X were discarded. SnpEff was used for SNP annotation using  
161 the H37Rv annotation reference (AL123456.2). Finally, SNPs falling in genes annotated as  
162 PE/PPE/PGRS, 'maturase', 'phage', '13E12 repeat family protein'; those located in insertion  
163 sequences; those within InDels or in higher density regions (>3 SNPs in 10 bp) were removed  
164 due to the uncertainty of mapping. Heterogenous SNPs (hSNPs) were obtained from filtered  
165 SNP files, they were classified as positions where > 5% and < 90% of the reads were the  
166 alternative allele [22]; we looked for hSNPs only for positions for which at least one sample  
167 harbor the SNP fixed, i.e. where > 90% of reads were the alternative allele. Lineages were  
168 determined by comparing called SNPs with specific phylogenetic positions established [23,24].  
169 An in-house R script was used to detect mixed infections based on the frequency of lineage-  
170 and sublineage-specific positions [15]. Pairwise distances were computed with the R *ape*  
171 package, statistical analysis and graphics were also performed with R. Raw sequencing data  
172 are available under the accession number PRJEB50491 (ENA).

173 All WGS Illumina sequencing runs stored in the ENA with metadata identifying them as *M.*  
174 *tuberculosis* complex as of July 2018 were downloaded (N=38075). A de Bruijn graph (k=31)  
175 was built from each with Cortex v1.0.5.21  
176 (<https://github.com/iqbal-lab/cortex/releases/tag/v1.0.5.21>), and sequencing errors were  
177 removed by excluding low coverage unitigs from the graph (threshold is sample dependent and  
178 automatically chosen by the software). The remaining kmers from these graphs were then used  
179 to build a Bitsliced Genomic Signature Index (BIGSI)[25]. The index was used to query these  
180 ENA samples for the outbreak-related SNPs, by first creating two 61 base-pair probes for each  
181 SNP (30bp flanking each side of the SNP, from the reference genome, and one probe for each  
182 SNP allele), returning binary information as to which sequencing runs contained all the kmers in  
183 the probes.

#### 184 *Phylogenetic analysis*

185 Multisequence alignment (MSA) files were constructed with concatenated SNPs discarding well-  
186 known drug resistance and invariant positions. Maximum likelihood trees were inferred with  
187 RAxML v8.2.11 [26] using GTRCATI model and 1000 fast-bootstrap replicates. Tree  
188 visualization and editing were conducted in ITOL (<https://itol.embl.de/>). Specific SNPs were  
189 identified using likelihood ancestral reconstruction of Mesquite software v3.61  
190 (<http://www.mesquiteproject.org>). Genomic network was constructed with the MSA files with a  
191 median joining network inference method implemented in PopArt Software [27].

#### 192 *Time and geographic origin of the outbreak*

193 A set of 200 samples including the outbreak, closest clades and a representative subset of  
194 global samples selected with Treemer [28] was used to estimate the time of the most common  
195 ancestor (tMRCA) of the outbreak and deeper nodes. The evolutionary history of the outbreak  
196 was reconstructed with BEAST2 v2.5.1 [29] using GTR (gamma 4) as site model, coalescent  
197 constant population as tree prior and relaxed clock log-normal with gamma distribution as prior.  
198 We use a tip dating method and also set the clock rate (mean =  $4.6 \times 10^{-8}$  substitutions per site  
199 per year with 95% HPD:  $3.3 \times 10^{-8}$  to  $6.2 \times 10^{-8}$ ), based on previous publications [30].  
200 Ascertainment bias was corrected by adjusting the clock rate based on the size of the alignment  
201 [14]. MCMC's chain length of 10M with sampling every 1000 steps for the posterior distribution.  
202 Three independent runs were performed to reach convergence. Log and tree files were  
203 combined with LogCombiner tool discarding 10% of burn-in. Combined files were inspected

204 with Tracer v1.7.1 [31], all parameters reached effective sample size (ESS) > 200 and well  
205 mixing. Tree was annotated with TreeAnnotator and visualized with FigTree v1.4.3.

206 Spatial dispersal dynamics of the outbreak and closest clades was reconstructed using a  
207 phylogeographic diffusion in discrete space approach implemented in BEAST v2.5.1 (Beast-  
208 classic package) [32]. A set of 142 samples, including the whole outbreak and closest clades,  
209 and those global strains with available information of year and location were used. Same clock  
210 and tree model, as well as priors, as in dating analysis were used. Same MCMC chains and  
211 further analysis as dating were conducted. In addition, SPREAD3 tool was used for geographic  
212 visualization of outbreak evolutionary history [33], only dispersion patterns with posterior values  
213 higher than 0.9 were plotted.

#### 214 *Phylodynamics of the outbreak*

215 The dynamics of the outbreak was studied with the Bayesian Birth-Death Skyline model  
216 implemented in BEAST2 v2.5.1. This model allows estimating parameters related to the  
217 evolution of the outbreak [34]. The analysis was conducted with all the samples included in the  
218 outbreak (64 samples, Table S1). We use a GTR (gamma 4) as site model, a strict6000 clock  
219 with LogNormal prior with the mean previously published [30] as prior. A Birth death skyline  
220 serial tree model was used. Specific parameters were set as follow (i) *becoming uninfected*  
221 (LogNormal; M=0, S=1.25), considering the global estimation of 1 year as TB infectiousness  
222 period and not taking into account latency [35] (ii) *origin* (LogNormal; M= 30; S=0.04) it was  
223 taking into account that the likely index case has arrived to the archipelago in the beginning of  
224 1993 starting immediately the outbreak; (ii) *reproductive number* (LogNormal; M=0; S= 1.25;  
225 dimension = 13) it was considered that in settings with low TB burden each person can  
226 generate a maximum of 12 secondary cases [35]; (iv) *sampling proportion* (beta; alpha = 3; beta  
227 = 50) it was calculated as the total number of TB cases in Canary Islands between between  
228 1993-2017 ([https://www3.gobiernodecanarias.org/sanidad/scs/listaImagenes.jsp?  
229 idDocumento=fe9e3e94-fee-11e0-ab85-376c664a882a&idCarpeta=0f67aaf7-9d88-11e0-b0dc-  
230 e55e53ccc42c](https://www3.gobiernodecanarias.org/sanidad/scs/listaImagenes.jsp?idDocumento=fe9e3e94-fee-11e0-ab85-376c664a882a&idCarpeta=0f67aaf7-9d88-11e0-b0dc-e55e53ccc42c)) and considering that 20-30% of cases belong to the outbreak resulting in a  
231 sampling range of 8-13%. MCMC's chain length of 10M was run with sampling every 1000 steps  
232 for the posterior distribution. Log file was inspected with Tracer v1.7.1 [31], all parameters  
233 reached ESS>200 and well mixing. Results were inspected and plotted with R package  
234 *bdskytools* (<https://github.com/laduplessis/bdskytools>).

235 Clocklike structure of the outbreak dataset was first evaluated with a linear regression analysis  
236 of tip dates vs root-to-tip distances with TempEst [36] and by applying a date randomization test  
237 (DRT) with 100 randomized replicates of the BDSKY. The clocklike structure was evaluated with  
238 the analysis proposed by [37] considering that DRT is passed when the clock rate estimate for  
239 the observed data does not overlap with the range of estimates obtained from the randomized  
240 sets; intermediate DRT is passed when the clock rate estimate for the observed data does not  
241 overlap with the confidence intervals of the estimates obtained from the randomized sets, and  
242 stringent DRT is passed when the confidence interval of the clock rate estimate for the observed  
243 data does not overlap with the confidence intervals of the estimates obtained from the  
244 randomized sets.

