

1 **Genomic determinants of sympatric speciation of the *Mycobacterium***
2 ***tuberculosis* complex across evolutionary timescales**

3

4 Álvaro Chiner-Oms^{1, #}, Leonor Sánchez-Busó^{2, #}, Jukka Corander^{2, 3, 4}, Sebastien
5 Gagneux^{5, 6}, Simon Harris², Douglas Young⁷, Fernando González-Candelas^{1, 8}
6 Iñaki Comas^{8, 9*}

7

8 ¹Unidad Mixta “Infección y Salud Pública” FISABIO-CSISP/Universidad de
9 Valencia, Instituto de Biología Integrativa de Sistemas-I2SysBio, Valencia,
10 Spain.

11 ²Pathogen Genomics, Wellcome Trust Sanger Institute, Cambridge CB10 1SA,
12 UK.

13 ³Department of Biostatistics, University of Oslo, 0317 Oslo, Norway.

14 ⁴Helsinki Institute of Information Technology (HIIT), Department of Mathematics
15 and Statistics, University of Helsinki, 00014 Helsinki, Finland.

16 ⁵Swiss Tropical and Public Health Institute, Basel, Switzerland.

17 ⁶University of Basel, Basel, Switzerland.

18 ⁷The Francis Crick Institute, 1 Midland Road, London NW1 1AT, UK.

19 ⁸CIBER en Epidemiología y Salud Pública, Valencia, Spain

20 ⁹Instituto de Biomedicina de Valencia, IBV-CSIC, Valencia, Spain.

21

22 # Contributed equally.

23 *Corresponding author.

24

25 Word-count: 4,016

26

27 **ABSTRACT**

28

29 Models on how bacterial lineages differentiate increase our understanding on
30 early bacterial speciation events and about the relevant niche-specific loci
31 involved. In the light of those models we analyse the population genomics events
32 leading to the emergence of the *Mycobacterium tuberculosis* complex (MTBC)
33 from related mycobacteria species. Emergence is characterised by a combination
34 of recombination events involving multiple core pathogenesis functions and
35 purifying selection on early diverging loci. After the separation from closely related
36 mycobacteria we identify the *phoR* gene, a transcriptional regulator involved in
37 multiple aspects of MTBC virulence, as a key functional player subject to
38 pervasive positive selection. First, during the early diversification of the MTBC,
39 PhoR played a central role defining the host range of the various MTBC
40 members. Later, following adaption to the human host, PhoR mediates host-
41 pathogen interaction during human-to-human transmission. We thus show that
42 linking pathogen evolution across evolutionary and epidemiological timescales
43 lead to the identification of past and present virulence determinants of interest for
44 biomedical research.

45 The increasing availability of population genomics data has allowed an improved
46 understanding of genotypic and ecological differentiation among closely related
47 bacteria¹. Models emanating from these data predict that sympatric speciation
48 among bacteria occupying the same ecological niche is more likely leaving many
49 times measurable genetic signatures in extant genomes^{2,3}. However, is still
50 underexplored how these models apply to professional pathogens, particularly
51 those characterized by an obligate association with their host species.

52 Species of the *Mycobacterium tuberculosis* complex (MTBC) cause devastating
53 morbidity and mortality in humans and animals, which also leads to important
54 economic losses⁴. The MTBC comprises a group of bacteria with genome
55 sequences having an average nucleotide identity of greater than 99% and sharing
56 a single clonal ancestor. This includes the predominantly human pathogens
57 referred to as *Mycobacterium tuberculosis* and *Mycobacterium africanum* as well
58 as a series of pathogens isolated from other mammalian species known as *M.*
59 *bovis*, *M. pinnipedi*, *M. antelope*, *M. microti*, etc. Human-adapted tuberculosis
60 bacilli show a strong geographic association, with some lineages and even
61 sublineages being globally distributed (e.g. lineage 4) and others geographically
62 restricted (e.g. lineage 5, 6, 7)^{5,6}. It is assumed that the causes of this variable
63 geographical distribution are both historical (e.g. trade, conquest, globalization)
64 and biological (e.g. interactions with different human genetic backgrounds)⁶.
65 There is limited transmissibility of animal-adapted strains in humans and,
66 conversely, human-adapted strains transmit poorly among animals⁷. Despite the
67 wide range of host species infected by the different members of the MTBC, there
68 is a maximum of ~2,500 single nucleotide polymorphisms (SNPs) separating any
69 two MTBC genomes⁸. The most closely-related bacteria that fall outside of the
70 MTBC include isolates initially referred to as “smooth tubercle bacilli” and now
71 referred to as *Mycobacterium canettii* (MCAN). MCAN strains differ from MTBC
72 isolates by tens of thousands of SNPs⁹. MCAN strains have been isolated from
73 the Horn of Africa, predominantly from children and often in association with
74 extrapulmonary tuberculosis¹⁰. Genomic comparisons have identified gene
75 content differences between MTBC, MCAN and other mycobacteria^{9,11} as well as
76 genetic differences in virulence-related loci¹².

77
78 Our understanding about the population genetics events mediating the
79 divergence of the ancestor of the MTBC from an MCAN-like ancestral pool is far
80 from complete. The availability of genome sequences from thousands of MTBC
81 clinical strains as well as of closely relatives like MCAN enable us not only to
82 identify molecular signatures of MTBC speciation events, but also to reveal
83 known and new targets for biomedical research.

84

85 RESULTS

86 To explore this, we first analyzed the differentiation between MTBC and MCAN
87 by searching for any hallmark of on-going recombination between and within
88 these groups of strains. Previous reports have suggested that there might be
89 limited but significant recombination among MTBC strains¹³. To maximize the
90 chances of identifying potential ongoing recombination events within the MTBC,
91 we screened a data set of 3,475 complete genome sequences of strains from
92 global sources¹⁴ and kept those that maximized the within-MTBC diversity (n =

93 1,591). These genomes are representative of the known geographic and genetic
94 diversity of the MTBC (Supplementary Fig. 1). Among those genomes, we
95 identified all the biallelic variant positions that were not related with known drug
96 resistance genes and called them in all the strains. In addition, we identified
97 potential false positive variant calls due to mapping errors. We detected 239
98 positions possibly involved in mapping errors due to duplicated sequences,
99 leaving a final set of 94,780 core variant positions in the MTBC.

100 From the 94,780 variant positions, we identified potential homoplastic sites, i.e.
101 polymorphic sites showing signs of convergent evolution, by parsimony mapping.
102 A total of 2,360 homoplastic sites were identified (2.5% of all variable sites).
103 Convergence can arise as a result of recombination, selection or neutral
104 processes. With these variant positions, we looked at consecutive runs of
105 homoplastic sites in the genomes. The 2,360 homoplastic positions defined 97
106 homoplastic runs (see Methods for definition of a homoplastic run). If the
107 accumulation of homoplasies was due to recombination, we expect the variant
108 positions involved in each consecutive run to share the same phylogenetic
109 mapping. From the 97 homoplastic runs, we detected only 2 cases in which two
110 variant positions shared phylogenetic congruence. The two regions accounted
111 for 4 convergent variants (Supplementary Table 1) and affected strains from
112 different MTBC lineages. Variants in positions 2195896 and 2195899 fell in the
113 primary regulatory region of *mazE5*¹⁵. On the other hand, variants in positions
114 2,641,161 and 2,641,163 fell in the intergenic region of *glyS* and Rv2358.
115 Although we cannot discard possible recombination events, it is more likely that
116 the two regions have been under positive selection, a mechanism known to lead
117 to homoplastic variants accumulation in the MTBC¹⁶. In summary, this large-scale
118 variant-by-variant analysis failed to identify on-going recombination between any
119 of the 1,591 MTBC strains analyzed.

120
121 As an additional method to identify possible on-going recombination, we
122 evaluated linkage disequilibrium (LD) as a function of the distance between the
123 94,780 core variant positions. R^2 values were slightly higher at shorter distances,
124 but high values of this parameter can arise when comparing variants with very
125 different frequencies¹⁷ as is revealed by the skewing of the site frequency spectra
126 of the MTBC towards low frequency values (Fig. 1B). We also calculated D' . In
127 this data set, as expected for a mostly clonal organism, LD measured by D'
128 remained at its maximum value, even when focusing at distant variant positions
129 more than 5 Kb apart, suggesting very little or no ongoing recombination (Fig.
130 1A).

131 To further validate these findings, we ran Gubbins with the same data set.
132 Gubbins detects the accumulation of a higher than expected number of variants
133 as a hallmark of possible recombination. We partitioned the 1,591 strain dataset
134 into the different lineages and screened for possible tracks of recombination. This
135 lineage-by-lineage analysis identified no potential recombination regions within
136 lineages. Therefore, the three approaches used agreed in assigning a negligible
137 role to recombination in the ongoing evolution of the MTBC.

138 Having established that recombination is not currently acting within the MTBC,
139 we compared a representative data set of MTBC genomes (Comas 2013, n =

140 219) with 7 MCAN genomes to identify and quantify ongoing recombination within
141 MCAN and between MCAN and the MTBC. From the 93,922 polymorphic sites
142 identified, 22,718 were biallelic homoplasies (24.2%). The genomic distribution
143 of variant positions and homoplasies was traced for both groups, showing very
144 different landscapes (Fig. 2). A total of 22,464 (98.9%) of those homoplasies were
145 only found among MCAN strains, representing almost half of the variability within
146 this group (22,464/52,392 biallelic sites, 42.9%) which points to recombination as
147 a main source of variability in MCAN. This fact is consistent with previous
148 reports⁹. By contrast and as expected, a flat homoplastic profile was found for the
149 219 MTBC strains, with a low occurrence of homoplasies (488/30,056 biallelic
150 sites, 1.6%).

151 To test for ongoing recombination between MCAN strains and MTBC strains, we
152 identified runs of homoplasies involving both groups. From the total 93,922
153 variants when we put together MCAN and MTBC strains we found 522 variable
154 sites (0.05% of the total) that were polymorphic in both groups, including 254
155 biallelic homoplasies in the MTBC group. Phylogenetic mapping showed that
156 these homoplasies mapped to the ancestral branch leading to the MTBC clade
157 and did not involve any extant MTBC strain. These results indicate that
158 recombination events were common between MCAN and the ancestor of the
159 MTBC, but were absent during subsequent diversification of the MTBC.
160 Consistently, Gubbins identified 990 potential recombinant segments, most of
161 which mapped on branches involving only MCAN strains ($n = 907$), thus
162 corroborating that most homoplastic variants only involved MCAN strains. The
163 remaining events mapped to the common branch of all the MTBC strains
164 (Supplementary Fig. 2).

