

# 1 Large contribution of repeats to genetic variation in a 2 transmission cluster of *Mycobacterium tuberculosis*

3 Christoph Stritt<sup>1,2\*</sup>, Michelle Reitsma<sup>1,2</sup>, Galo Goig<sup>1,2</sup>, Anna Dötsch<sup>1,2</sup>, Sonia Borrell<sup>1,2</sup>,  
4 Christian Beisel<sup>3</sup>, Daniela Brites<sup>1,2</sup>, and Sebastien Gagneux<sup>1,2\*</sup>

5 <sup>1</sup>Swiss Tropical and Public Health Institute, Allschwil, Switzerland

6 <sup>2</sup>University of Basel, Basel, Switzerland

7 <sup>3</sup>Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

8 \*Corresponding authors: [sebastien.gagneux@swisstph.ch](mailto:sebastien.gagneux@swisstph.ch), [christoph.stritt@swisstph.ch](mailto:christoph.stritt@swisstph.ch)

10 Running title: Repeat-associated mutations in the MTBC

## 11 Abstract

12 *Repeats are the most diverse and dynamic, but also the least well understood component of microbial*  
13 *genomes. For all we know, repeat-associated mutations such as duplications, deletions, inversions,*  
14 *and gene conversion might be as common as point mutations, but because of short-read myopia*  
15 *and methodological bias they have received much less attention. Long-read sequencing opens the*  
16 *perspective of resolving repeats and systematically investigating the mutations they induce. For this*  
17 *study, we assembled the genomes of 16 closely related strains of the bacterial pathogen *Mycobacterium**  
18 *tuberculosis from PacBio HiFi reads, with the aim of characterizing the full spectrum of DNA*  
19 *polymorphisms. We find that complete and accurate genomes can be assembled from HiFi reads, with*  
20 *read size being the main limitation in the presence of duplications. By combining a reference-free*  
21 *pangenome graph with extensive repeat annotation, we identified 110 variants, 58 of which can*  
22 *be assigned to repeat-associated mutational mechanisms such as strand slippage and homologous*  
23 *recombination. While recombination events are less frequent than point mutations, they can affect*  
24 *large regions and introduce multiple variants at once, as shown by three gene conversion events and a*  
25 *duplication of 7.3 kb that involve *ppe18* and *ppe57*, two genes possibly involved in immune subversion.*  
26 *Our study shows that the contribution of repeat-associated mechanisms of mutation can be similar to*  
27 *that of point mutations at the microevolutionary scale of an outbreak. A large reservoir of unstudied*  
28 *genetic variation in this “monomorphic” bacterial pathogen awaits investigation.*

29

## Introduction

30 DNA repeats are the most diverse and dynamic components of any genome, not counting  
31 viruses. They comprise a veritable zoo of elements of different origins and complexities, ranging  
32 from short tandem repeats and subtle palindromes to autonomously replicating transposable  
33 elements, inteins, and members of multigene families (for a comprehensive review see Treangen  
34 et al. 2009). While the molecular evolution of repeats is highly variable, they share the property  
35 of providing a substrate for homologous or illegitimate recombination. This makes them the  
36 principal cause of genome instability and DNA polymorphisms in the form of duplications,  
37 deletions, inversions, and non-reciprocal transfers between homologs through gene conversion  
38 (Darmon and Leach 2014).

39 Whole genome sequencing has shown that repeat-associated variation is ubiquitous. This  
40 insight stems from dedicated studies (e.g. Achaz 2002, Schmid et al. 2018), but maybe even  
41 more from frustrated attempts to assemble genomes or identify variants using reads that are  
42 too short to span repeats, which results in fragmented assemblies and ambiguous mappings.  
43 Along with short-read myopia comes a methodological bias towards point mutations, which  
44 are simpler to model and underlie downstream analyses such as phylogenetic inference and  
45 selection scans. For most organisms, little remains therefore known about the types, rates, and  
46 phenotypic effects of repeat-associated mutations. Recently, a systematic investigation of repeats  
47 has come within reach thanks to long-read sequencing and analytical tools such as pangenome  
48 graphs (Garrison, Guarracino, et al. 2023) and hierarchical alignment (Armstrong et al. 2020).  
49 For the streamlined genomes of prokaryotes, base-perfect assemblies can be created (Wick et al.  
50 2023), and pangenome graphs can be used to obtain a concise representation of all variant types  
51 (Yang et al. 2023).

52 One area in which a full characterization of genetic diversity would be particularly useful  
53 is the study of bacterial pathogens. Some of the most deadly of them – including the agents  
54 of anthrax, typhoid, plague, leprosy, and tuberculosis – have been designated “monomorphic”  
55 because of their low levels of genetic diversity (Achtman 2012). Lack of horizontal gene transfer  
56 and evolution under extreme clonality contribute to the phenomenon of “monomorphy”, possibly  
57 through strong background selection (Stritt and Gagneux 2023). In the absence of horizontal  
58 gene transfer, intrachromosomal recombination between repeats might be a key mutational  
59 mechanism in these organisms. As in other organisms, however, the study of repeats and

60 structural variants has been neglected since the advent of short-read sequencing, and we remain  
61 largely ignorant about their contribution to genetic and phenotypic variation.