#### 245 *Data and code availability*

246 All sequences have been deposited on the ENA repository under the project number  
247 PRJEB50491. All the scripts, tools and reference sequences are included in the repository  
248 section of the laboratory webpage ([http://tgu.ibv.csic.es/?page\\_id=1794](http://tgu.ibv.csic.es/?page_id=1794)) and the gitlab page  
249 (<https://gitlab.com/tbgenomicsunit/ThePipeline>). Project numbers of the additional sequences  
250 from other studies are listed in Tables S2 and S4. Any additional information required to  
251 reanalyze the data reported in this paper is available from the lead contact upon request.

## 252 **Results**

### 253 *Genomic delineation of the GC outbreak*

254 We retrieved 86 samples from outbreak patients previously typed by molecular methods, plus  
255 3 additional samples from the likely index case collected in different years. Out of those, we  
256 discarded 24 samples which did not achieve sufficient DNA quality for sequencing, leaving us  
257 with a total of 65 usable samples (62 patients, ~10% of the outbreak, Table S1). The GC  
258 outbreak was delineated following the diagram detailed in Figure 1.

259 First, all sequences were queried for the 3 SNP previously defined as GC outbreak markers  
260 (partial SNPs profile)[10]. We observed that 2 isolates (GC466, GC515), previously included in  
261 the outbreak based on molecular typing methods, did not actually harbor those marker SNPs,  
262 thus both were excluded from further analysis. In addition, we queried our collection of  
263 sequences and found 7 additional samples, 2 from Valencia Region and 5 from Liberia (Table  
264 S2). We then constructed a maximum likelihood (ML) tree with all the samples with the marker  
265 SNPs, and observed a differentiated and well-supported monophyletic clade (bootstrap = 99,

266 Figure 2A) including the index case, most of the isolates previously assigned to outbreak and  
267 the Valencian cases (60 single TB cases, Table S1). The phylogeny also revealed that the  
268 SNPs initially considered as GC outbreak-specific, were not exclusive, since additional global  
269 strains also harbour them.

270 Pairwise distance was also estimated considering the phylogenetic outbreak definition based on  
271 WGS data (Table 1), in order to evaluate if distances were in agreement with the accepted  
272 recent transmission thresholds of 0-10 SNPs. Mean and median distance within-outbreak were  
273 4.3 and 4 SNPs (range 0-17) respectively, while distance between outbreak and non-outbreak  
274 samples were 33 and 30 SNPs (range 18-63).

Outbreak	SNP pairwise distance					
	Min	Max	Q1	Median	Mean	Q2
GC	0	17	2	4	4.3	6
BA	0	31	5	7	8.2	10
Bern	0	18	3	5	5.4	7
Denmark	0	33	6	8	9.9	13
Thailand	0	25	9	12	11.8	15

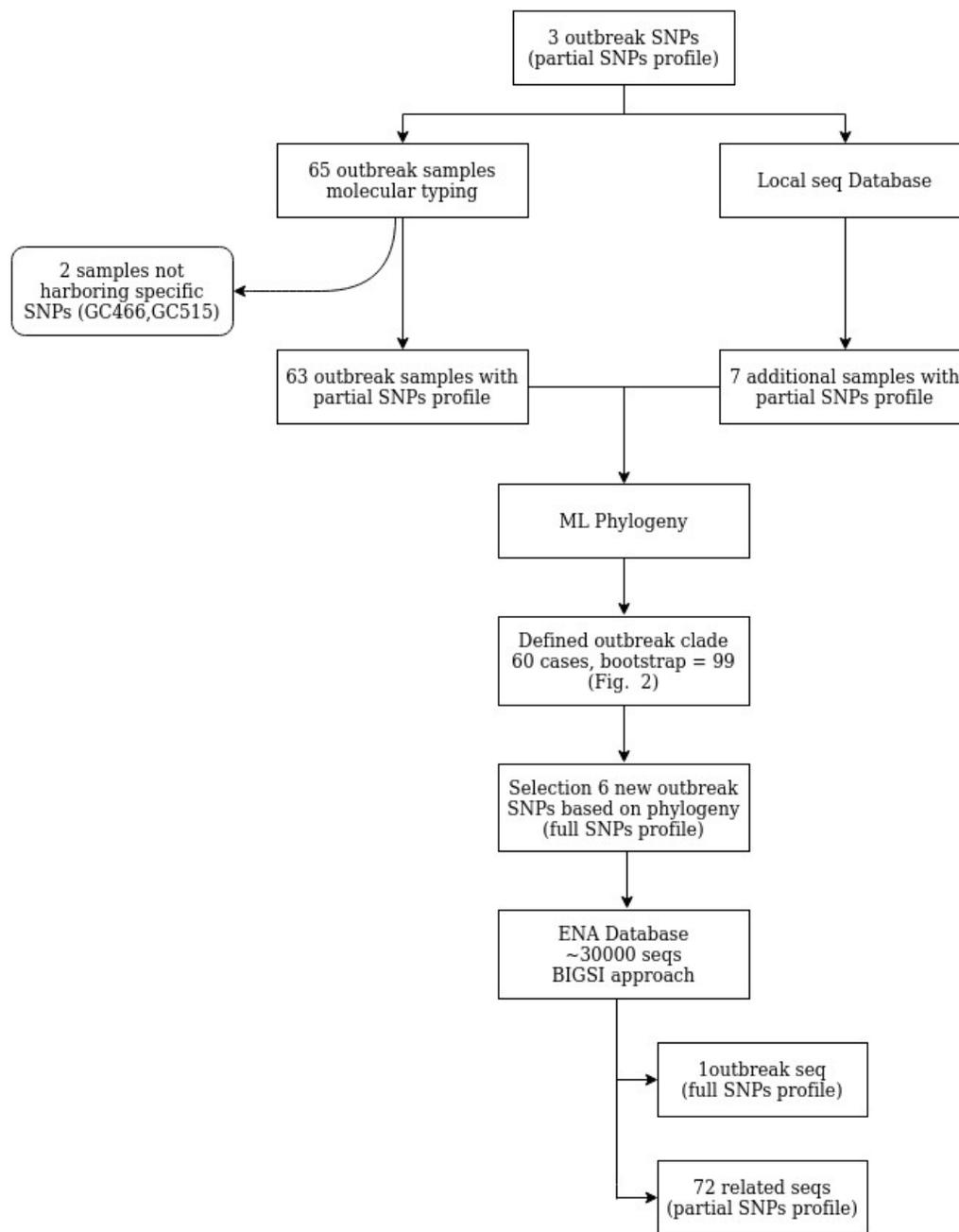
275 **Table 1.** Pairwise distance comparison among different outbreaks. Minimum (Min), maximum  
276 (Max), mean, median, 1<sup>st</sup> (Q1) and 2<sup>nd</sup> (Q2) quartile are provided

277 When compared to other known outbreak strains, the GC one displayed the lowest within-  
278 outbreak mean and median pairwise distance, though similar to the Bern outbreak (Table 1,  
279 Figure 2B-C). GC and Bern outbreaks exhibited a unimodal right-skewed pairwise genetic  
280 distance distribution, pointing out that most of the samples have low pairwise SNP distance  
281 among them. On the contrary, Denmark and Buenos Aires (BA) had a slightly bimodal  
282 distribution, suggesting the existence of more than one distinctive clade within both outbreaks,  
283 as proposed elsewhere [4]. Thailand outbreak displayed a unimodal normal distribution  
284 indicating a high pairwise SNP dispersion.