165 **Sympatric and stepwise emergence of the MTBC ancestor**

166 Our results show that recombination with closely related mycobacteria played a
167 role in the emergence of the common ancestor of the MTBC. To gain a better
168 insight into the distribution of genome variability in MCAN and the ancestor of the
169 MTBC, we extracted all the variant positions that were homoplastic between the
170 MTBC ancestor and any of the MCAN strains (7,700 positions). The MTBC
171 ancestor genome showed a similar homoplasia profile to that of the MCAN strains
172 (Fig. 2), suggesting that the ancestor of the MTBC speciated in sympatry with
173 MCAN ancestral strains. Notably, both MCAN and the MTBC ancestor shared a
174 peak around the CRISPR region, highlighting the dynamic nature of this region
175 possibly as a result of common phage infections.

176 A Gubbins analysis including MCAN genomes and the most likely common
177 ancestor of the MTBC revealed a total of 65 recombination events mapping to
178 the branch leading to the MTBC (Supplementary Table 2). To explore whether
179 these fragments reflected real recombination, a phylogeny was constructed with
180 each of them. A comparison with the topology of the non-recombinant alignment
181 (whole genome alignment subtracting the recombinant regions) using those
182 recombinant regions with enough phylogenetic signal (Supplementary Fig. 3)
183 revealed significant incongruence (SH test; p -value <0.05 , Supplementary Fig. 4,

184 Supplementary Table 3). Thus, both Gubbins and phylogenetic approaches
185 indicated that these 65 regions are likely recombinant regions.

186 To test whether speciation of the MTBC ancestor occurred in one single episode
187 or in multiple episodes over time, we analyzed the relative age of divergence of
188 the recombination fragments from the MCAN closest clade using BEAST. Given
189 the uncertainties about the timescale of MTBC evolution, the substitution rate of
190 the non-recombinant fragment was estimated normalizing to an arbitrary age of
191 the time since the Most Recent Common Ancestor (tMRCA) of the MTBC of 1
192 (see Methods for details). We then used the inferred substitution rate to estimate
193 the relative age of each recombinant fragment by using only variants that
194 accumulated in the branch of the MTBC after the recombination event. Results
195 show that the MTBC ancestor differentiated from MCAN sequentially (Fig. 3A).
196 The estimated ages show large HPD intervals as expected from the low number
197 of variant positions per fragment. Although the distribution of tMRCA for the
198 fragments represents a continuum, we can still observe one peak marking
199 “recent” events, just before the whole pathogen population became clonal, and
200 another for “ancient” events, closer to the time of divergence from the MCAN
201 group (Fig 3B).

202 If recombination played a major role in shaping the MTBC ancestral genome with
203 regards to pathogenesis, we would expect some functions related to the
204 interaction with the host to be affected. Indeed, we observed an enrichment in
205 experimentally confirmed essential genes in the regions involved in the
206 recombination, suggesting that recombination targeted important cell functions
207 (Chi-square test; p -value < 0.01). An enrichment analysis of Gene Ontology terms
208 for the genes contained in these regions identified functions related with growth,
209 and most specifically with growth involved in symbiotic interactions inside a host
210 cell as significantly overrepresented (Binomial test; adj. p -value $< 0,05$) (Fig 3C).
211 Remarkably most of the genes involved have been implicated in cirulence using
212 animal models of infection (see Discussion, Supplementary Table 4).

213 A sympatric model of speciation predicts that some parts of the genome will be
214 involved in adaptation to a new niche. The hallmark would be the accumulation
215 of variants differentiating the emerging species, at the genome-wide level or in a
216 few loci, as a consequence of reduced recombination (i.e. incipient speciation of
217 the region). We identified all of the variants that mapped to the MTBC ancestral
218 branch and that had a different nucleotide in all the MCAN strains, the so called
219 divergent variants (divSNPs). The landscape of divergent variants ($n = 5,688$, Fig.
220 4) revealed that a total of 120 genes harbored more divergent variants than
221 expected by chance ($pFDR \leq 0.01$).

222 However, bacterial genomes are highly dynamic and different processes can
223 contribute to the genetic make-up of extant species. Consequently, not all the
224 detected regions necessarily result from pure divergence by accumulation of
225 substitutions. Our phylogenetic analysis identified several genes in which divSNP
226 were introduced by horizontal gene transfer ($n = 12$) or by recombination to a
227 MCAN not present in our dataset ($n = 54$).

228 A total of 53 genes in the MTBC ancestral genome were highly divergent with
229 respect to MCAN due to substitution events (Supplementary Table 5). While the
230 genome-wide identified divSNPs might result from genetic drift or hitchhiking
231 events associated with selection on other loci, the accumulation in only 53 genes
232 suggests that those regions might have played an important role during the
233 process of niche differentiation. In agreement, these 53 genes were significantly
234 more conserved than the rest of the genome ($dN/dS = 0.154$ vs genome average
235 $dN/dS = 0.279$, chi-squared p-value = 0.000). This result suggests that, despite
236 the increased divergence from the MCAN strains, those 53 regions have been
237 evolving under purifying selection. Alternatively, the accumulation of divergent
238 variants could also represent hotspot regions for mutation. None of the genes
239 showed a similar pattern of mutation accumulation in other MCAN (no overlap
240 between the divSNPs probabilities distributions for these 53 genes and the rest
241 of the genomes, t-test p-value < 0.05).

242 **Regions under positive selection after the transition to obligate pathogen**

243 Having established that some divSNPs accumulate in genes under purifying
244 selection, we screened for positive selection patterns to identify additional genes
245 important in the transition from a newly emerged pathogen to a globally
246 established pathogen. We first revisited the evolution of antigenic proteins. Those
247 regions are recognized by the immune system and most of them are
248 hyperconserved within the MTBC^{18,19}. Interestingly, and in agreement with
249 previous data from MCAN genomic analyses⁹, the dN/dS calculated in the branch
250 of the ancestor showed a very similar pattern, with essential genes being more
251 conserved than non-essential ones and T-cell epitopes being hyperconserved
252 (Supplementary Fig. 5). Only nine divSNPs (5 synonymous and 4 non-
253 synonymous) were found in T-cell epitopes regions, which is significantly less
254 than expected by chance (Poisson distribution, p-value < 0.001).

255 Thus, antigenic regions do not show an altered pattern or intensity of selective
256 pressure as one might expect after a speciation event. We then explored what
257 other regions of the genome changed significantly in selective pressure by
258 comparing the MTBC ancestor dN/dS and the actual dN/dS in extant populations.
259 To look for robust dN/dS measures, we only took into account those genes with
260 more than 3 synonymous variant positions and at least 1 non-synonymous variant
261 position for each of the two sets. Due to the low number of divSNPs in individual
262 genes, only 121 genes were evaluated. Consequently, although additional genes
263 to those shown in the ensuing analyses may have changed the selection pattern
264 or intensity, they cannot be evaluated properly (Supplementary Table 6). We
265 were particularly interested in those genes with a drastic change from purifying
266 ($dN/dS < 1$) to diversifying or positive selection ($dN/dS > 1$) or vice versa.

267 Most of the genes evaluated did not show any sign of changing selective pressure
268 or pattern. However, when looking at the dN/dS variation data, three genes

269 appeared as outliers (as defined by Tukey²⁰). Genes changing to evolve under
270 positive selection after divergence from the MTBC ancestor were Rv2464c, a
271 DNA glycosylase involved in DNA repair, and Rv0758, also known as *phoR*.
272 Conversely, Rv0202c was under a stronger negative selective pressure following
273 speciation. Notably, PhoR forms part of the PhoP/PhoR virulence regulation
274 system²¹. In the branch leading to the MTBC ancestor, this gene was as
275 conserved at the amino acid level, as other essential genes (Chi-square test; p-
276 value 0.4721), but when we looked within the extant MTBC diversity, the gene
277 was significantly less conserved at the amino acid level than essential genes
278 (Chi-square test; p-value < 0.001).

279 **Positive selection on *phoR* linked to on-going selective pressures**

280 Given the known central role of PhoPR in MTBC virulence, we focused our
281 attention on the new mutations found in PhoR by expanding our MTBC dataset
282 to 4,593 human and five animal genomes. Using this expanded data set, we
283 observed a total of 193 nonsynonymous mutations and 31 synonymous
284 mutations in *phoR* (Fig. 6A). The average dN/dS for this gene was well above 1
285 (dN/dS = 2.37), suggesting the action of positive selection. Codon-based tests of
286 positive selection for *phoR* identified a higher dN/dS than expected by chance
287 and at least two codons with strong evidence to be under positive selection
288 (Supplementary Table 7). Additional evidence for the action of positive selection
289 on this gene derives from the nonsynonymous mutations, among which we found
290 34 homoplastic variants, which are strong predictors of positive selection in
291 MTBC (Supplementary Table 8). Non-synonymous mutations significantly
292 accumulated in the sensor domain (chi-square, p-value < 0.01), further
293 corroborating that they are involved in the fine tuning of the PhoR sensitive
294 function to the changing environment during infection (Fig. 6C).

295 All the new mutations identified in our analysis were found in human clinical
296 isolates and mapped to relatively recent branches in the MTBC phylogeny (Fig
297 6A). Thus, we reasoned that most mutations were associated with recent
298 selective pressures as compared to the mutations found in the lineage 5, 6 and
299 animal clade reported previously²². We tested whether novel *phoR* mutations are
300 arising in clinical settings during infection and for their potential involvement in
301 ongoing transmission. We used a population-based data set from Malawi²³ where
302 more than 70% of the strains were collected during fifteen years and genome
303 sequenced. We found 14 mutations (13 nonsynonymous and 1 synonymous) in
304 *phoR* exclusive of the Malawi data set with *phoR* having a dN/dS of 3.93.
305 Moreover, the mean relative age of the nonsynonymous *phoR* variants was
306 significantly younger than that of other nonsynonymous variants in the dataset (t-
307 test, p-value << 0.01) and the *phoR* variants from the Malawi data set were more
308 recent than those *phoR* mutations from the reference dataset (t-test, p-value =
309 0.01)(Fig. 6B). From the 13 nonsynonymous mutations in the Malawi dataset, 8
310 were in terminal branches and were markers of recent transmission clusters.

311 Moreover, *phoR* mutations in the Malawi dataset involved larger transmission
312 clusters (permutations test, p-value < 0.001). Taken together, these data indicate
313 that novel *phoR* mutations arise during infection and propagate in on-going
314 human-to-human transmissions in clinical settings.