62 In this study, we use Pacific Biosciences (PacBio) HiFi sequencing to characterize the full  
63 spectrum of DNA polymorphisms in 16 strains of *Mycobacterium tuberculosis*, the agent of  
64 tuberculosis (TB), which with other closely related lineages forms the *Mycobacterium tuberculosis*  
65 complex (MTBC, Gagneux 2018). We focus on a transmission cluster in the city of Bern,  
66 Switzerland, previously characterized through RFLP (Genewein et al. 1993) and Illumina  
67 sequencing (Stucki et al. 2015; Kühnert et al. 2018). The cluster reflects mainly transmission  
68 among homeless and substance abusers in the 1990ies, with spillovers to the general population  
69 and reactivated TB diagnosed up to 2012 (Stucki et al. 2015). This is not the typical manifestation  
70 of TB, which today mainly affects poor countries and is, after COVID-19, the second most deadly  
71 infectious disease in the world (WHO 2023). How these bacteria manage to be so successful with  
72 so little genetic diversity remains puzzling. Part of the answer may lie in the genomic “dark  
73 matter”, the repetitive 10% in the genomes of these bacteria.

74 Here we unlock the repeatome of the MTBC by addressing a simple question: what types  
75 of variants are there? More specifically, we 1) evaluate the accuracy of assemblies constructed  
76 from PacBio HiFi reads; 2) characterize the repeat landscape of the MTBC; and 3) describe the  
77 different types of DNA polymorphisms and their underlying mutational mechanisms. Our  
78 results show that the contribution of neglected genomic regions and types of mutations can be  
79 similar to that of point mutations in non-repetitive regions.

## 80 Results

### 81 Complete and accurate assemblies from CCS reads

82 16 strains from the Bernese outbreak, isolated between 1988 and 2005, were selected for  
83 sequencing on a PacBio Sequel II in circular consensus sequencing (CCS) mode (Supplemental  
84 Fig. S1, Supplemental Table S1). Using the Flye assembly algorithm, all but one genome (P001-  
85 N1377) could be assembled into single circular chromosomes, despite considerable variation in  
86 read lengths and sequencing depths (Table 1).

87 To test whether assemblies are not only complete but also accurate, we aligned the long reads  
88 back against the respective assemblies and called variants to discover inconsistencies between the

89 two. Five sites in four assemblies were identified where the reads contradict the assembly, while  
 90 12 assemblies are free of inconsistencies and appear to be correct to the base. A closer inspection  
 91 of the inconsistencies suggests that they are due to different causes (Supplemental Fig. S2). In  
 92 P003-N1374, the assembly with the lowest mean coverage (16x), a total of five reads disagree on  
 93 a sequence of five versus four cytosines. A second case, in P028-N1362, seems to reflect genuine  
 94 heterogeneity as it might arise during culture or from a heterogeneous inoculum, with 28 reads  
 95 supporting an adenine versus 28 reads supporting a guanine. A third inconsistency arises from  
 96 a duplication in P001-N1377 (discussed below). This sequence was not resolved by Flye but  
 97 in the subsequent circularization step, where 12 bp evident in the reads went missing. Finally,  
 98 two nearby single-base insertions in P034-N1426 suggested by the assembly but not the reads  
 99 reflect misassembly: one single read shows the presence of the two additional bases, while 79  
 100 contradict it. The last two inconsistencies, which are clear assembly errors, were corrected: 12 bp  
 101 were added to the duplication in P001-N1377, and two bases were deleted from P034-N1426.

Strain	Reads			Assembly		
	Nr Reads	N50	Longest	Length	Coverage	Circular
P007-N1108	44,212	5,344	21,567	4,405,381	47	Y
P028-N1362	43,856	5,961	27,471	4,405,379	56	Y
P003-N1374	11,692	6,726	24,804	4,405,378	16	Y
P001-N1377	26,848	6,839	27,922	4,412,646	38	N
P008-N1380	24,578	7,534	27,155	4,405,372	38	Y
P010-N1385	53,811	5,391	21,423	4,405,379	58	Y
P020-N1386	29,983	7,884	27,587	4,405,380	48	Y
P059-N1392	76,540	4,781	15,395	4,405,268	70	Y
P073-N1394	38,388	7,422	27,900	4,405,232	59	Y
P074-N1402	31,862	7,594	29,206	4,405,232	50	Y
P066-N1411	26,719	7,075	28,501	4,405,380	39	Y
P034-N1426	66,552	6,740	27,519	4,405,492	90	Y
P052-N1429	26,244	7,941	27,868	4,405,378	43	Y
P042-N1430	18,475	7,657	26,514	4,405,378	29	Y
P022-N1431	40,133	7,474	27,860	4,405,379	62	Y
P006-N1591	42,575	7,713	27,040	4,405,380	68	Y

Table 1: Sequencing and assembly statistics

## 102 **The repeat landscape of *Mycobacterium tuberculosis***

103 Considering the key role of repeats in causing structural variation and gene conversion, we  
104 first sought to understand the repeat landscape in the studied genomes. Different types of  
105 repeats were annotated in one newly assembled genome, P034-N1426 (Fig. 1): homopolymers  
106 of at least 5 bp, short sequence repeats (SSRs, direct repeats of 3 to 9 bp), tandem repeats (TRs,  
107 direct repeats > 9 bp), and insertion sequences. Homopolymers were the most abundant type  
108 with 6,770 occurrences. 67 SSRs were identified, the large majority of them triplets of six or  
109 nine that do not shift reading frames (Fig. 2A). TRs were found at 47 locations, with quite  
110 some variation in repeat periodicity and length (Fig. 2A). Finally, 57 insertion sequences were  
111 identified, including 12 copies of IS6110 (IS3 family) and 9 copies of IS1081 (IS256 family), two  
112 IS families in the MTBC that are known to vary in copy number.