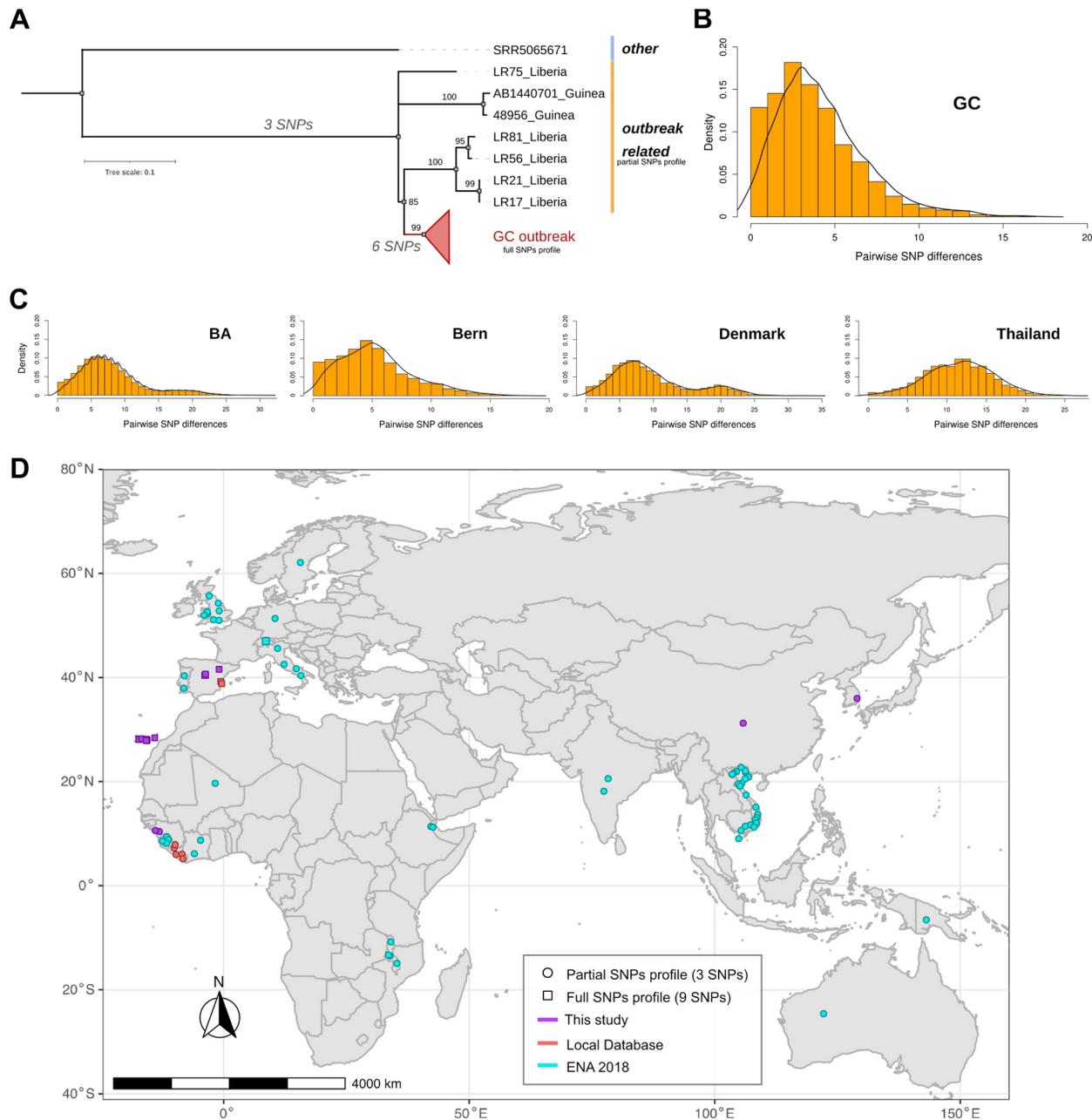
285 Thus, distance values obtained for GC also support the phylogenetic circumscription of the  
286 outbreak. With this delimitation, two additional samples from Madrid (AB1440701, 48956),  
287 belonging to recent Guinean immigrants, and previously identified as infected by the GC1237  
288 strain [10], were excluded (Table S1).

289 *Evaluation of success outside the outbreak GC setting*

290 In order to know if the GC strain can be found outside the Canary Islands, and if so if it has a  
291 similar epidemiological success, we designed an approach to evaluate all the MTB sequences  
292 available in public databases. This gives us the opportunity to study the spread of the outbreak,  
293 not only to the study area, but inspect its broader extension. First, the phylogenetic  
294 circumscription allows us to identify six additional markers defining the outbreak clade, most of  
295 which were missense mutations located in genes involved in metabolic processes and  
296 respiration, with effects likely unrelated with virulence or transmission (Table S3, Figure 2A).  
297 Thus, the full SNPs profile (9 SNPs) was used to query every *M. tuberculosis* sample deposited  
298 in ENA by July 2018, using the BIGSI index [25] (see methods). We identified 1 sequence  
299 meeting the full SNPs profile, from Switzerland, thus being part of the outbreak, and 72  
300 harboring the partial SNPs profile, thus being related to the outbreak (Table S2, Figure 2D). We  
301 consider these sequences in further analysis, as they can shed light on the remote origin of the  
302 outbreak strain.