315

316 **DISCUSSION**

317 We present evidence that the MTBC ancestor transitioned to an obligate
318 pathogenic lifestyle in sympatry from a common genetic pool including the
319 ancestor of extant MCAN strains. Specifically, we found common patterns of
320 genome-wide recombination between the ancestral MTBC genome and the
321 ancestors of extant MCAN strains. The high recombination rate between MCAN
322 strains, including the MTBC ancestor, stands in sharp contrast to the strictly
323 clonal population structure of extant MTBC strains. By analyzing events leading
324 to the transition from a recombinogenic to a clonal organism, we have also been
325 able to identify genomic regions under different selective pressures. Our results
326 suggest that mutations in *phoR* have allowed the bacteria to adapt to different
327 mammalian hosts and they still play an important role during infection and
328 transmission in current clinical settings.

329 Population genomics data has led to the development and testing of different
330 models of how different genetic clusters of the same species can arise in
331 sympatry^{1,3,24}. In the case of *Vibrio cholera*, an appropriate combination of certain
332 virulence-associated variants, ecological opportunity and additional virulence
333 factors mediated the successful transition of particular clones from an
334 environmental to a pathogenic lifestyle²⁵. Other known cases such as pathogenic
335 *Salmonella*²⁶ or *Yersinia* species²⁷ may have followed a similar scheme. The
336 MTBC represents an extreme case of clonal emergence associated to its obligate
337 pathogenic lifestyle. Here, we have shown that, despite the high average
338 nucleotide identity between MCAN and the MTBC, there is complete genomic
339 isolation between them. There is now experimental evidence in the laboratory
340 that genetic exchange between MCAN strains occurs easily but this is not the
341 case between MCAN and the MTBC²⁸. We have shown that there is no
342 measurable on-going recombination among the MTBC strains based on our
343 analysis of 1,591 genomes and in agreement with other recent evidence^{29,30}.
344 Recombination in natural populations depends both on the capacity of
345 chromosomal DNA exchange between the two groups involved and on the
346 ecological opportunity. The mechanisms, if any, on how the MTBC bacilli lost their
347 capacity to recombine when the ancestral genetic pool showed very similar
348 recombination patterns to MCAN strains remains to be elucidated. Ecological
349 opportunity may also influence on the lack of opportunities of exchange between
350 MTBC strains. Despite the occurrence of co-infections, the bacilli occupy mainly
351 an intracellular lifestyle, thereby reducing the opportunities for genetic exchange.

352 We can only speculate on how the transition from a likely environmental or
353 opportunistic pathogen to an obligate pathogen occurred, but our analysis has
354 identified a series of non-random evolutionary events. Notably those events
355 involve core pathogenesis genes. We have identified highly divergent regions in
356 the MTBC ancestor compared to MCAN. The pattern of SNP accumulation
357 suggests that those regions were important in the transition to a closer
358 association with the host. In addition, recombination events in the branch leading
359 to the MTBC ancestor affected essential genes as well as genic regions known
360 to be involved in host-pathogen interaction. The *mymA* operon (Rv3083-Rv3089)
361 is related to the production of mycolic acids and its disruption leads to an aberrant
362 cell-wall structure. Importantly, knock-out studies³¹ have shown that this operon
363 is essential for growth in macrophages and spleen. Furthermore, the deletion of
364 genes in this operon leads to a higher TNF-alpha production, highlighting their
365 role on regulating host-pathogen interaction³². The other major operon identified
366 in our analysis is the *mce1* operon³³. *mce1* knock-out mutants are hypervirulent
367 in a mice model of infection and lose the capacity of a proper pro-inflammatory
368 cytokine production that is needed for the establishment of the infection³⁴ and
369 granuloma³³. How all these processes are mediated by *mce1* is still not clear
370 pointing to *mce1* as a priority target for biomedical research.

371 Our analysis identified one gene, *phoR*, which is under positive selection in extant
372 MTBC strains although it was under purifying selection in the MTBC ancestor.
373 PhoR is the sensor component of the PhoPR two-component systems, which
374 play a major role in MTBC pathogenesis³⁵. Experiments with mutations identified
375 in this study showed that the lipid composition of the membrane and the ESAT-6
376 secretion, a major virulence factor in the MTBC, are affected²². Early mutations
377 in *phoR* were linked to adaptation of the MTBC to different animal species. Here,
378 we show that *phoR* continues to play a role in the ongoing adaptation of MTBC
379 strains during human-to-human transmission. We speculate that recent *phoR*
380 mutations help to fine-tune the immunogenicity of the pathogen during infection,
381 allowing it to manipulate the host response and increase the chances of
382 transmission. Given that PhoPR is involved in membrane composition³⁶,
383 mutations in this regulator might also be involved in the susceptibility to some
384 antibiotics. However, antibiotic selection is an unlikely explanation for the oldest
385 mutations in PhoPR as they likely predate antibiotic usage.

386 Thus, a model can be proposed in which recombination, together with acquisition
387 of new genetic materia^{11,37}, generated a favorable genetic background for the
388 MTBC ancestor to occupy or increase its association to the mammalian host.
389 Contrarily to *Vibrio cholerae*, in which pandemic strains have emerged from the
390 environment multiple times²⁵, we see this emergence only once in the MTBC,
391 perhaps because the right combination of multiple, fortuitous genetic events and
392 ecological conditions has only happened once. More provocative is the idea that
393 MTBC might just be part of a spectrum of association to the host occupied by the

394 different MCAN-MTBC groups. The fact that the so-called Clone A MCAN strains
395 are more common in the clinic may suggest differences in ecological niches within
396 the MCAN group itself³⁸. In agreement, previous publications^{9,39} and our own
397 analysis (Figure 2) have identified Clone A strains as the evolutionary closest
398 MCAN group to MTBC.

399 In the MTBC, the stronger, obligate association with new host(s) was
400 accompanied by new selective pressures. In accordance, we have identified
401 genes in the MTBC genome highly diverging from MCAN and evolving under
402 purifying selection, suggesting that they have become essential following MTBC's
403 transition to an obligate pathogenic life-style. In the final stages of adaptation,
404 positive selection on genes such as *phoR* and others^{39,40} led to a narrowing of
405 the host-range and later still to a further fine-tuning during the spread of the
406 bacteria within the new host species.

407

408 **ACKNOWLEDGEMENTS**

409 We thank Alberto Marina for advice in the interpretation of the PhoR molecular
410 structure. This work was funded by projects the European Research Council
411 (ERC) (638553-TB-ACCELERATE) and Ministerio de Economía y
412 Competitividad (Spanish Government) research grant SAF2016-77346-R (to IC).
413 BFU2014-58656-R and BFU2017-89594-R from Ministerio de Economía y
414 Competitividad (Spanish Government) and PROMETEO/2016/122 from
415 Generalitat Valenciana (to FGC). ACO is recipient of a FPU fellowship from
416 Ministerio de Educación y Ciencia FPU13/00913 (Spanish Government). Swiss
417 National Science Foundation (grants 310030_166687, IZRJZ3_164171,
418 IZLSZ3_170834 and CRSII5_177163), the European Research Council (309540-
419 EVODRTB) and SystemsX.ch (to SG).

420

421 **AUTHOR'S CONTRIBUTION**

422 IC conceived this work. ACO, IC, LSB, SH, JC and FGC analyzed the data. ACO,
423 IC, LSB, FGC wrote the first version of the draft. All authors critically reviewed
424 and contributed to the final version of the manuscript.

425 **DATA AVAILABILITY**

426 The data that support the findings of this study are available from the
427 corresponding author upon reasonable request.

428 **REFERENCES**

- 429 1. Vos, M. A species concept for bacteria based on adaptive
430 divergence. *Trends Microbiol.* **19**, 1–7 (2011).
- 431 2. Shapiro, B. J. & Polz, M. F. Microbial speciation. *Cold Spring Harb.*
432 *Perspect. Biol.* **7**, a018143 (2015).
- 433 3. Marttinen, P. & Hanage, W. P. Speciation trajectories in
434 recombining bacterial species. *PLoS Comput. Biol.* **13**, e1005640
435 (2017).
- 436 4. WHO | Global tuberculosis report 2017. *WHO* (2017).
- 437 5. Comas, I. *et al.* Out-of-Africa migration and Neolithic coexpansion of
438 *Mycobacterium tuberculosis* with modern humans. *Nat Genet* **45**,
439 1176–1182 (2013).
- 440 6. Stucki, D. *et al.* *Mycobacterium tuberculosis* lineage 4 comprises
441 globally distributed and geographically restricted sublineages. *Nat.*
442 *Genet.* **48**, 1535–1543 (2016).
- 443 7. Davies, P. D. O. Tuberculosis in humans and animals: are we a
444 threat to each other? *J. R. Soc. Med.* **99**, 539–40 (2006).
- 445 8. Coscolla, M. & Gagneux, S. Consequences of genomic diversity in
446 *Mycobacterium tuberculosis*. *Semin. Immunol.* **26**, 431–444 (2014).
- 447 9. Supply, P. *et al.* Genomic analysis of smooth tubercle bacilli
448 provides insights into ancestry and pathoadaptation of
449 *Mycobacterium tuberculosis*. *Nat. Genet.* **45**, 172–9 (2013).
- 450 10. Aboubaker Osman, D., Bouzid, F., Canaan, S. & Drancourt, M.
451 Smooth tubercle bacilli: neglected opportunistic tropical pathogens.
452 *Front. public Heal.* **3**, 283 (2015).
- 453 11. Veyrier, F., Pletzer, D., Turenne, C. & Behr, M. A. Phylogenetic
454 detection of horizontal gene transfer during the step-wise genesis of
455 *Mycobacterium tuberculosis*. *BMC Evol. Biol.* **9**, 196 (2009).
- 456 12. Brennan, P. J. Bacterial evolution: Emergence of virulence in TB.
457 *Nat. Microbiol.* **1**, 15031 (2016).
- 458 13. Namouchi, A., Didelot, X., Schöck, U., Gicquel, B. & Rocha, E. P. C.
459 After the bottleneck: Genome-wide diversification of the
460 *Mycobacterium tuberculosis* complex by mutation, recombination,
461 and natural selection. *Genome Res.* **22**, 721–34 (2012).
- 462 14. Coll, F. *et al.* A robust SNP barcode for typing *Mycobacterium*
463 *tuberculosis* complex strains. *Nat. Commun.* **5**, 4812 (2014).