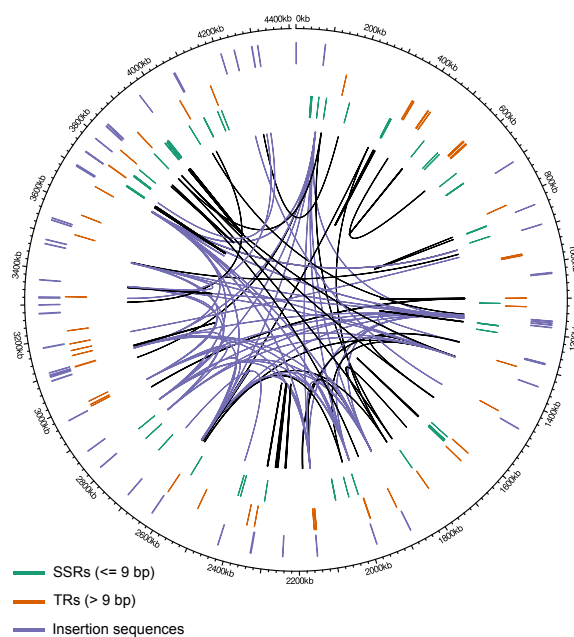


Figure 1: Repeat landscape in P034-N1426. Annotated short sequence repeats (SSRs, <=9 bp), tandem repeats (>9 bp), and insertion sequences are shown in the three outer tracks. Homopolymers are not shown because of their large number. The links inside the circle show pairs of homologous sequences of at least 90% identity and 50 bp length. Links that connect copies of insertion sequences are shown in purple, others in black.

113 As a second approach to characterize the repeat landscape, we identified pairs of homologous  
114 sequences of at least 50 bp and 90% identity across the genome. Sequence homology is the  
115 substrate for homologous recombination and thus informative about where in the genome we

116 might expect recombination-associated mutations. The thresholds are somewhat arbitrary as  
117 the minimal efficient processing segment (MEPS) for homologous recombination is not known  
118 in the MTBC; this point is addressed in the discussion. Excluding homology pairs within TRs  
119 to avoid redundancy in repeat annotation, we identified 136 pairs of homologous sequences,  
120 making up a total of 93,685 bp or 2.1% of the 4.4 Mb genome (Fig. 1).

121 To better understand which genetic elements share sequence homology, we intersected the  
122 homology segments with the gene and IS annotations. The repeat landscape is dominated by 45  
123 and 15 homology pairs of highly similar IS6110 and IS1081 copies, respectively (Fig. 1). Members  
124 of the ESX, PE, and PPE gene families constitute a second prominent feature: 9 ESX genes (*esxI*,  
125 *esxJ*, *esxK*, *esxL*, *esxN*, *esxO*, *esxP*, *esxV*, *esxW*), 8 PPE genes (*ppe18*, *ppe19*, *ppe34*, *ppe38*, *ppe46*,  
126 *ppe57*, *ppe59*, *ppe60*), and 7 PE\_PGRS genes (*pe\_pgrs17*, *pe\_pgrs18*, *pe\_pgrs19*, *pe\_pgrs20*, *pe\_pgrs27*,  
127 *pe\_pgrs28*, *pe\_pgrs45*) are part of homology pairs (Fig. 2B). Most of the remaining sequence  
128 homology is found between pairs of genes of unknown function, designated by gene names  
129 beginning with "Rv" in Figure 2B.

### 130 **Types, frequencies, and genomic context of the 110 identified variants**

131 Equipped with a basic understanding of what types of repeats are located where in the genome,  
132 we constructed a pangenome graph from the 16 assemblies and identified "bubbles" in the graph.  
133 110 variants at 109 sites were identified (Fig. 3A), including 75 single nucleotide polymorphisms,  
134 11 multinucleotide polymorphisms (MNPs, simultaneous changes of two or three base pairs), 17  
135 deletions and 7 insertions. Deletion and insertion lengths range from 1 to 7,346 bp, the majority  
136 being indels smaller than 10 bp (Fig. 3B). Regarding the frequency of the variants, 88 of 110 are  
137 singletons, that is, are present in only one of the sampled strains, 14 are shared between up to  
138 five strains, and eight are present in all but P034-N1426, the strain that diverged early from the  
139 rest of the sample (Fig. 3C, Supplemental Fig. S1).

140 To identify repeat-associated mutations, we intersected the variant sites with our repeat  
141 annotation. 62 of the 110 variants are associated with repeats, while 48 occur in non-repetitive  
142 genic or intergenic regions (Fig. 3D). For most repeat-associated variants, the annotation directly  
143 suggests an underlying mechanism: of the seven insertions and deletions larger than 50 bp, five  
144 locate to tandem repeats, with the length of the variants corresponding to a multiple of the  
145 repeat periodicity. Four small indels are located in SSRs and homopolymers, suggesting strand