303 **Figure 1.** Workflow detailing the GC outbreak delineation procedure



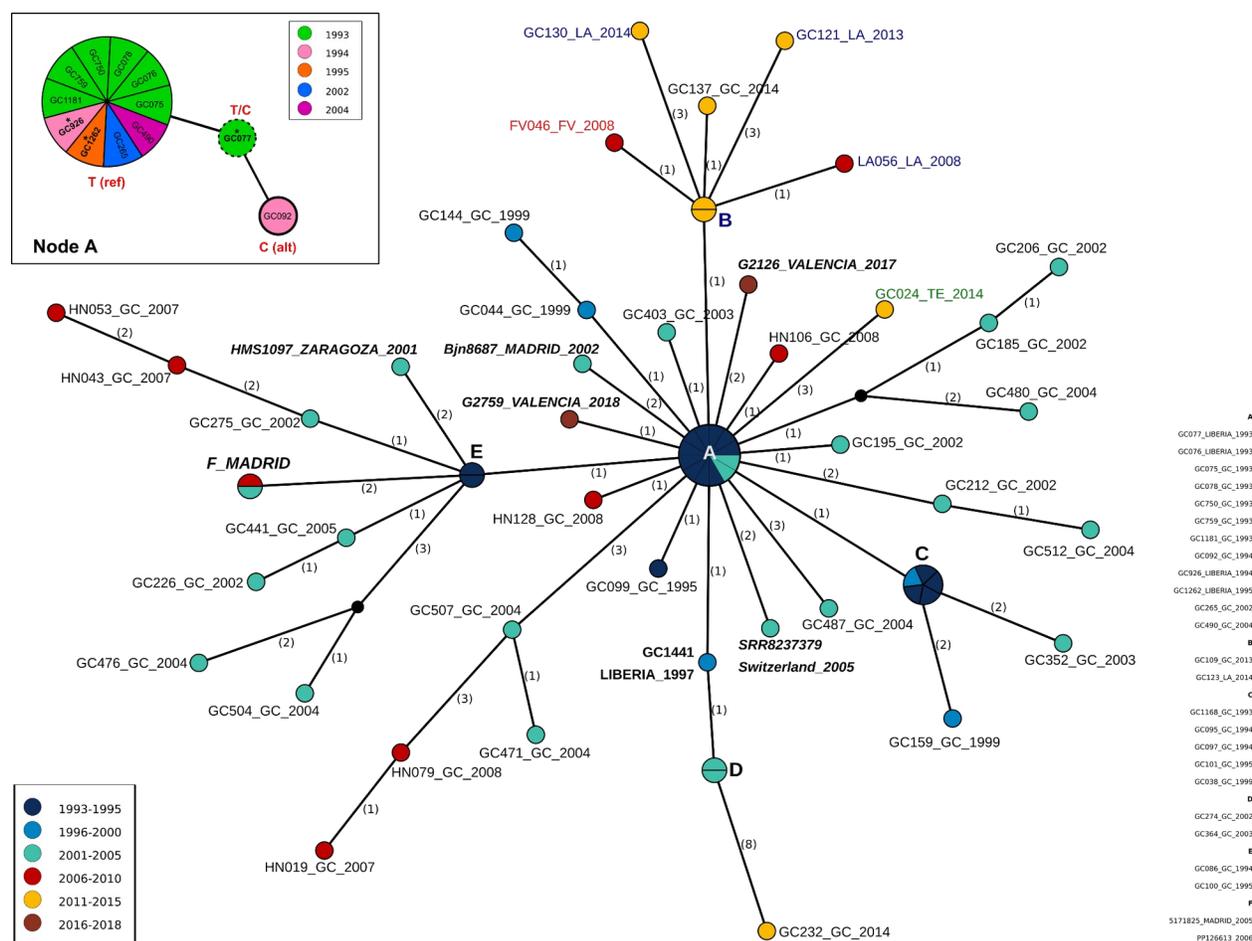
304 **Figure 2.** Genomic delineation of the outbreak. **A.** ML tree highlighting the phylogenetic  
 305 circumscription of the outbreak and related strains identified with the partial SNPs profile. **B-C.**  
 306 Density graphics of pairwise number of SNPs between samples of GC outbreak and Buenos  
 307 Aires (BA) [5]; Bern [6]; Denmark [4]; Thailand [3]. **D.** Study of the success of GC outbreak  
 308 outside the studied area, maps indicating the origin of the sequences identified by the different

309 sources and their particular SNPs profiles; this study includes all samples sequenced as  
310 considered part of the outbreak by different molecular typing methods.

### 311 *GC outbreak topology*

312 To study the epidemiological characteristics of the GC strain in the Canary Island and abroad,  
313 we constructed a median joining network with all samples, including the three additional  
314 samples from the likely index case (Table S1, Figure 3). The genomic network displays a star-  
315 like structure. The likely index case is located in the center (including the oldest sample and two  
316 additional samples collected in 1994, 1995) along with 9 other isolates; most of them from 1993,  
317 and all collected from Gran Canaria. In depth investigation of central node (node A) revealed  
318 one heterogenous SNP (hSNP), at genomic position 3190007, for which the first sample of the  
319 likely index case (GC077) harbor both the reference (T: 82.7%) allele, shared by the rest of  
320 samples of the outbreak, and the alternative allele (C: 17.3%), which is fixed in sample GC092  
321 from 1994 (Figure 3). The other hSNPs identified were uninformative since only one allele was  
322 fixed, either reference or alternative.

323 Connected with node A, three additional smaller star-like structures were observed,  
324 corresponding to nodes B, C and E (Figure 3), and resembling secondary outbreaks. Node B  
325 represents the later spreading (around 2008) from Gran Canaria to the other islands in the  
326 archipelago; with an additional independent introduction in Tenerife. Node C, the smallest  
327 secondary outbreak, also occurred in Gran Canaria. Notably, the biggest secondary outbreak,  
328 node E, started early and includes samples from the peninsula (Madrid and Zaragoza).  
329 Furthermore, the network reveals that the index case initiated a new transmission chain in 1997  
330 (isolate GC1441), also in Gran Canaria. Surprisingly, the two samples from Valencia were  
331 directly linked to the origin of the outbreak, without connection between them. Epidemiological  
332 data agree with the absence of link between both cases, but there is also no evidence of trip to  
333 Canary Islands for either patient, suggesting an additional missing case (or cases) linking them.  
334 For the three isolates collected in Madrid, there is some link with Gran Canaria; case Bjn8687  
335 stayed there in prison two years before diagnosis and the others (PP126613, 5171825) also  
336 visited the island at some point before diagnosis (Table S1). For the isolate identified from  
337 Switzerland, no epidemiological information was obtained. Overall, the extent of secondary  
338 transmissions after exportation events seems very limited, since no additional sequences from  
339 other places were retrieved from ENA.