- 464 15. Cortes, T. *et al.* Genome-wide mapping of transcriptional start sites
465 defines an extensive leaderless transcriptome in *Mycobacterium*
466 *tuberculosis*. *Cell Rep.* **5**, 1121–31 (2013).
- 467 16. Mortimer, T. D., Weber, A. M. & Pepperell, C. S. Signatures of
468 selection at drug resistance loci in *Mycobacterium tuberculosis*.
469 *mSystems* **3**, e00108-17 (2018).
- 470 17. Hedrick, P. & Kumar, S. Mutation and linkage disequilibrium in
471 human mtDNA. *Eur. J. Hum. Genet.* **9**, 969–972 (2001).
- 472 18. Comas, I. *et al.* Human T cell epitopes of *Mycobacterium*
473 *tuberculosis* are evolutionarily hyperconserved. *Nat. Genet.* **42**,
474 498–503 (2010).
- 475 19. Coscolla, M. *et al.* *M. tuberculosis* T cell epitope analysis reveals
476 paucity of antigenic variation and identifies rare variable TB
477 antigens. *Cell Host Microbe* **18**, 538–48 (2015).
- 478 20. Tukey, J. W. (John W. *Exploratory data analysis*. (Addison-Wesley
479 Pub. Co, 1977).
- 480 21. Gonzalo-Asensio, J. *et al.* PhoP: a missing piece in the intricate
481 puzzle of *Mycobacterium tuberculosis* virulence. *PLoS One* **3**,
482 e3496–e3496 (2008).
- 483 22. Gonzalo-Asensio, J. *et al.* Evolutionary history of tuberculosis
484 shaped by conserved mutations in the PhoPR virulence regulator.
485 *Proc. Natl. Acad. Sci.* **111**, 11491–11496 (2014).
- 486 23. Guerra-Assunção, J. A. *et al.* Large-scale whole genome
487 sequencing of *M. tuberculosis* provides insights into transmission in
488 a high prevalence area. *Elife* **4**, (2015).
- 489 24. Shapiro, B. J. *et al.* Population genomics of early events in the
490 ecological differentiation of bacteria. *Science* **336**, 48–51 (2012).
- 491 25. Shapiro, B. J., Levade, I., Kovacicova, G., Taylor, R. K. & Almagro-
492 Moreno, S. Origins of pandemic *Vibrio cholerae* from environmental
493 gene pools. *Nat. Microbiol.* **2**, 16240 (2016).
- 494 26. Bäumler, A. & Fang, F. C. Host specificity of bacterial pathogens.
495 *Cold Spring Harb. Perspect. Med.* **3**, a010041 (2013).
- 496 27. McNally, A., Thomson, N. R., Reuter, S. & Wren, B. W. ‘Add, stir
497 and reduce’: *Yersinia* spp. as model bacteria for pathogen evolution.
498 *Nat. Rev. Microbiol.* **14**, 177–190 (2016).
- 499 28. Boritsch, E. C. *et al.* *pks5*-recombination-mediated surface
500 remodelling in *Mycobacterium tuberculosis* emergence. *Nat.*
501 *Microbiol.* **1**, 15019 (2016).

- 502 29. Mortimer, T. D. & Pepperell, C. S. Genomic signatures of distributive
503 conjugal transfer among Mycobacteria. *Genome Biol. Evol.* **6**, 2489–
504 2500 (2014).
- 505 30. Boritsch, E. C. *et al.* Key experimental evidence of chromosomal
506 DNA transfer among selected tuberculosis-causing mycobacteria.
507 *Proc. Natl. Acad. Sci. U. S. A.* **113**, 9876–81 (2016).
- 508 31. Cheruvu, M., Plikaytis, B. B. & Shinnick, T. M. The acid-induced
509 operon Rv3083–Rv3089 is required for growth of *Mycobacterium*
510 *tuberculosis* in macrophages. *Tuberculosis* **87**, 12–20 (2007).
- 511 32. Olsen, A. *et al.* Targeting *Mycobacterium tuberculosis* tumor
512 necrosis factor alpha-downregulating genes for the development of
513 antituberculous vaccines. *MBio* **7**, e01023-15 (2016).
- 514 33. Casali, N., White, A. M. & Riley, L. W. Regulation of the
515 *Mycobacterium tuberculosis mce1* operon. *J. Bacteriol.* **188**, 441–
516 449 (2006).
- 517 34. Shimono, N. *et al.* Hypervirulent mutant of *Mycobacterium*
518 *tuberculosis* resulting from disruption of the *mce1* operon. *Proc.*
519 *Natl. Acad. Sci.* **100**, 15918–15923 (2003).
- 520 35. Broset, E., Martín, C. & Gonzalo-Asensio, J. Evolutionary landscape
521 of the *Mycobacterium tuberculosis* complex from the viewpoint of
522 PhoPR: implications for virulence regulation and application to
523 vaccine development. *MBio* **6**, e01289-15 (2015).
- 524 36. Walters, S. B. *et al.* The *Mycobacterium tuberculosis* PhoPR two-
525 component system regulates genes essential for virulence and
526 complex lipid biosynthesis. *Mol. Microbiol.* **60**, 312–330 (2006).
- 527 37. Orgeur M, Brosch R. Evolution of virulence in the *Mycobacterium*
528 *tuberculosis* complex. *Curr Opin Microbiol.* **41**, 68-75 (2018).
- 529 38. Blouin, Y. *et al.* Progenitor *Mycobacterium canettii* clone responsible
530 for lymph node tuberculosis epidemic, Djibouti. *Emerg. Infect. Dis.*
531 **20**, 21–28 (2014).
- 532 39. Ates, L. S. *et al.* Mutations in *ppe38* block PE_PGRS secretion and
533 increase virulence of *Mycobacterium tuberculosis*. *Nat. Microbiol.* **3**,
534 181–188 (2018).
- 535 40. Malone, K. M. *et al.* Comparative 'omics analyses differentiate
536 *Mycobacterium tuberculosis* and *Mycobacterium bovis* and reveal
537 distinct macrophage responses to infection with the human and
538 bovine tubercle bacilli. *Microb. Genomics* **4**, (2018).
- 539 41. Walker, T. M. *et al.* Whole-genome sequencing for prediction of

- 540 Mycobacterium tuberculosis drug susceptibility and resistance: a
541 retrospective cohort study. *Lancet. Infect. Dis.* **15**, 1193–1202
542 (2015).
- 543 42. Schmieder, R. & Edwards, R. Quality control and preprocessing of
544 metagenomic datasets. *Bioinformatics* **27**, 863–4 (2011).
- 545 43. Li, H. & Durbin, R. Fast and accurate long-read alignment with
546 Burrows-Wheeler transform. *Bioinformatics* **26**, 589–95 (2010).
- 547 44. Li, H. *et al.* The sequence alignment/map format and SAMtools.
548 *Bioinformatics* **25**, 2078–9 (2009).
- 549 45. Koboldt, D. C. *et al.* VarScan 2: somatic mutation and copy number
550 alteration discovery in cancer by exome sequencing. *Genome Res.*
551 **22**, 568–76 (2012).
- 552 46. Cingolani, P. *et al.* A program for annotating and predicting the
553 effects of single nucleotide polymorphisms, SnpEff: SNPs in the
554 genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*
555 (*Austin*). **6**, 80–92 (2012).
- 556 47. Feuerriegel, S. *et al.* PhyResSE: a web tool delineating
557 *Mycobacterium tuberculosis* antibiotic resistance and lineage from
558 whole-genome sequencing data. *J. Clin. Microbiol.* **53**, 1908–14
559 (2015).
- 560 48. Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based
561 phylogenetic analyses with thousands of taxa and mixed models.
562 *Bioinformatics* **22**, 2688–2690 (2006).
- 563 49. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool
564 for the display and annotation of phylogenetic and other trees.
565 *Nucleic Acids Res.* **44**, W242-5 (2016).
- 566 50. Maddison, W. P. & Maddison, D. R. Mesquite: a modular system for
567 evolutionary analysis. (2017).
- 568 51. Purcell, S. *et al.* PLINK: a tool set for whole-genome association
569 and population-based linkage analyses. *Am. J. Hum. Genet.* **81**,
570 559–75 (2007).
- 571 52. Team, R. C. R: a language and environment for statistical
572 computing. (2015).
- 573 53. Darling, A. E., Mau, B. & Perna, N. T. progressiveMauve: multiple
574 genome alignment with gene gain, loss and rearrangement. *PLoS*
575 *One* **5**, e11147 (2010).
- 576 54. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples
577 of recombinant bacterial whole genome sequences using Gubbins.

- 578 *Nucleic Acids Res.* **43**, e15 (2015).
- 579 55. Schmidt, H. A., Strimmer, K., Vingron, M. & von Haeseler, A. TREE-
580 PUZZLE: maximum likelihood phylogenetic analysis using quartets
581 and parallel computing. *Bioinformatics* **18**, 502–4 (2002).
- 582 56. Bos, K. I. *et al.* Pre-Columbian mycobacterial genomes reveal seals
583 as a source of New World human tuberculosis. *Nature* **514**, 494–
584 497 (2014).
- 585 57. Rambaut, A., Suchard, M., Xie, D. & Drummond, A. Tracer v1.6.
586 (2014).
- 587 58. Maere, S., Heymans, K. & Kuiper, M. BiNGO: a Cytoscape plugin to
588 assess overrepresentation of Gene Ontology categories in biological
589 networks. *Bioinformatics* **21**, 3448–3449 (2005).
- 590 59. Shannon, P. *et al.* Cytoscape: a software environment for integrated
591 models of biomolecular interaction networks. *Genome Res.* **13**,
592 2498–2504 (2003).
- 593 60. Storey, J. D. The positive false discovery rate: a Bayesian
594 interpretation and the q-value . 2013–2035 (2003).
- 595 61. Fedrizzi, T. *et al.* Genomic characterization of nontuberculous
596 Mycobacteria. *Sci. Rep.* **7**, 45258 (2017).
- 597 62. Lechner, M. *et al.* Proteinortho: detection of (co-)orthologs in large-
598 scale analysis. *BMC Bioinformatics* **12**, 124 (2011).
- 599 63. Sievers, F. *et al.* Fast, scalable generation of high-quality protein
600 multiple sequence alignments using Clustal Omega. *Mol. Syst. Biol.*
601 **7**, 539 (2011).
- 602 64. Ota, T. & Nei, M. Variance and covariances of the numbers of
603 synonymous and nonsynonymous substitutions per site. *Mol. Biol.*
604 *Evol.* **11**, 613–9 (1994).
- 605 65. Murrell, B. *et al.* FUBAR: a fast, unconstrained bayesian
606 approximation for inferring selection. *Mol. Biol. Evol.* **30**, 1196–1205
607 (2013).
- 608 66. Murrell, B. *et al.* Gene-wide identification of episodic selection. *Mol.*
609 *Biol. Evol.* **32**, 1365–1371 (2015).
- 610 67. Saitou, N. & Nei, M. The neighbor-joining method: a new method for
611 reconstructing phylogenetic trees. *Mol. Biol. Evol.* **4**, 406–425
612 (1987).
- 613 68. Pennell, M. W. *et al.* geiger v2.0: an expanded suite of methods for
614 fitting macroevolutionary models to phylogenetic trees.