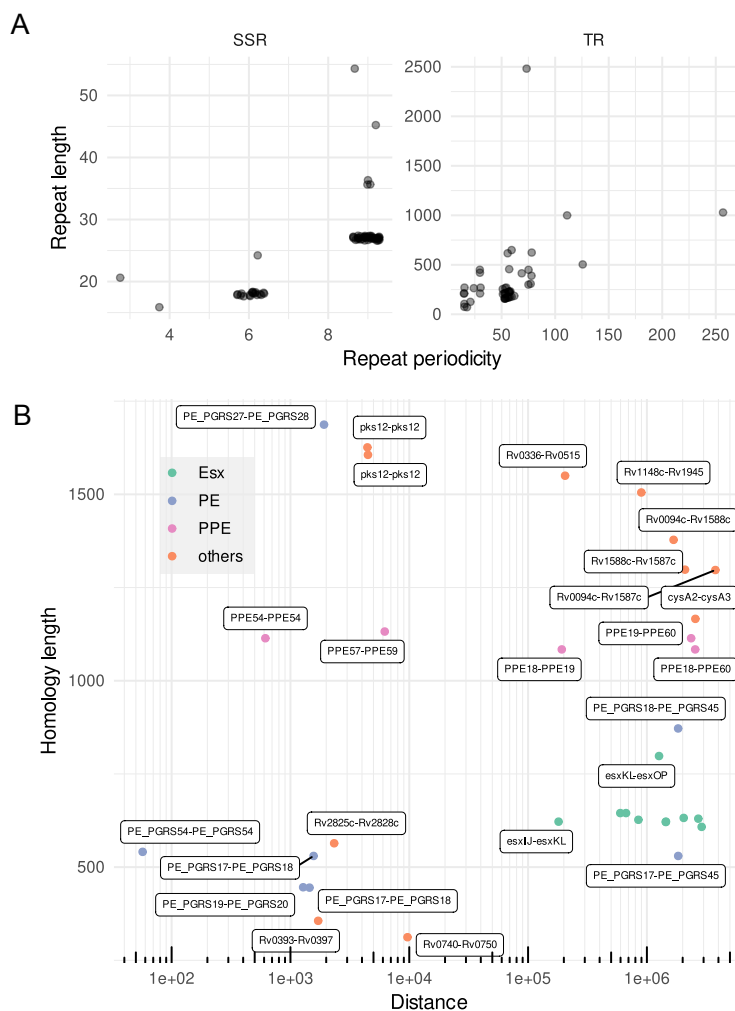


Figure 2: Characteristics of SSRs, TRs, and homology pairs. A) Periodicity and total length of short sequence repeats (SSRs) and tandem repeats (TRs). B) Gene pairs sharing substantial homology (over at least 20% of their length). The x-axis shows the distance between the pairs, the y-axis the length of the homology segment. Colors indicate the three main repetitive gene families in the MTBC. Not all *esx* pairs are labeled due to lack of space.

146 slippage as underlying mechanism. The most striking pattern, however, is the large number of  
 147 variants that intersect with homology segments. 49 variants (45%), including all MNPs, locate  
 148 to homology segments. This is a more than 20-fold over-representation, considering that these  
 149 segments make up 2.1% of the genome.

150 **Gene conversion between PE/PPE genes accounts for more than a third of all variants**

151 A closer look at the variants occurring in homology segments shows that they occur as dense  
 152 clusters of variants in single strains and are located in PE/PPE genes, two multigene families that

153 are characteristic of pathogenic mycobacteria and play a role in host-pathogen interactions. 34 of  
154 the variants identified in P059-N1392 cluster in the repetitive C-terminal domain of *pe\_pgrs28*  
155 (Fig. 4). Similarly, six variants in P003-N1374 cluster in *ppe18*, and four variants in P052-N1429  
156 in *ppe19*. Since *ppe18* was annotated as a surface antigen, we further investigated the variants in  
157 this genes and found that two non-synonymous mutations affect two distinct epitope regions in  
158 the gene (Supplemental Fig. S3).

159 Given that these variant clusters occur in segments of homology, we hypothesized that they  
160 were caused by gene conversion between close paralogs. To test this, we blasted the suspected  
161 conversion tracts against the genome, expecting two exact hits in the two paralogs involved  
162 versus only a single exact hit in the source gene in strains where no gene conversion occurred.  
163 Indeed the conversion tract in *pe\_pgrs28* yields two exact matches in *pe\_pgrs27* (Fig. 4). The  
164 matching regions are separated by 286 bp, suggesting that the conversion tract is not continuous  
165 but interrupted by a stretch of the target gene. The suspected conversion tracts in *ppe18* and  
166 *ppe19* yield exact matches in *ppe19* and *ppe18*, respectively – gene conversion has worked both  
167 ways in different strains between these two genes, which are 190 kb apart.

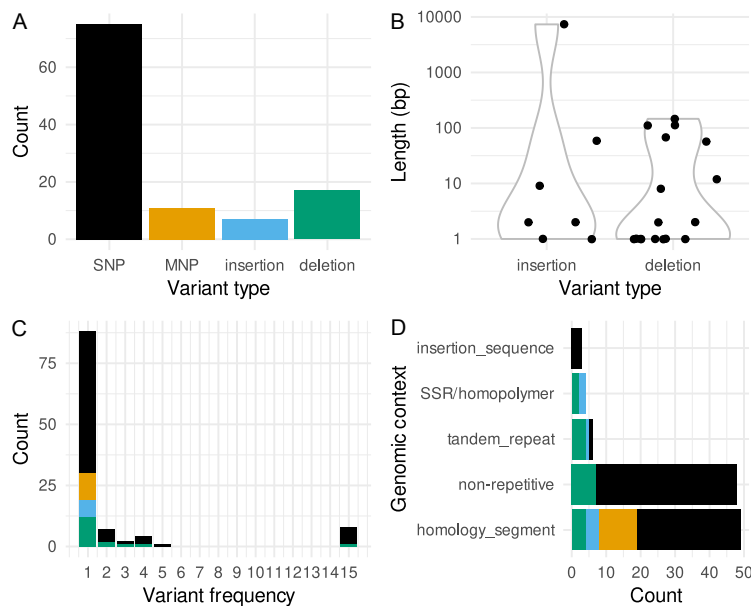


Figure 3: Types of DNA polymorphisms in the 16 sampled strains. A) Counts of SNPs, MNPs, insertions and deletions. B) Length of insertions and deletions, with a log-scale on the y axis. C) Frequency of the variants in the 16 strains. D) Intersection of the variant sites with the annotated repeats.