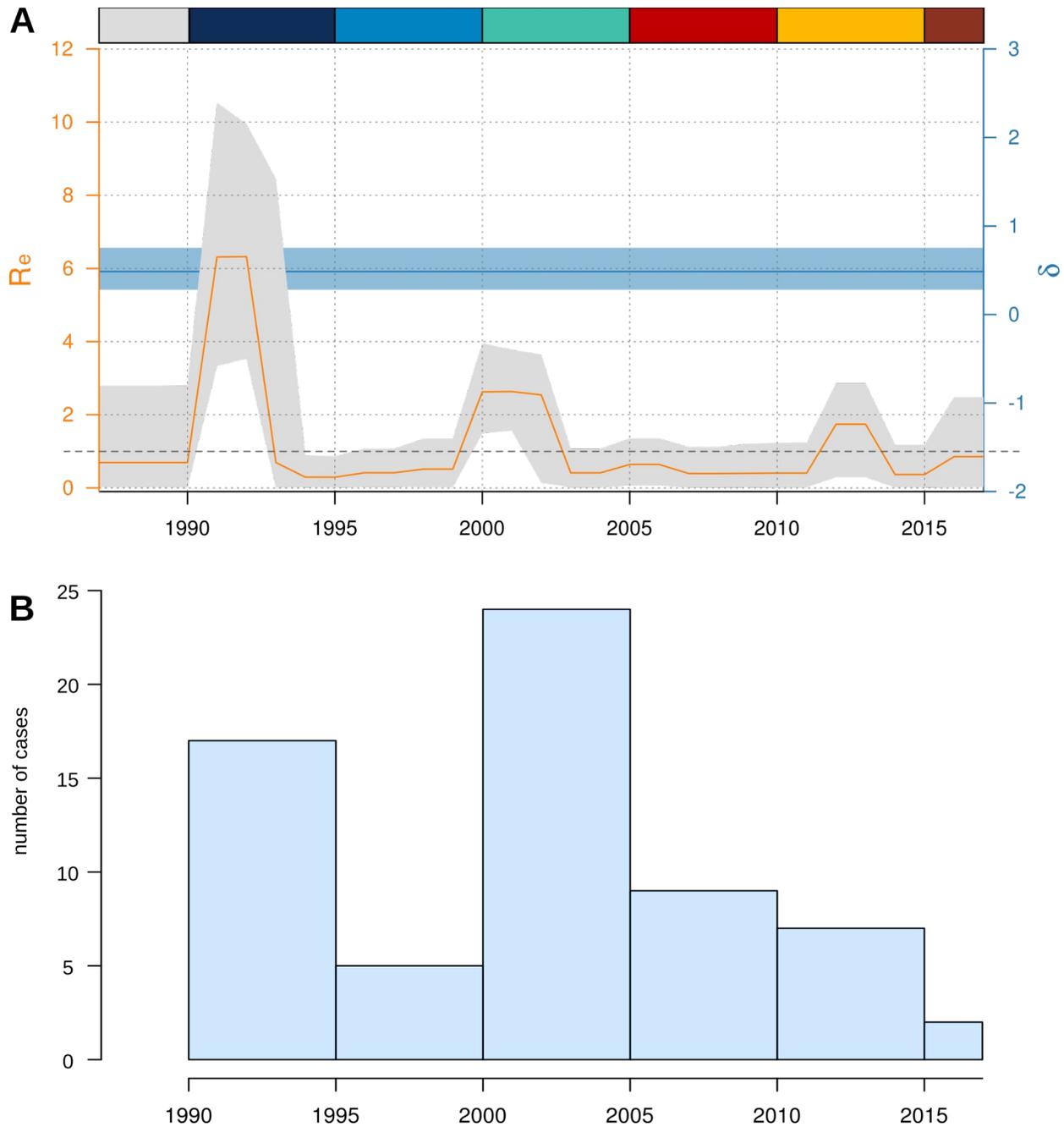


340 **Figure 3.** Median joining network analysis. Numbers in parentheses indicate the number of  
 341 SNPs between nodes. Node size indicates the number of samples with the same genome, node  
 342 color denotes sampling time; name color indicates different islands (blue: Lanzarote; red:  
 343 Fuerteventura; green: Tenerife; black: Gran Canaria). Samples from the continent are indicated  
 344 in italic and bold letters. GC1441 is an additional sample from the index case. Node A resolution  
 345 with hSNPs pos: 3910007 is detailed; reference (ref) and alternative (alt) alleles distribution  
 346 among samples is indicated. Names with asterisks indicate samples of the index case.

### 347 *Phylogenetics analysis*

348 We reasoned that if the GC strain had any transmission advantage this could be reflected in its  
 349 natural history. Since we found evidence of temporal structure in our dataset (see Methods), it is  
 350 suitable for applying a Birth-death skyline (BDSKY) model to estimate the epidemiological

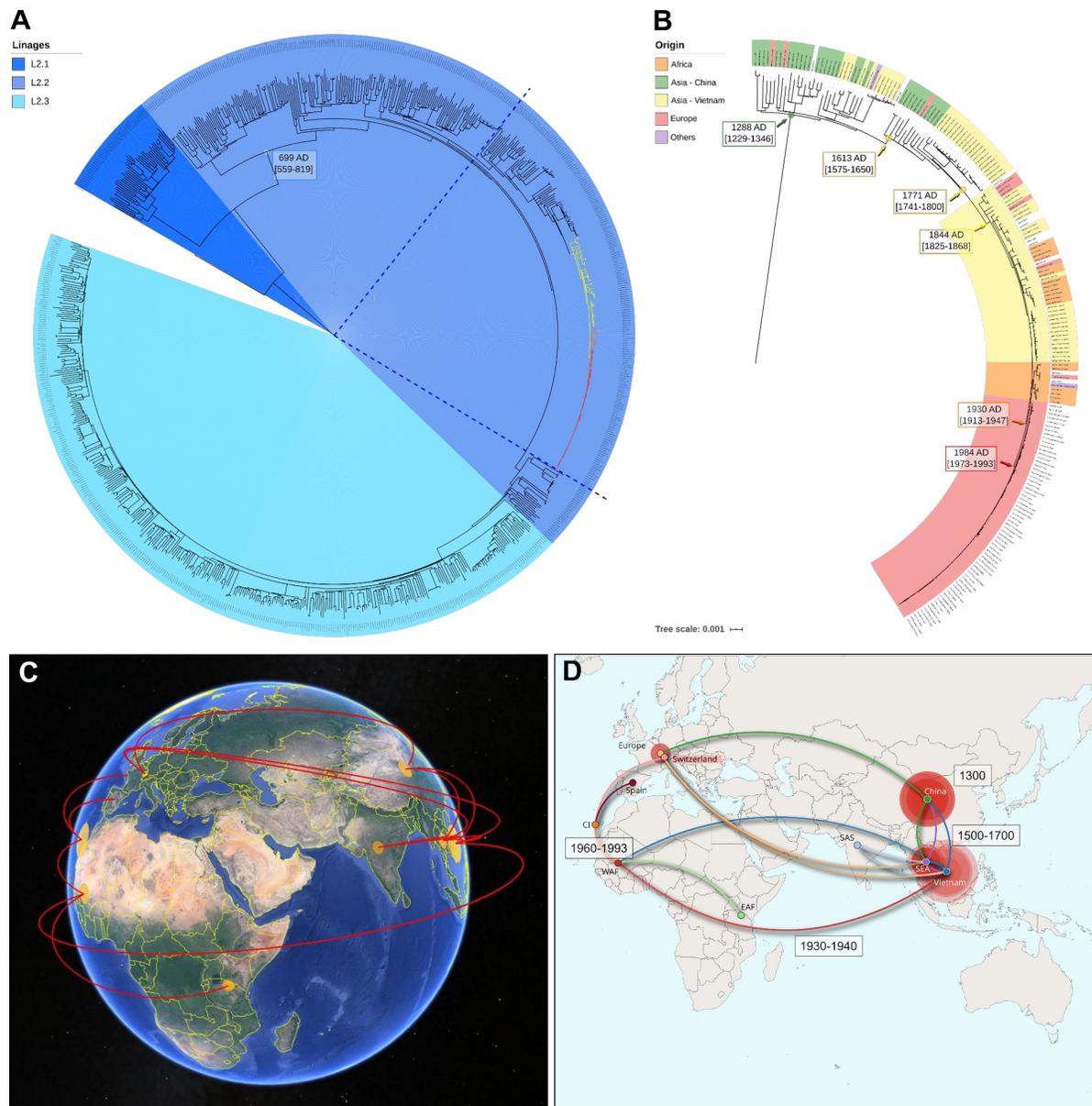
351 parameters of the GC outbreak (Supplementary Notes).  
352 The becoming uninfected or recovery rate ( $\delta$ ) resulted in a median value of 0.49 [0.28 - 0.75,  
353 95% HPD] suggesting an infection period of 2 years [1.3 - 3.6 y] in agreement with the global  
354 estimation of 1-3 years previously proposed [35]. The effective reproduction number varies  
355 through the period of the outbreak, displaying a particular profile with peaks of expansion ( $R_e > 1$ )  
356 and reduction ( $R_e < 1$ ). The peaks were observed approximately every ten years, and have an  
357 extension of three years (Figure 4A). Peaks matched with the secondary outbreak nodes  
358 observed in the network, and the periods of outbreak reduction coincide with chains instead of  
359 star-like transmission patterns (Figure 3). As observed in the histogram, peaks do not reflect the  
360 sampling effort, since periods 2005-2010 and 2010-2015 display a similar number of samples  
361 (Figure 4B). While in the first period isolates appeared in terminal nodes (Figure 3), in the  
362 second the cases are part of a secondary outbreak (Figure 3, node B). On average, the  
363 outbreak is decreasing, since the last period showed a  $R_e$  close to 1. Besides its particular  
364 profile,  $R_e$  never exceeds a value of 12 [6, 0.01 - 11 95% HPD], the maximum number of  
365 secondary cases caused by an infected person per year estimated for TB [35]. Overall, there is  
366 no indication from phylodynamics analysis that the GC strain has a transmission advantage, at  
367 least in terms of secondary cases generated.  
368 On the contrary, if ecological drivers like host or social determinants were behind the initial  
369 success of the GC strain, it will be mostly expected among the vulnerable population. In this  
370 sense, from all 61 patients belonging to the GC outbreak, 64% presented at least one risk factor  
371 including; PDU (Parenteral Drug Users), NPDU (Non-Parenteral Drug Users), imprisonment,  
372 HIV infection, alcohol abuse, indigence and no TB treatment adherence (Table S1). All of them  
373 are known risk factors for TB infection in low burden countries. Furthermore, considering only  
374 patients with available information, the value increases to 76%.