- 615 *Bioinformatics* **30**, 2216–2218 (2014).
- 616 69. Finn, R. D. *et al.* The Pfam protein families database: towards a
617 more sustainable future. *Nucleic Acids Res.* **44**, D279–D285 (2016).
- 618 70. Letunic, I. & Bork, P. 20 years of the SMART protein domain
619 annotation resource. *Nucleic Acids Res.* **46**, D493–D496 (2018).

620 FIGURE LEGENDS

621 **Fig. 1. No ongoing recombination within the MTBC** A) Linkage disequilibrium
622 as a function of genetic distance detected in a representative sample of
623 *Mycobacterium tuberculosis* complex strains (n = 1,591). B) Site frequency
624 spectrum of MTBC strains using the 94,781 variant positions.

625 **Fig. 2. Genome-wide variant profiles vary between *M. canetti*, *M.*
626 *tuberculosis* and the MTBC ancestor** A) Number of homoplasies (grey) as a
627 function of the total number of variants detected (orange) in the MCAN dataset,
628 in the branch leading to the MTBC most recent common ancestor and within the
629 MTBC. Black dots indicate recombination events detected in the most recent
630 common ancestor of the MTBC. B) Homoplastic variant positions mapping to the
631 branch of the most recent common ancestor of the MTBC coincide with events
632 detected by Gubbins (highlighted in yellow) and show correlated phylogenetic
633 patterns within each event.

634 **Fig. 3. Past recombination between *M. canetti* strains and the MTBC
635 ancestor** A) Relative age of the recombination fragments detected (black boxes).
636 The red error bars represent the 95% highest probability density (HPD). The age
637 was relativized from 0 to 1, being 0 the age of the non-recombinant fraction of the
638 genome and 1 the highest value in the confidence intervals. B) Histogram
639 distribution of the recombination fragments relative ages. C) Gene Ontology
640 terms overrepresented in the coding regions contained in the recombinant
641 fragments.

642 **Fig. 4. Divergent positions between the MTBC ancestor and *M. canetti*
643 clade.** Average of divSNPs per 10 kb positions (green) as compared to the
644 average of homoplastic variants (gray). Blue arrows above the distribution are
645 genes that significantly accumulate more divSNPs.

646 **Fig. 5. Genes with differential selective pressures across the MTBC
647 speciation stages.** Genes changing selective pressure in the branch of the
648 MTBC ancestor as compared to extant MTBC strains. Red lines mark those
649 genes being outliers of the dN/dS variation distribution.

650 **Fig. 6. *phoR* is under positive selection in human affecting strains.** A)
651 Genome-based phylogeny calculated from a total of 4,598 clinical samples
652 obtained from different sources. The synonymous and nonsynonymous variants
653 found in are *phoR* mapped to the corresponding branch. Variants in internal
654 branches affect complete clades which are coloured in the phylogeny.
655 Homoplasies are marked in the outer circle of the phylogeny. B) Relative ages
656 distribution of the *phoR* variants in the reference dataset from Coll *et al.*¹⁴ and the
657 transmission dataset²³ in comparison with the rest of the genome variants. C)
658 Schematic view of PhoR with the aminoacid changes found across the 4,598
659 samples dataset marked on it. Aminoacid changes are significantly more
660 abundant in the sensor domain (p-value < 0,01).

661

662 METHODS

663 Datasets used

664 **Global collection 1 (1,591 isolates).** This collection was obtained from Coll *et*
665 *al.*¹⁴ who gathered it from different sources. We downloaded all the available
666 FASTQ files associated to the following ENA accession numbers: ERP000192,
667 ERP000276, ERP000520, ERP001731, ERP000111, SRP002589, ERP002611,
668 ERP000436, SRA065095, ERP001885, ERP001567. A total of 3,475 genomes
669 were downloaded, aligned to the reference and their variants extracted. The
670 transmission clusters were analyzed (transmission event defined by less than 15
671 SNPs) and a representative genome of each transmission cluster was kept. We
672 ended with 1,591 strains, representative of the 7 lineages.

673
674 **Transmission data set (1,646 isolates).** This data set was obtained from
675 Guerra-Assunção *et al.*²³ It includes samples taken over a 15-year period in a
676 district of Malawi. We downloaded all the available FASTQ files associated to the
677 following ENA accession numbers: ERP000436 and ERP001072.

678
679 **Global collection 2 (219 isolates).** This dataset from Comas *et al.*⁵ represent a
680 selected set of strains representing the known phylogenetic diversity of the
681 human MTBC strains.

682
683 **Global collection 3 (4,762 isolates).** This dataset was obtained by joining the
684 global collection 1 and the transmission dataset described above, as well as the
685 isolates obtained from Walker *et al.*⁴¹. To the fastq files obtained previously, we
686 added the ones obtained from the accession numbers found in the supplementary
687 material from Walker *et al.* A total of 7,977 genomes were put together, aligned
688 to the reference and their variants extracted. In this case, the transmission
689 clusters (<15 SNPs) were filtered out and only one representative strain of each
690 cluster was kept. Also, we removed potential coinfections. Coinfections were
691 assessed by looking at lineage/sub-lineage coexisting markers¹⁴. Samples with
692 evidence of more than one variant were considered as possible co-infections and
693 removed from ensuing analyses.

694
695 ***Mycobacterium canettii* dataset.** Nine *M. canettii* draft genomes were
696 downloaded from Genbank (CIPT 140010059, NC_015848.1; CIPT 140070010,
697 NC_019951.1; CIPT 140060008, NC_019950.1; CIPT 140070017,
698 NC_019952.1; CIPT 140070008, NC_019965.1; CIPT 140070002,
699 NZ_CAOL00000000.1; CIPT 140070005, NZ_CAOM00000000.1; CIPT
700 140070013, NZ_CAON00000000.1 and CIPT 140070007,
701 NZ_CAOO00000000.1). However, genomes from strains CIPT 140070010 and
702 CIPT 140070017 were discarded because they showed a larger proportion of
703 SNPs than expected, many of them potential sequencing errors.

704 **MTBC most likely ancestral genome.** The MTBC ancestor was derived in a
705 previous publication¹⁸. This ancestor is H37Rv-like in terms of genome
706 structural variants, but H37Rv alleles were substituted by those present in the
707 inferred common ancestor of all MTBC lineages.

708

709 **FASTQ mapping and variant calling for the MTBC strains**

710 Fastq files were trimmed to remove low quality reads using prinseq⁴² and aligned
711 to the MTBC most likely ancestral genome¹⁸ using BWA-mem algorithm⁴³.
712 Alignments with less than 20x mean coverage per base were filtered out. The
713 variant calling was performed using samtools⁴⁴ and VarScan⁴⁵. Due to the low
714 variability found in *M. tuberculosis*, to avoid mapping errors and false SNPs a
715 variant was filtered out if: i) it was supported by less than 20 reads; ii) it was found
716 in a frequency of less than 0.9; iii) it was found near indel areas (10 bp window);
717 or iv) it was found in areas of high accumulation of variants (more than 3 variants
718 in a 10 bp defined window). Variants were annotated using SnpEff⁴⁶. Variants
719 present in PE/PPE genes, phages or repeated sequences were also filtered-out,
720 as they tend to accumulate SNPs due to mapping errors. High quality variant calls
721 were combined in a non-redundant variant list and used to retrieve the most likely
722 allele at each strain to generate a variant alignment.

723

724 **Phylogenetic inference and parsimony mapping of SNPs**

725 In the Global-1 dataset we identified 140,239 variants following the steps defined
726 above. As we wanted to identify nucleotide variants due to recombination events,
727 a stricter filtering was applied to remove putative recombination signal due to
728 polymorphisms introduced by other causes. Variants related with antibiotic
729 resistance were obtained from PhyResSe⁴⁷ and were removed from the analysis.
730 Also, non-biallelic variants were removed from the analysis. To avoid false
731 positives, we also removed positions in which a variant was called in at least one
732 strain but also with a gap in at least other strain. To identify variants coming from
733 mapping errors we generated fragments of 50 bp downstream, upstream and
734 midstream of the variant positions in the reference genome. With these
735 fragments, we performed a BLAST search over the reference genome to check
736 whether they mapped to other regions. Variants identified in reads that mapped
737 to more than one region of the reference genome (query coverage per HSP over
738 98% and percentage of identical matches between the query and the reference
739 genome of 98%) were removed from the analysis.

740 The remaining variants (94,780) were used to infer a phylogenetic tree using
741 RAxML⁴⁸ with the GTRCATI (GTR + optimization of substitution rates +
742 optimization of site-specific evolutionary rates) model of evolution and
743 represented with the iTOL software⁴⁹. Variants were mapped to the phylogeny
744 using the Mesquite suite⁵⁰. Homoplastic variants were identified based on
745 parsimony criteria. Using these homoplastic variant positions, we looked for
746 homoplastic runs. A homoplastic run was defined if two (or more) homoplastic
747 variants were found in the genome in correlative positions or with at least one
748 variant between them. Variants present in the same homoplastic run were
749 mapped on the phylogeny using Mesquite to look for coincident phylogenetic
750 patterns.

751

752 **Linkage-disequilibrium calculation**

753 Using the filtered variant positions (94,780), we used the PLINK software⁵¹ to
754 calculate the linkage-disequilibrium statistics D' and R^2 . To estimate these values,
755 we took into account variants with a minimum frequency of 0.01 and we used a
756 sliding window of 10 Kb. The results obtained were processed with R⁵². To plot
757 the D' and R^2 pattern by variant distance, we calculated average D' and R^2 values
758 for 50 bp sliding windows.

759

760 **Multiple alignment of *M. canettii* and MTBC**

761 Seven *M. canettii* draft genomes were aligned to each other and to the ancestor
762 of MTBC using progressiveMauve⁵³. The segmented alignment obtained in
763 XMFA format was converted to a plain Fasta format using the MTBC ancestor as
764 reordering reference with a custom Perl script. Positions with gaps in the
765 reference sequence were removed from the final alignment, so the resulting
766 aligned genomes had the same size than the reconstructed MTBC ancestor
767 (4,411,532 Mb). The MTBC pseudogenomes reconstructed from mapping to the
768 MTBC ancestor from the different datasets described above were concatenated
769 to the *M. canettii* alignment obtained in the previous step for further analyses.