## 168 **Birth of a new PPE gene through homologous recombination**

169 As noted above, for one strain (P001-N1377) the Flye assembly step resulted in a single non-  
170 circularized contig. To understand why the assembly failed in this strain, we aligned the reads of  
171 P001-N1377 against a close relative where this region posed no problems. Double coverage and  
172 clipped reads, where the clipped parts feed back into the repeat on the opposite site, suggest  
173 that this region is duplicated in P001-N1377 (Supplemental Fig. S4). A comparatively long read  
174 would be required to resolve this  $2 * 7,346 = 14,692$  bp region, given an N50 read length of 6,839  
175 bp for this strain. Indeed there is one read of 18,797 bp that spans the region (Supplemental Fig.  
176 S4). The short overlap on the 3' side (140 bp) might explain why Flye failed to close the gap.

177 The duplication occurred in a region of the genome that contains multiple PE/PPE genes  
178 and insertion sequences (Fig. 5A) and where nested repeats testify to past duplication events  
179 (Fig. 5D). According to our de novo annotation, the duplication contains eight reading frames,  
180 five of them coding for PE/PPE genes. A comparison with the H37Rv reference annotation and  
181 our insertion sequence annotation shows that three CDS are part of a IS21 insertion sequence,  
182 while the unnamed gene in front of the IS is *ppe58* (Fig. 5A). Three of the five annotated PE/PPE  
183 CDS are not annotated in H37Rv; their sequences are similar to the closely related *ppe57*, *ppe58*,  
184 and *ppe59* (Fig. 4D), suggesting that these are leftovers from a previous duplication event.

185 The gene models (Fig. 4A) proved not very helpful when trying to understand the  
186 convoluted graph for the duplicated region and how the duplication might have originated.  
187 More informative were the segments of homology (see above). A homology pair is present in the  
188 region, HS-156A and HS-156B (Fig. 5C); the two segments comprise *ppe57* and *ppe59*, respectively,  
189 and an additional stretch at the 5' side of these genes. A blast search of these segments against  
190 the graph reveals a third segment at the very center of the duplication, highlighted by the red  
191 lines in Figure 4D, which is a recombinant between the two parental segments (Fig. 6A). The  
192 location of the recombination breakpoint is fuzzy but has to be located before or within the first  
193 70 bp the new, identical copy of *ppe57*, as *ppe57* and *ppe59* are identical for the first 70 bp but then  
194 differ substantially. The presence of a third homology segment at the center of the duplication  
195 suggests a mechanism of duplication through homologous recombination (Fig. 6B).

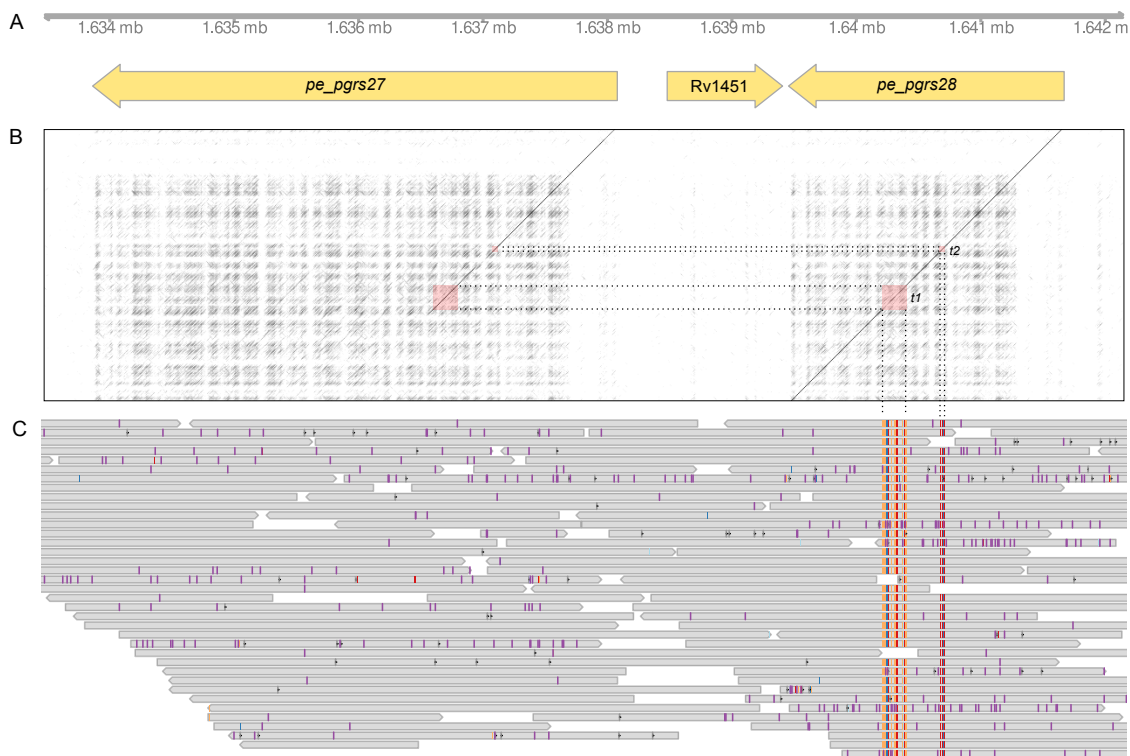


Figure 4: Gene conversion event between *pe\_pgrs27* and *pe\_pgrs28* in the strain P059-N1392. A) Genes annotated in the affected region. B) A sequence dotplot of the whole region in P059-N1392 on the x-axis and the sequence of the target gene *pe\_pgrs28* on the y-axis. The red squares show the two conversion tracts t1 and t2. C) PacBio reads of P059-N1392 aligned against P034-N1426, highlighting the 34 variants in the two conversion tracts.