375 **Figure 4.** Phylodynamics analysis of GC outbreak. **A.** Birth-death serial skyline results showing  
376 reproductive number ( $R_e$ ) variation across outbreak period and recovery rate ( $\delta$ ). Rectangles in  
377 the top indicate periods of time as colored in Figure 3. **B.** Histogram of number of cases  
378 sequenced in the different time periods.

379 *Tracking the GC outbreak strain over centuries*

380 The GC outbreak strain belongs to the MTBC lineage 2. A global ML phylogeny was  
381 constructed including 740 L2 global strains (MSA with 43401 concatenated SNPs) from different  
382 studies (Table S4). Placed in the context of L2 diversity and considering the classification  
383 previously proposed [38,39], the GC outbreak strain belongs to L2.2.3 (Figure 5A). The closest  
384 clade is mainly composed of African strains (Figure 5B, orange clade) followed by a mixed  
385 African-Vietnamese clade (Figure 5B, yellow clade) and more distantly related to China (Figure  
386 5B). The emergence of the outbreak was estimated around 1984 AD [1973-1993 AD, 95% HPD]  
387 in agreement with the arrival time of the likely index case to the Archipelago (Figure S1). The  
388 MRCA of the African clade was estimated at 1930 AD [1913-1947 AD, 95% HPD] and the  
389 MRCA of the African-Vietnamese within the middle of 19th century (Figure 5B). Origin of L2.2  
390 was estimated between 559-819 AD in agreement with previous estimates using the same  
391 approach [40]. Phylogeographic and dispersal reconstruction of the outbreak ancestors, placed  
392 its oldest ancestor in China around the middle 1300 AD; from there it spread to Vietnam  
393 between 1500-1700 AD. From Vietnam it dispersed to Africa between 1930-1940 AD, then to  
394 Canary Islands between 1960-1993 AD to later reach Europe between 1993-2002 (Figure 6,  
395 File S1).



396 **Figure 5.** Dating, phylogeographic and dispersal analysis. **A.** ML global phylogeny of L2  
 397 highlighting the sublineages. Time of the most recent common ancestor (tMRCA) of L2.2. is  
 398 indicated. Colored branches correspond to GC-outbreak (red), African related clade (orange)  
 399 and African-Vietnamese clade (yellow). **B.** tMRCA and origin of GC outbreak, closest nodes  
 400 and oldest Chinese ancestor. Colors indicate origin of either samples (labels) and the ancestor  
 401 of different clades. **C.** Results of the phylogeographic analysis in Google earth view. **D.** Routes  
 402 of dispersion and time, lines are colored following the destination of the migration, and times  
 403 from the oldest origin in China, to Vietnam and Canary Islands are denoted.

## 404 **Discussion**

405 *Mycobacterium tuberculosis* is an obligate human pathogen and its success in the population  
406 deeply depends on human movement. By combining WGS and phylodynamics analyses, here  
407 we describe the dynamics of a 25 years outbreak in detail. After the index case arrival to Gran  
408 Canaria, the strain spread quickly, since many cases -with identical sequences- were observed  
409 in 1993. Two close secondary outbreaks appeared soon, one connected with the samples in the  
410 archipelago, and the other extended in the same island, probably as an early independent  
411 transmission chain. Lately, a differentiated secondary outbreak spread to the rest of the islands,  
412 reaching high representativeness [12]. By evaluating serial samples from different years, we  
413 could identify secondary outbreaks associated with the likely index case, probably due to poor  
414 treatment adherence. The deep analysis of hSNPs resolved the relations among samples in the  
415 central core of the outbreak, adding molecular evidence to support GC077 as the index case  
416 (first sample). A similar approach has been described in Lee et al. [22] showing the potential of  
417 using hSNP for further genetic discrimination within outbreaks. This approach also provided  
418 evidence that most of the observed diversity within the outbreak was already present in the  
419 index case samples. Furthermore, this approach allowed us to circumscribe the outbreak,  
420 discarding distantly related isolates linked by molecular methods that only query a restricted  
421 portion of the genome.

422 With the added value of querying a large collection of WGS [13,25], our comprehensive  
423 approach allowed us to link the GC strain to cases in continental Europe and thus determine the  
424 extent of an outbreak beyond the geographic limits where it was studied. Thus, while the GC  
425 outbreak strain had a very high local success, reaching a frequency of 27% of all isolates in  
426 Gran Canaria [12], it had limited expansion outside the archipelago; with multiple introductions  
427 in the continent not resulting in secondary cases [10].

428 By combining the finding of a wide SNP profile and by querying a large dataset, we were able to  
429 shed light on the remote origin of the outbreak strain. Our analysis places the ancestors of the  
430 GC strain hundreds of years ago, in China around 1300 AD. Thanks to the vast amount of  
431 MTBC genomic data available we have traced its initial movements from China to Vietnam, later  
432 to Liberia and finally to Gran Canaria by the end of the 20<sup>th</sup> century. Once in Gran Canaria the  
433 strain generated a large number of secondary cases there, but was rarely seen outside the  
434 archipelago. Thus, here we were able to document how past epidemiological events impact  
435 current TB epidemics.

436 In addition to inferring its remote origin, we had the unique opportunity to study its proximal  
437 origin using a phylodynamic approach. Proximal origins of large TB outbreaks are rarely  
438 identified by epidemiological information, mainly because of the difficulty in identifying the index  
439 patient in an area with many local cases and, in the case of tuberculosis, because of latency  
440 periods. Since the GC outbreak was undoubtedly linked to the migration of an infected person  
441 from Liberia to Gran Canaria in 1993 [9], this represents a unique opportunity to validate the  
442 commonly applied Bayesian approach for tuberculosis. Our results (1984 AD [1973-1993 AD,  
443 95% HPD]) largely agreed with the available epidemiological information.