770

771 **Recombination analyses and phylogenetic evaluation**

772 Recombination was evaluated in the alignment containing 219 strains from MTBC
773 and 7 *M. canettii* and in the one containing the MTBC ancestor and 7 *M. canettii*.
774 First, repetitive regions (i.e. PPE/PGRS) were masked from both alignments and,
775 second, recombination events were inferred using Gubbins⁵⁴, which identifies
776 clusters of high SNP density as markers.

777 Gubbins identify 70 potential recombinant regions in the alignment containing the
778 7 *M. canettii* strains and the MTBC ancestor. Four of these regions were obviated
779 as they fell in regions deleted in several *M. canettii* strains. One more region was
780 removed from the analysis because it was extremely short (41 bp) and we did not
781 obtain reliable results in the subsequent analysis.

782 For the remaining 65 fragments a phylogeny was calculated using RAxML⁴⁸ and
783 applying the GTRCATI model. Also, a reference phylogeny was calculated with
784 the same method using the complete genomes after subtracting these 65 regions.
785 This reference phylogeny had the same topology as the one obtained from the
786 complete genomes. To test for phylogenetic incongruence between the putative
787 recombination fragments and the genome phylogeny we applied the Shimodaira-
788 Hasegawa and Expected Likelihood Weight tests implemented in TREE-
789 PUZZLE⁵⁵.

790

791 **Dating analyses**

792 To infer the age of the 65 recombinant fragments we first reasoned that most of
793 the mutations found were contributed by recombination and not by mutation once
794 the fragment had been integrated in the genome. Thus, before dating the

795 fragments we first removed all the homoplastic variants with other MCAN strain
796 found in the fragments. The final alignments for the 65 fragments consisted of
797 only those variants accumulated after the recombination event. We then used the
798 non-recombinant part of the genome to infer a substitution rate assuming two
799 different dating scenarios published for the tMRCA^{5,56}. We run BEAST for each
800 fragment pre-specifying monophyletic groups and substitution rate based on the
801 non-recombinant genome phylogenetic reconstruction. We used an uncorrelated
802 log-normal distribution for the substitution rate in all cases and a skyline model
803 for population size changes. We ran several chains of up to 10E6 generations
804 sampling every 1E3 generations to ensure independent convergence of the
805 parameters. Convergence was assessed using Tracer⁵⁷. For both evolutionary
806 scenarios, the results obtained were largely congruent and proportional to the
807 age limit imposed for the MTBC ancestor. As there is controversy about the
808 correct MTBC ancestor age, the results were transformed to relative ages for
809 plotting the final results.

810

811 **Gene ontology enrichment analysis**

812 Genes present in the recombinant regions between *MCAN* and the MTBC
813 ancestor were annotated using SnpEff⁴⁶. A Gene Set Enrichment analysis (GSE)
814 was performed to look for enriched gene functions in these regions. The BiNGO
815 tool⁵⁸ was used to study the enrichment in certain functional categories
816 comparing the most abundant terms in the recombinant regions in comparison
817 with those contained in the complete annotation. The tool uses a hypergeometric
818 test (sampling without replacement) and the BH correction for multiple testing
819 comparisons. The Cytoscape program⁵⁹ was used to visualize the results.

820

821 **divSNP analysis**

822 From the *M.canetti* and MTBC ancestor alignment, we extracted those positions
823 having one variant in all the *M.canetti* strains and another variant in the MTB
824 ancestor. The divSNP frequency by nucleotide was calculated by dividing the
825 total number of divSNPs (5688) by the total number of bases in the alignment.
826 Next, the expected abundance of divSNPs for each gene was calculated by
827 multiplying the nucleotide divSNP frequency by the number of nucleotides in each
828 gene. From the expected and the observed divSNP abundance, we used a
829 Poisson distribution to calculate the probability of having the observed divSNPs
830 by chance for each gene. We selected genes having a pFDR $\leq 0,01$ using the
831 q-value Storey method⁶⁰.

832 Complete mycobacterial genomes for reference strains⁶¹ (Supplementary table 9)
833 were downloaded from RefSeq and GenBank. The orthologous genes were
834 calculated from the amino acid sequences and using the Proteinortho tool⁶². A
835 gene was considered as orthologous if the BLAST analysis showed a minimum
836 identity of 25%, a query coverage of 50% and a maximum e-value of 1E-05. The
837 orthologous genes were aligned using Clustal-omega⁶³ and the phylogenies were
838 constructed using RAxML and applying the PROTCATIAUTO model. The
839 reference phylogeny was constructed using only the core genome (proteins

840 having orthologous in all the mycobacterial genomes downloaded) with RAxML
841 using the same options as above. The reference and alternative phylogenies
842 calculated with the orthologous for the divSNPs enriched genes were manually
843 inspected to check for congruence.

844

845 **dN/dS analysis**

846 The dN/dS statistics were calculated using the R statistical language⁵². The
847 potential synonymous and non-synonymous substitution sites for each region
848 were calculated using the SNAP tool⁶⁴. The dN/dS ratio for each region was
849 calculated as follows:

850

$$851 \frac{\text{Non – synonymous variants} / \text{Non – synonymoussites}}{\text{Synonymous variants} / \text{Synonymoussites}}$$

852

853 The dN/dS for the MTBC ancestor was calculated using the divSNPs while the
854 dN/dS for the MTBC were calculated using the 94,780 SNPs defined as the core
855 variant set. To look for a robust comparison between both ratios only genes
856 having at least 3 synonymous and 1 non-synonymous variants were taken into
857 account. To compare the dN/dS ratios, both were normalized by the genomic
858 dN/dS for each category (0,24 for the MTBC ancestor and 0,62 for the MTBC).
859 The difference between the dN/dS ratio was calculated by subtracting the MTBC
860 dN/dS to that of the MTBC ancestor. The genes that account for the largest
861 differences in the dN/dS were identified as outliers ($Q2 - 1.5 * IQR$, $Q3 + 1.5 * IQR$)²⁰
862 of the differences distribution.

863

864 ***phoR* positive selection analysis**

865 Positive selection over *phoR* was tested using FUBAR⁶⁵ and BUSTED⁶⁶. FUBAR
866 was run with 5 MCMC chains of length 10,000,000. 1,000,000 states were used
867 as burn-in and a Dirichlet prior of 0.5. BUSTED was run with default parameters.

868 To study the potential effect of *phoR* mutations on transmission efficacy we used
869 the data set from Guerra-Assunção *et al.*²³. We identified SNPs in branches
870 leading either to leaves or to transmission clusters. Transmission clusters were
871 categorized in large, medium or small according to the number of isolates in the
872 cluster (large = over 75th percentile, medium=between 25th and 75th percentile,
873 small = under 25th percentile). Each gene was scored to check for accumulation
874 of mutations in branches leading to large transmission clusters according to the
875 expression:

876

$$877 \text{Score} = \text{Largeclusters} * 3 + \text{Mediumclusters}$$

878

879 Genes with high mutation rates have a higher number of polymorphisms that
880 could lead to a larger score by chance. To test the probability of obtaining the
881 observed score by chance, a permutation test was carried out 10,000 times. Each
882 of the SNPs identified was reassigned randomly to the same branches and the
883 score was recalculated for each gene. The expected score distribution for each
884 gene was compared to the observed score to calculate the probability. This test
885 was performed for transmission events defined at 10 SNPs.

886 The relative ages for the variant positions were calculated as node-to-tip
887 distances. In order to have a common framework, a phylogeny was constructed
888 including all the samples from the transmission dataset and from the reference
889 dataset. The phylogeny was constructed using the Neighbour-Joining
890 algorithm⁶⁷. For each variant position we first identified the node in which the
891 variant appeared. The node-to-tip distance was calculated afterwards for each
892 node using the geiger package⁶⁸. Distances were normalized to obtain a relative
893 distance. Later, all the non-synonymous variants except the *phoR* polymorphisms
894 were used as a reference set. The nonsynonymous *phoR* variants to be
895 compared were categorized in two groups, those exclusive to the reference
896 dataset¹⁴ and those derived from the transmission dataset²³.