196

## Discussion

197 The promise of third generation sequencing technologies is to resolve repeats and thus enable  
198 a systematic study of repeat-associated genetic variation. In this study we show that virtually  
199 base-perfect MTBC genomes can be assembled from PacBio HiFi reads. These provide the  
200 basis to explore all genetic variation present in a set of genomes, in particular in the repeatome.  
201 In the following, we discuss homologous recombination as a mechanism underlying genetic  
202 diversity and argue that this mechanism is frequent and prone to have fitness consequences for  
203 the bacteria. We then situate our study among other recent attempts to obtain complete genomes  
204 from long reads and discuss read size as the limiting factor in the presence of duplications.  
205 Finally, we consider the use of complete genomes in genome-scale inference and point out  
206 caveats when including repeat-associated mutations for estimating phylogenies and selection  
207 scans.

## 208 **Recombination re-emerges as an important mechanism in the MTBC**

209 Our analysis of a sample of 16 strains in a single transmission cluster does not lend itself  
210 to generalizations. Still, the observation of three gene conversion events, a duplication, and  
211 several insertions and deletions in tandem repeats in a sample of closely related strains suggests  
212 that homologous recombination operates frequently. While point mutations remain the most  
213 frequent events (Fig. 3A), recombination can affect large regions and contribute disproportionately  
214 to genetic variation: three gene conversion events between members of the PE/PPE families  
215 (*pe\_pgrs27/28*, *ppe19* and *ppe18*) account for 44 of the 110 variants identified, while the largest  
216 variant is a duplication of 7.3 kb that arose from recombination between a homology pair  
217 encompassing *ppe57* and *ppe59*.

218 Evidence for gene conversion and homologous recombination was presented in a series  
219 of studies in the early 2000s, before the topic largely disappeared with the advent short read  
220 sequencing. Karboul et al. (2006) first discussed gene conversion as a potentially important  
221 mechanism in the MTBC and described recurrent conversion between the adjacent *pe\_pgrs17* and  
222 *pe\_pgrs18*. Evidence for gene conversion between members of the PPE (McEvoy, Van Helden,  
223 et al. 2009) and ESX (Uplekar et al. 2011) families followed, as well as for *pe\_pgrs27* and *pe\_pgrs28*  
224 (Delogu et al. 2008), the gene most strongly affected by gene conversion in our study, and *ppe18*  
225 and *ppe60* (McEvoy, Cloete, et al. 2012).

226 These and similar studies were based on PCR amplification of few genes. With short read  
227 sequencing, the study of repeats and homologous recombination was marginalized. Several  
228 studies still investigated “recombination”, but used this general term to denote horizontal gene  
229 transfer (HGT) rather than intrachromosomal recombination (e.g. Godfroid et al. 2018; Chiner-  
230 Oms et al. 2019). The apparent absence of HGT in the MTBC is expressed in the paradigm “TB  
231 does not recombine”. As repeats are now coming into focus again, it is important to recall that  
232 homologous recombination is a many-sided fundamental mechanism involved not only in HGT,  
233 but also in replication and repair, and its underlying pathways are present in the MTBC (Gupta  
234 et al. 2011).

235 Some expectations regarding the occurrence of repeat-associated mutations, and particularly  
236 homologous recombination, can be formed by considering the distribution of repeats in a single  
237 genome (Fig. 1). The rate of recombination, and thus gene conversion, decreases exponentially  
238 as sequences diverge (Shen and Huang 1986). Indeed, to our knowledge all examples of gene