444 The main limitations of our study is the low proportion of outbreak cases sequenced from  
445 Canary Islands, it accounts for approximately 10% of total cases, however we show that the  
446 dataset has enough phylogenetic signal for the analysis presented. In accordance, our  
447 phylodynamics results do not correlate to sampling effort (Figure 4). Spain also lacks a systemic  
448 WGS program, however in our study we use systematic typing data from three of the largest  
449 regions in Spain: Madrid, Aragón and Valencia Region to confirm that the strain had no success  
450 beyond the archipelago. Finally, by querying the whole ENA database we could retrieve  
451 additional cases sequenced elsewhere and could evaluate the expansion of the outbreak.

452 A common theme when it comes to MTBC strains is whether some are more transmissible than  
453 others or if its success is driven by ecological factors, or both. Outbreak strains are good  
454 candidates to look for intrinsic transmissibility as they are responsible for a large number of  
455 cases with sustained transmission over decades. For the GC strain, we estimate an infection  
456 period close to two years and less than 12 secondary cases per infected individual, same  
457 values proposed for TB in general [35], suggesting that there is no transmissibility advantage of  
458 the GC strain, at least in terms of shorter latency periods or higher contagion rate. In addition, it  
459 was demonstrated that GC strain did not display higher virulence, nor accumulated mutations in  
460 sequential samples or acquired resistance mutations [11]. Similarly, our extensive query of  
461 outbreak sequences in public repositories did not identify that this strain is responsible for  
462 outbreaks outside the Canary Islands. For example, an in-depth genotyping effort did not  
463 identify secondary cases associated with a case in Madrid that had prolonged disease [11] or in  
464 Aragón where systematic genotyping is applied to all TB positive cases. In addition, mutations  
465 associated with the GC strain are unlikely to have a role in transmission. For other outbreak  
466 strains, notably the Danish C2 [41] and the Toronto [42] strains, mutations with a functional role  
467 have been proposed but no experimental evidence is available. On the contrary, all those  
468 outbreak strains share the fact that they thrive in populations with significant risk factors,

469 particularly during the '80s and '90s of the past century. Overall, data suggests that the success  
470 of the long-lived GC outbreak strain, and probably others, is related to ecological factors  
471 associated with founder effects linked to the host and social determinants of the disease rather  
472 than an intrinsic transmissibility difference.

### 473 **Acknowledgements**

474 This project has been funded by the European Research Council (101001038-TB-  
475 RECONNECT), the Ministerio de Economía, Industria y Competitividad (PID2019-104477RB-  
476 I00) and European Commission –NextGenerationEU (Regulation EU 2020/2094), through  
477 CSIC's Global Health Platform (PTI Salud Global) to IC. This project has been funded by the  
478 Instituto de Salud Carlos III (FIS18/0336), Gobierno de Aragón/Fondo Social Europeo  
479 "Construyendo Europa desde Aragón" to SS.

### 480 **Conflicts of interest**

481 IC received consultancy fees from Foundation for innovative new diagnostics. The author has  
482 no other competing interests to declare. The remaining authors declare no competing interests.

483 **Contributors.** M.G.L and I.C conceived the study; M.I.C.H, F.C, R.C, L.S, B.P, L.P.L, D.G.V,  
484 S.S, J.C., E.L. collected the samples, obtained and curated the epidemiological data; M.T.P.  
485 processed and sequenced the samples; P.W., Z.I performed BIGSI analysis; M.M.M. performed  
486 the scripts for DST analysis; M.G.L performed all analyses; M.G.L and I.C wrote the first draft of  
487 the manuscript, with revisions from all co-authors.

### 488 **References**

- 489 [1] Global tuberculosis report 2020 n.d. <https://www.who.int/publications/i/item/9789240013131>  
490 (accessed April 7, 2021).
- 491 [2] Glaziou P. Predicted impact of the COVID-19 pandemic on global tuberculosis deaths in  
492 2020. medRxiv 2020:2020.04.28.20079582.
- 493 [3] Coscolla M, Barry PM, Oeltmann JE, Koshinsky H, Shaw T, Cilnis M, et al. Genomic  
494 epidemiology of multidrug-resistant Mycobacterium tuberculosis during transcontinental  
495 spread. J Infect Dis 2015;212:302–10.
- 496 [4] Folkvardsen DB, Norman A, Andersen ÅB, Michael Rasmussen E, Jelsbak L, Lillebaek T.  
497 Genomic Epidemiology of a Major Mycobacterium tuberculosis Outbreak: Retrospective

- 498 Cohort Study in a Low-Incidence Setting Using Sparse Time-Series Sampling. *J Infect Dis*  
499 2017;216:366–74.
- 500 [5] Eldholm V, Monteserin J, Rieux A, Lopez B, Sobkowiak B, Ritacco V, et al. Four decades of  
501 transmission of a multidrug-resistant *Mycobacterium tuberculosis* outbreak strain. *Nat*  
502 *Commun* 2015;6:7119.
- 503 [6] Stucki D, Ballif M, Bodmer T, Coscolla M, Maurer A-M, Droz S, et al. Tracking a  
504 tuberculosis outbreak over 21 years: strain-specific single-nucleotide polymorphism typing  
505 combined with targeted whole-genome sequencing. *J Infect Dis* 2015;211:1306–16.
- 506 [7] Casali N, Broda A, Harris SR, Parkhill J, Brown T, Drobniowski F. Whole Genome  
507 Sequence Analysis of a Large Isoniazid-Resistant Tuberculosis Outbreak in London: A  
508 Retrospective Observational Study. *PLoS Med* 2016;13:e1002137.
- 509 [8] Roetzer A, Diel R, Kohl TA, Rückert C, Nübel U, Blom J, et al. Whole genome sequencing  
510 versus traditional genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a  
511 longitudinal molecular epidemiological study. *PLoS Med* 2013;10:e1001387.
- 512 [9] Caminero JA, Pena MJ, Campos-Herrero MI, Rodríguez JC, García I, Cabrera P, et al.  
513 Epidemiological evidence of the spread of a *Mycobacterium tuberculosis* strain of the  
514 Beijing genotype on Gran Canaria Island. *Am J Respir Crit Care Med* 2001;164:1165–70.
- 515 [10] Pérez-Lago L, Herranz M, Comas I, Ruiz-Serrano MJ, Roa PL, Bouza E, et al. Ultrafast  
516 Assessment of the Presence of a High-Risk *Mycobacterium tuberculosis* Strain in a  
517 Population. *J Clin Microbiol* 2016;54:779–81.
- 518 [11] Pérez-Lago L, Navarro Y, Montilla P, Comas I, Herranz M, Rodríguez-Gallego C, et al.  
519 Persistent Infection by a *Mycobacterium tuberculosis* Strain That Was Theorized To Have  
520 Advantageous Properties, as It Was Responsible for a Massive Outbreak. *J Clin Microbiol*  
521 2015;53:3423–9.
- 522 [12] Pérez-Lago L, Campos-Herrero MI, Cañas F, Copado R, Sante L, Pino B, et al. A  
523 *Mycobacterium tuberculosis* Beijing strain persists at high rates and extends its geographic  
524 boundaries 20 years after importation. *Sci Rep* 2019;9:1–6.
- 525 [13] Chiner-Oms Á, Sánchez-Busó L, Corander J, Gagneux S, Harris SR, Young D, et al.  
526 Genomic determinants of speciation and spread of the complex. *Sci Adv* 2019;5:eaaw3307.
- 527 [14] Cancino-Muñoz I, López MG, Torres-Puente M, Villamayor LM, Borrás R, Borrás-Máñez M,  
528 et al. Population-based sequencing of *Mycobacterium tuberculosis* reveals how current  
529 population dynamics are shaped by past epidemics. *medRxiv* 2022:2022.01.24.22269736.
- 530 [15] López MG, Dogba JB, Torres-Puente M, Goig GA, Moreno-Molina M, Villamayor LM, et al.  
531 Tuberculosis in Liberia: high multidrug-resistance burden, transmission and diversity  
532 modelled by multiple importation events. *Microb Genom* 2020;6.  
533 <https://doi.org/10.1099/mgen.0.000325>.
- 534 [16] Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor.  
535 *Bioinformatics* 2018;34:i884–90.
- 536 [17] Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact  
537 alignments. *Genome Biology* 2014;15:R46. <https://doi.org/10.1186/gb-2014-15-3-r46>.
- 538 [18] Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
539 *Bioinformatics* 2009;25:1754–60.
- 540 [19] Li H. A statistical framework for SNP calling, mutation discovery, association mapping and  
541 population genetical parameter estimation from sequencing data. *Bioinformatics*  
542 2011;27:2987–93.
- 543 [20] Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2: somatic  
544 mutation and copy number alteration discovery in cancer by exome sequencing. *Genome*  
545 *Res* 2012;22:568–76.