897

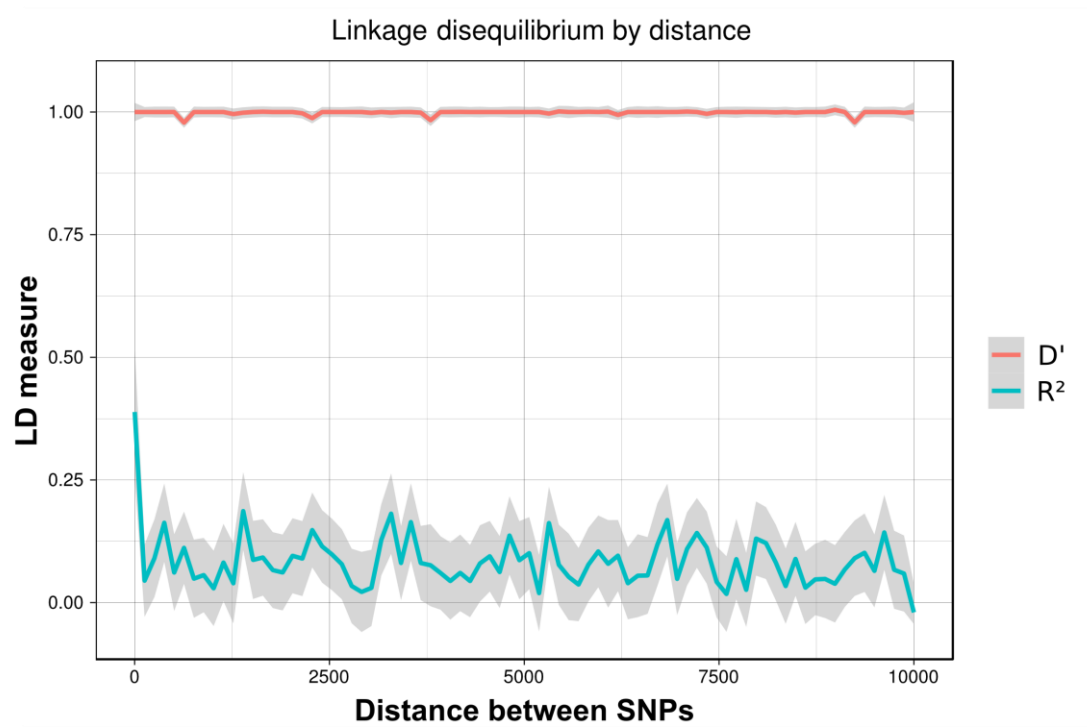
898 **PhoR structure representation**

899 The PhoR structure was inferred by using PFAM⁶⁹ and SMART⁷⁰.

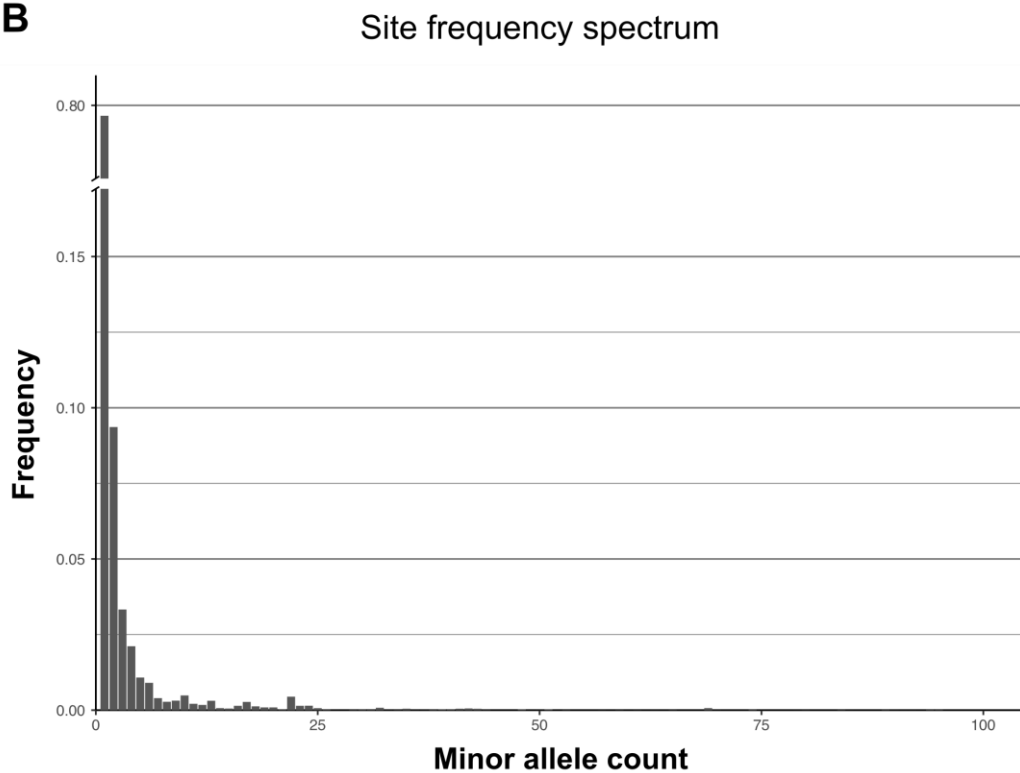
900

901 **FIGURES**

A

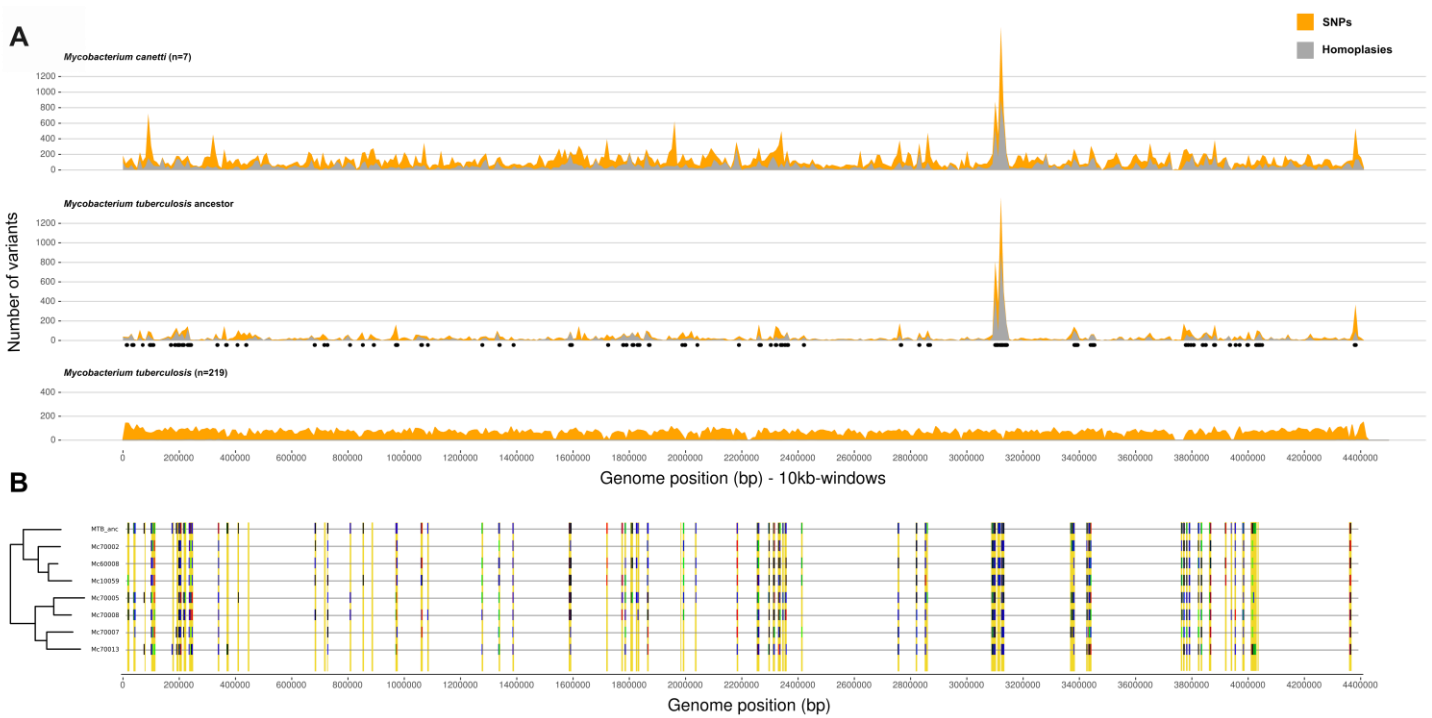


B



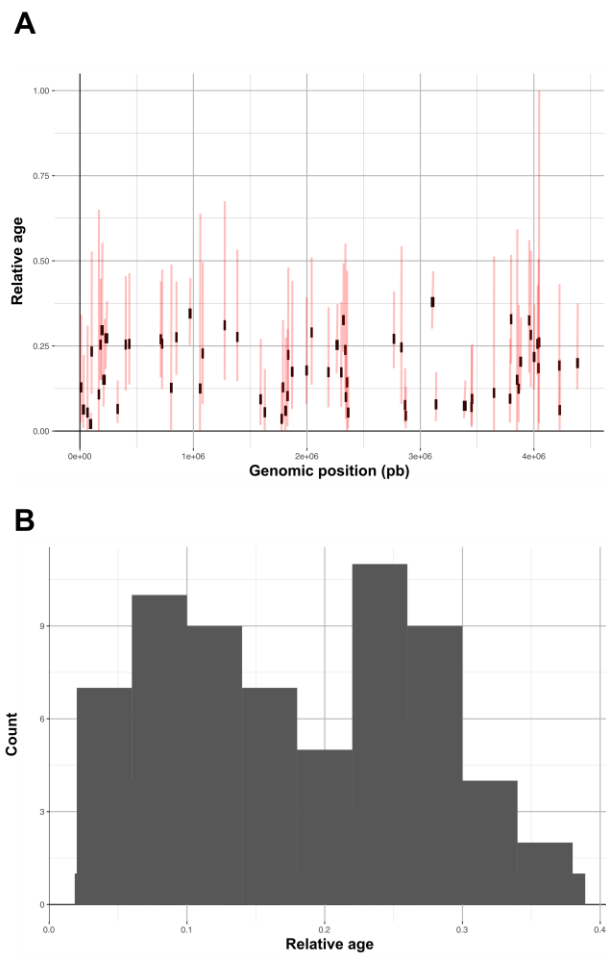
902 Figure 1

903