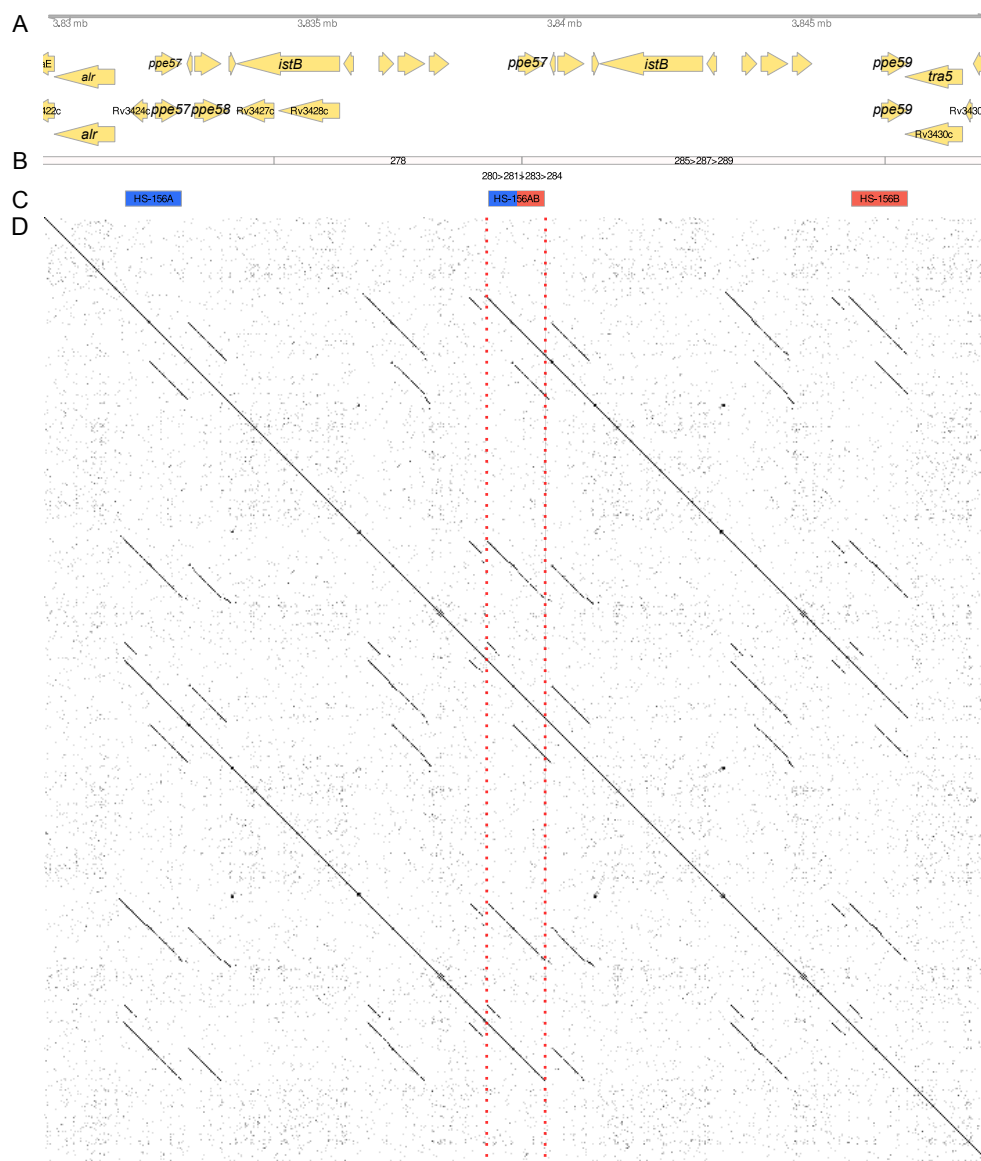


Figure 5: The duplicated region in P001-N1377. A) Genes annotated in the duplicated region, above from our *de novo* annotation, below from a liftover from the H37Rv reference annotation. B) Path in the genome graph corresponding to the duplicated region. C) Homology segments. D) Sequence dotplot of the duplicated region against itself. The red vertical lines highlight the center of the duplication where the duplicates overlap and the recombinant homology segment is located.

239 conversion in the MTBC involve closely related genes. Even in the large PE and PPE gene  
240 families, with 99 and 68 members in the reference strain H37Rv, close homology is restricted  
241 to relatively few pairs and triplets of genes, and to some larger repeats within single genes  
242 (*pe\_pgrs24*; Fig. 2B). These numbers would be higher with more permissive thresholds for

243 homology search; e.g. *ppe58* would appear as a paralog of *ppe57* and *ppe59* (Fig. 4 Fishbein et al.  
244 2015).

245 While the number of genes involved in recombination might be small relative to the size of  
246 these gene families, at least some of the variants in these genes will have phenotypic consequences  
247 and affect traits of primary interest in MTBC research, including antimicrobial resistance, immune  
248 response, and virulence. The three main repetitive gene families in the MTBC (PE, PPE, ESX) play  
249 important roles in host-pathogen interactions; they code for secretion systems, surface receptors  
250 or secreted proteins that interact with the human immune system and were instrumental in the  
251 evolutionary transition to an obligately pathogenic lifestyle (Gey Van Pittius et al. 2006). Even  
252 in our small sample of closely related strains we found two genes affected by gene conversion  
253 and duplication, *ppe18* and *ppe57*, for which there is experimental evidence for a role in immune  
254 subversion (Nair et al. 2009; Xu et al. 2015). Gene conversion has been shown to be a cause  
255 of antigenic variation in different prokaryote species (reviewed by Santoyo and Romero 2005).  
256 Given that both *ppe18* and *ppe57* are part of TB vaccines in active development (Guo et al. 2023),  
257 a better understanding of their molecular evolution could be of practical relevance.

258 One type of variant we did not observe are insertion sequence polymorphisms, even though  
259 copies of IS6110 and IS1081 are the most conspicuous features of the repeat landscape in the  
260 studied genomes. IS6110 is an active element that varies in copy number in the MTBC, from  
261 zero to more than 25 (McEvoy, Falmer, et al. 2007)). Copies of this element flank large inversions  
262 and deletions (Roychowdhury et al. 2015), suggesting that they are involved in homologous  
263 recombination. Because IS copies are dispersed through the genome, recombination between  
264 copies with other outcomes than gene conversion is likely to be disruptive and would be rarely  
265 observed. It will be interesting to test whether high copy numbers of IS6110 lead to more  
266 deletions and inversions or, on the opposite, high copy numbers are made possible by mutations  
267 that decrease the recombination rate and thus the rate of deleterious large-scale structural  
268 variants.