- 546 [21] McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The  
547 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA  
548 sequencing data. *Genome Res* 2010;20:1297–303.
- 549 [22] Lee RS, Proulx J-F, McIntosh F, Behr MA, Hanage WP. Previously undetected super-  
550 spreading of revealed by deep sequencing. *Elife* 2020;9.  
551 <https://doi.org/10.7554/eLife.53245>.
- 552 [23] Coll F, McNerney R, Guerra-Assunção JA, Glynn JR, Perdigão J, Viveiros M, et al. A robust  
553 SNP barcode for typing *Mycobacterium tuberculosis* complex strains. *Nat Commun*  
554 2014;5:4812.
- 555 [24] Stucki D, Brites D, Jeljeli L, Coscolla M, Liu Q, Trauner A, et al. *Mycobacterium*  
556 *tuberculosis* lineage 4 comprises globally distributed and geographically restricted  
557 sublineages. *Nat Genet* 2016;48:1535–43.
- 558 [25] Bradley P, den Bakker HC, Rocha EPC, McVean G, Iqbal Z. Ultrafast search of all  
559 deposited bacterial and viral genomic data. *Nat Biotechnol* 2019;37:152–9.
- 560 [26] Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
561 phylogenies. *Bioinformatics* 2014;30:1312–3.
- 562 [27] Leigh JW, Bryant D. popart : full-feature software for haplotype network construction.  
563 *Methods in Ecology and Evolution* 2015;6:1110–6. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210x.12410)  
564 [210x.12410](https://doi.org/10.1111/2041-210x.12410).
- 565 [28] Menardo F, Loiseau C, Brites D, Coscolla M, Gygli SM, Rutaihwa LK, et al. Treemmer: a  
566 tool to reduce large phylogenetic datasets with minimal loss of diversity. *BMC*  
567 *Bioinformatics* 2018;19. <https://doi.org/10.1186/s12859-018-2164-8>.
- 568 [29] Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, et al. BEAST 2: A Software  
569 Platform for Bayesian Evolutionary Analysis. *PLoS Computational Biology*  
570 2014;10:e1003537. <https://doi.org/10.1371/journal.pcbi.1003537>.
- 571 [30] Bos KI, Harkins KM, Herbig A, Coscolla M, Weber N, Comas I, et al. Pre-Columbian  
572 mycobacterial genomes reveal seals as a source of New World human tuberculosis. *Nature*  
573 2014;514:494–7.
- 574 [31] Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior Summarization in  
575 Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol* 2018;67:901–4.
- 576 [32] Lemey P, Rambaut A, Drummond AJ, Suchard MA. Bayesian Phylogeography Finds Its  
577 Roots. *PLoS Comput Biol* 2009;5:e1000520.
- 578 [33] Bielejec F, Baele G, Vrancken B, Suchard MA, Rambaut A, Lemey P. SpreaD3: Interactive  
579 Visualization of Spatiotemporal History and Trait Evolutionary Processes. *Mol Biol Evol*  
580 2016;33. <https://doi.org/10.1093/molbev/msw082>.
- 581 [34] Stadler T, Kühnert D, Bonhoeffer S, Drummond AJ. Birth-death skyline plot reveals  
582 temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc Natl Acad*  
583 *Sci U S A* 2013;110:228–33.
- 584 [35] Dowdy DW, Dye C, Cohen T. Data needs for evidence-based decisions: a tuberculosis  
585 modeler’s “wish list” [Review article] 2013. <https://doi.org/10.5588/ijtld.12.0573>.
- 586 [36] Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of  
587 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol*  
588 2016;2:vev007.
- 589 [37] Menardo F, Duchêne S, Brites D, Gagneux S. The molecular clock of *Mycobacterium*  
590 *tuberculosis*. *PLoS Pathog* 2019;15:e1008067.
- 591 [38] Shitikov E, Kolchenko S, Mokrousov I, Bespyatykh J, Ischenko D, Ilina E, et al. Evolutionary  
592 pathway analysis and unified classification of East Asian lineage of *Mycobacterium*  
593 *tuberculosis*. *Sci Rep* 2017;7:1–10.

- 594 [39] Rutaihwa LK, Menardo F, Stucki D, Gygli SM, Ley SD, Malla B, et al. Multiple Introductions  
595 of Mycobacterium tuberculosis Lineage 2–Beijing Into Africa Over Centuries 2019.  
596 <https://doi.org/10.3389/fevo.2019.00112>.
- 597 [40] Liu Q, Ma A, Wei L, Pang Y, Wu B, Luo T, et al. China’s tuberculosis epidemic stems from  
598 historical expansion of four strains of Mycobacterium tuberculosis. *Nature Ecology &*  
599 *Evolution* 2018;2:1982–92. <https://doi.org/10.1038/s41559-018-0680-6>.
- 600 [41] Folkvardsen DB, Norman A, Andersen ÅB, Rasmussen EM, Lillebaek T, Jelsbak L. A Major  
601 Mycobacterium tuberculosis outbreak caused by one specific genotype in a low-incidence  
602 country: Exploring gene profile virulence explanations. *Sci Rep* 2018;8:1–8.
- 603 [42] Mehaffy C, Guthrie JL, Alexander DC, Stuart R, Rea E, Jamieson FB. Marked  
604 microevolution of a unique Mycobacterium tuberculosis strain in 17 years of ongoing  
605 transmission in a high risk population. *PLoS One* 2014;9:e112928.