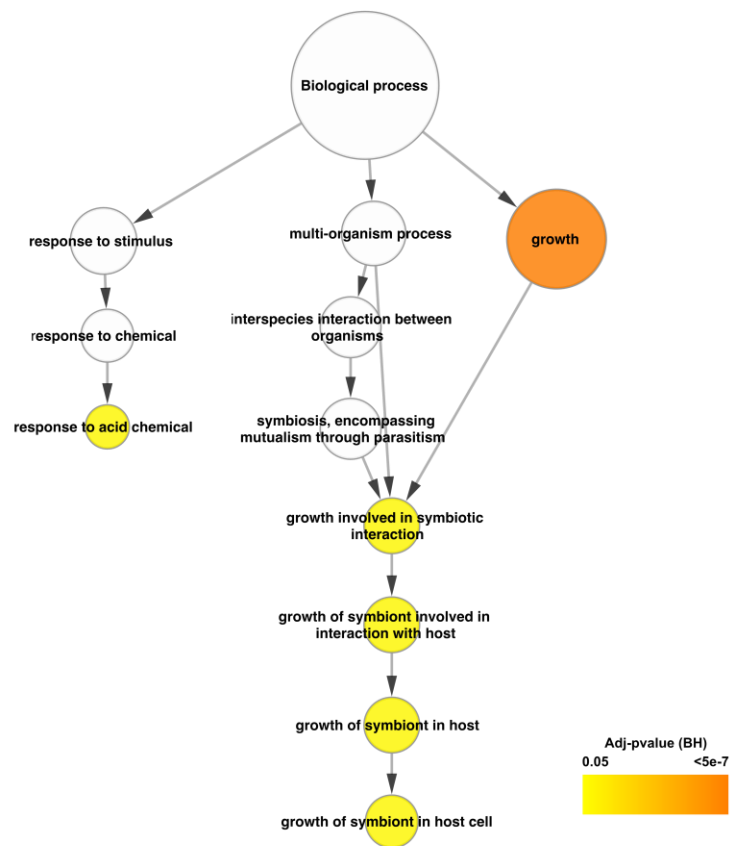
904 Figure 2

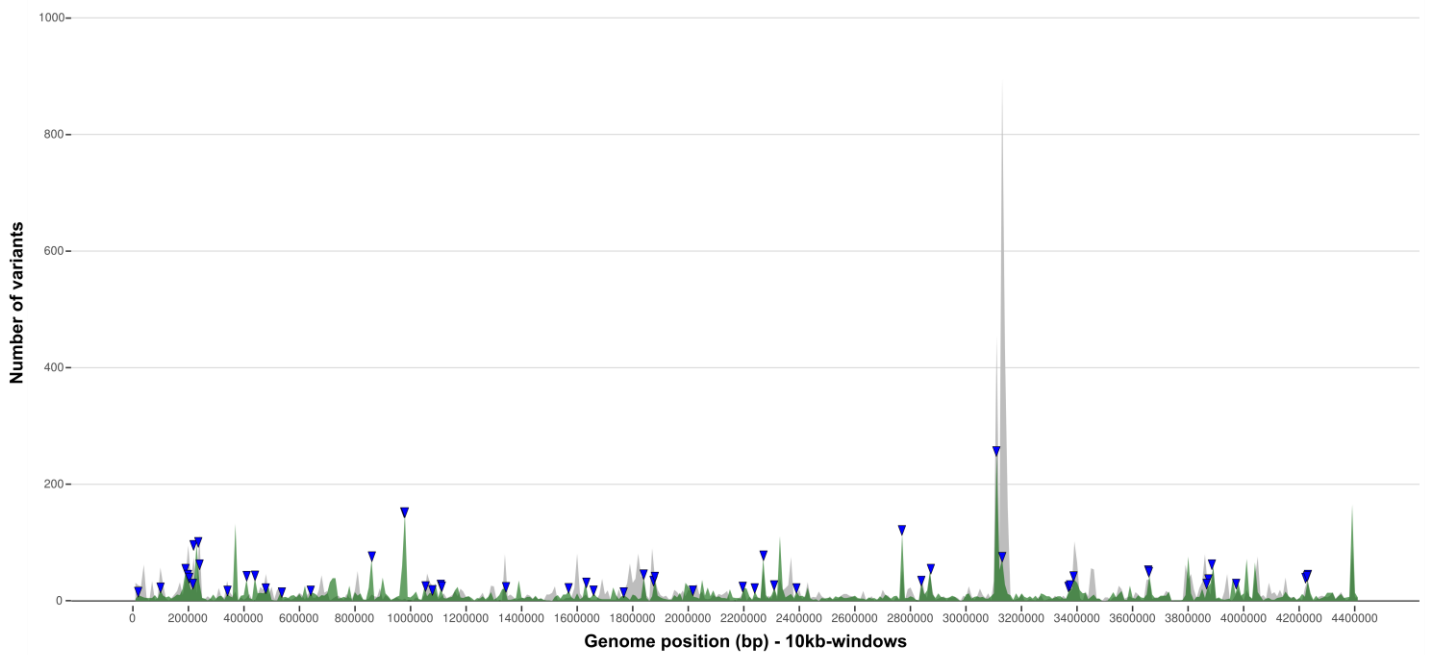
905



906 Figure 3

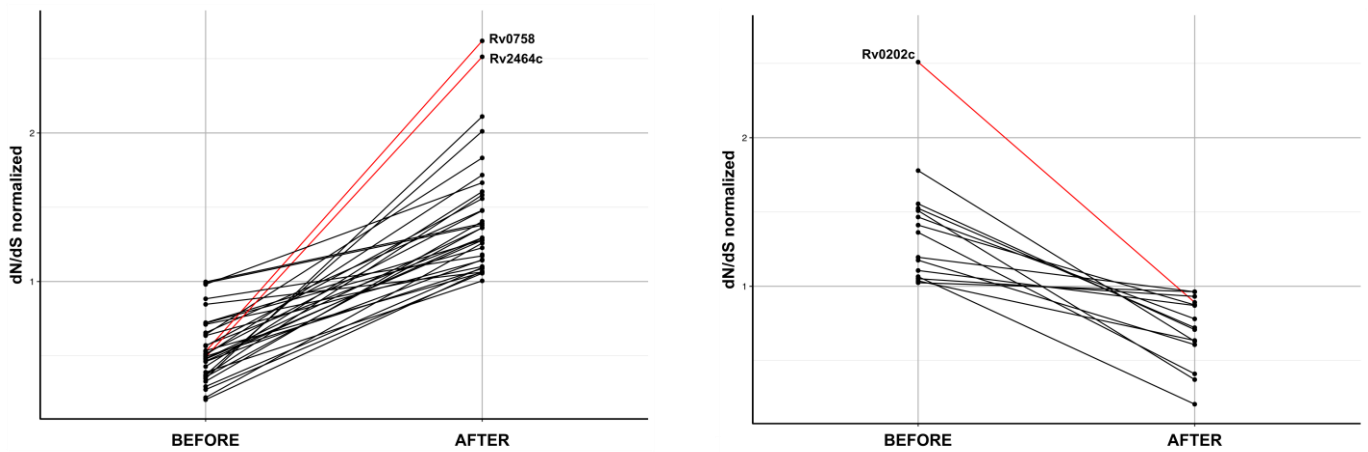
C





907 Figure 4

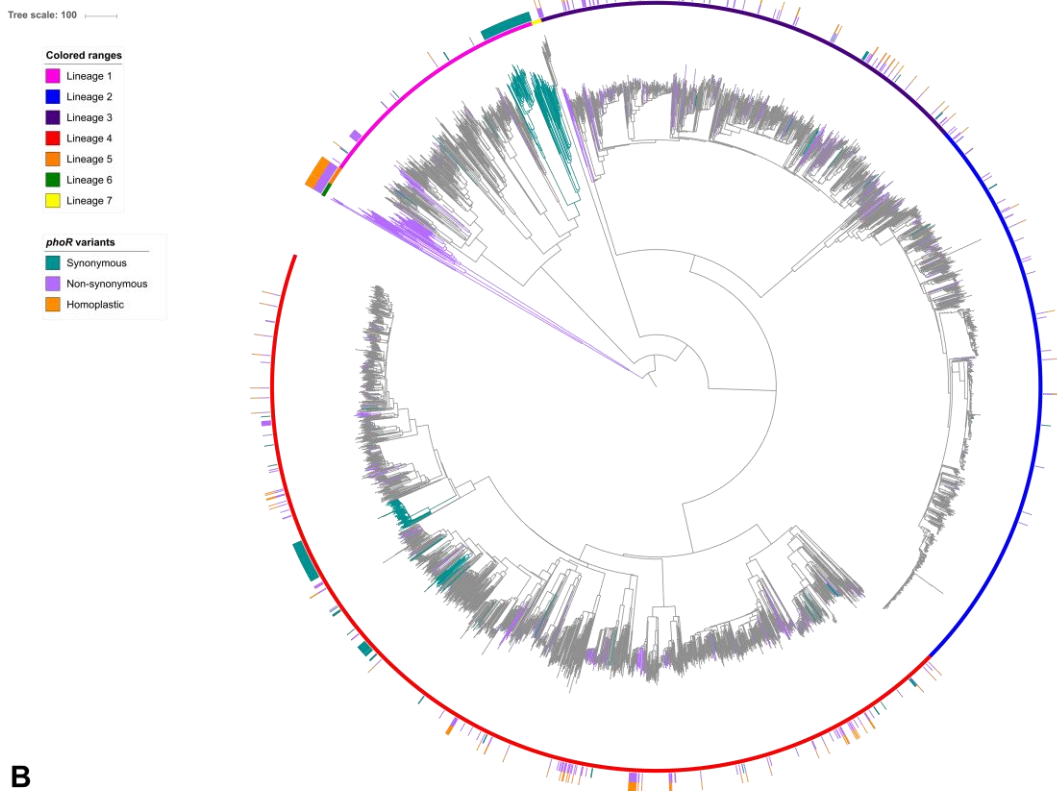
908



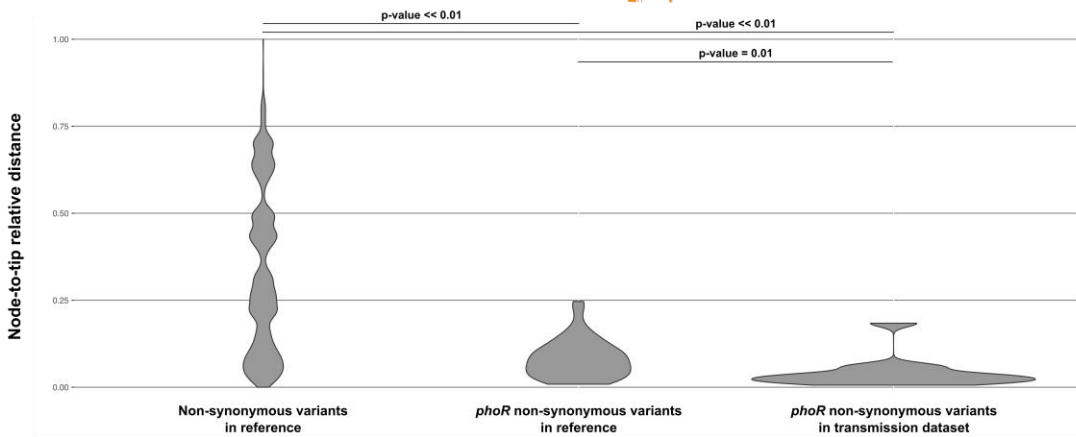
909 Figure 5

910

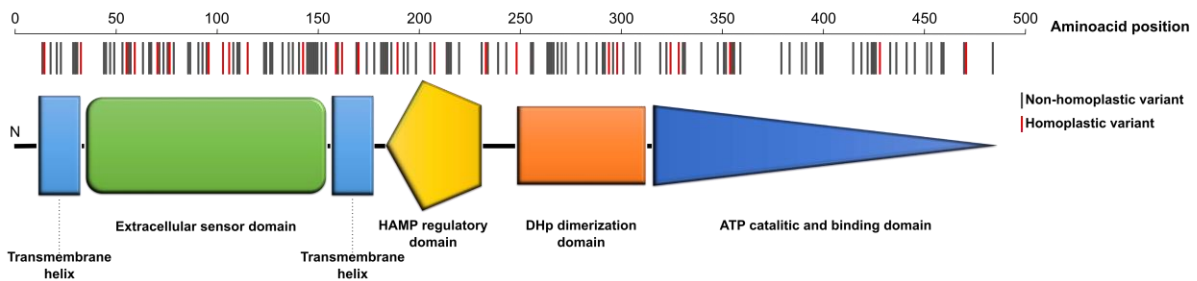
A



B



C



911 Figure 6