## 269 **Duplications are the new dark matter**

270 Our results show that assemblies from PacBio CCS reads are complete and essentially error-free.  
271 One single of the 16 genomes could not be closed in the primary assembly step, and only five  
272 sites in a total of 70,493,034 assembled bases were found to be ambiguous or inconsistent with

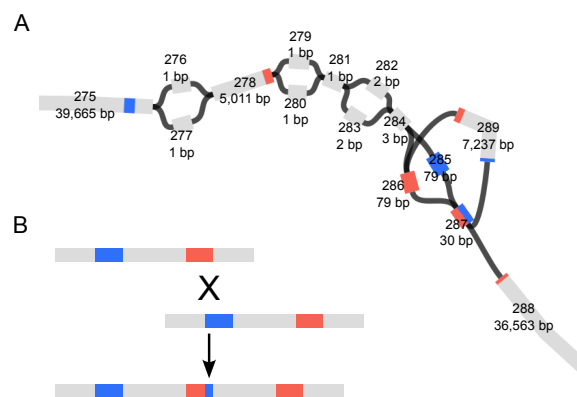


Figure 6: A) Graph representation of the duplication. Nodes are numbered and their length in bp is indicated. The colors indicate sequence similarity as determined through a blast search of the duplicated and adjacent parts. B) A model of homologous recombination underlying the duplication.

273 the underlying reads. Furthermore, surprisingly low sequencing depths are required to obtain  
274 complete MTBC genomes: a depth of 16x was sufficient to obtain a closed genome with only  
275 one single ambiguous position. These assemblies are thus an appropriate point of departure to  
276 study the molecular evolution of repeats and dynamic gene families.

277 The accuracy of CCS reads allowed us to forego many of the methodological complexities  
278 of previous approaches to long-read assembly. Noisy long reads necessitate high sequencing  
279 depths and hybrid approaches that combine long and short reads and multiple assembly and  
280 polishing steps (e.g. De Maio et al. 2019, Wick et al. 2023). This results in complicated pipelines  
281 and directs attention away from biological to technical questions. Our results show that a simple  
282 approach that combines CCS reads and the Flye assembly algorithm performs well, resulting in  
283 assemblies that can be trusted.

284 The principal limitation of CCS reads, at present, is their size and thus the ability to bridge  
285 duplications and amplifications (Tvedte et al. 2021). Duplications, thus, are the new dark matter  
286 and the reason why the assembly problem is not solved for the MTBC. The proportion of long  
287 reads might be increased through improved DNA extraction and library preparation (Wick et al.  
288 2023). But even with longer reads, whether an assembly can be closed also depends on the  
289 frequency and length of duplications and amplifications.

290 In the strains analyzed in this study, we stumbled upon a duplication of 7,346 bp that was  
291 bridged by a single read. Compared to other duplications that have been described in the MTBC,  
292 this is not particularly large, but its assembly still requires reads longer than 14,692 bp. Well

293 described duplications are the 30 kb DU1 and 36 kb DU2 duplications in the BCG vaccine  
294 clade (Brosch et al. 2000), or the massive duplication of 350 kb that includes the DosR regulon  
295 (Domenech et al. 2010) and has appeared repeatedly in lineages two and four (Weiner et al.  
296 2012). More recent examples are a 38- to 60-fold amplification of *esxR/esxS* and flanking PE/PPE  
297 genes in H37Rv mutant strains where the ESX-3 excretion system had been deleted (Wang et al.  
298 2022); and a 120 kb duplication that evolved twice independently during experimental evolution  
299 (Smith et al. 2022). These genomes would not have been completed with the approach presented  
300 here. The most extensive investigation of amplifications in the MTBC so far, based on short-read  
301 coverage, suggests that amplifications are frequent but restricted to few genomic regions: 590  
302 amplifications were found in 1,000 diverse MTBC genomes, the large majority of them in 24  
303 hotspots regions (Abrahams et al. 2022).

### 304 **Making use of the repeatome**

305 Of more immediate practical concern than the fitness effects of repeat-associated mutations  
306 is the use of the repeatome in genome-scale inference. In an epidemiological context, a key  
307 consideration is how repeats can be included to increase the genetic resolution for the inference  
308 of transmission chains (e.g. Modlin et al. 2021, Marin et al. 2021). It is evident that using  
309 complete assemblies and all types of polymorphisms increases the available information; it is  
310 less obvious exactly how this information should be used. Much of phylogenetic inference  
311 and selection scans is based on substitution models that apply to sequentially introduced point  
312 mutations (Yang 2014). As shown above (Fig. 4), gene conversion can introduce many variants  
313 in a single event.

314 While a gene conversion event can be informative about tree topology when it is shared  
315 between strains, modeling it through a substitution model will result in exaggerated branch  
316 lengths and biased time estimates. For the same reason, variants due to gene conversion  
317 should be excluded in clustering analyses based on SNP thresholds. Multiple variant types and  
318 mutational models could be used simultaneously in a Bayesian setting, as has been shown with  
319 SNPs and short indels (Redelings and Suchard 2007).

320 Caveats also apply to inferring positive selection on genes affected by gene conversion. On  
321 the one hand, codon models as well as the popular  $d_N/d_S$  statistic are based on the same  
322 substitution models as phylogenetic inference. Applying these methods to close paralogs can

323 reveal false signatures of positive selection (Casola and Hahn 2009), and we suspect that previous  
324 reports of positive selection acting on a large number of PE/PPE genes (e.g. Zhang et al. 2011  
325 Namouchi et al. 2013, Phelan et al. 2016) might have been exaggerated. Also homoplasmy, a  
326 second popular signature of selection, has to be interpreted with care since gene conversion  
327 mimics convergent evolution (Green et al. 2023).

## 328 **Methods**

### 329 **Samples and genome sequencing**

330 Strains for PacBio sequencing were selected to represent the “Bernese outbreak”, a small  
331 transmission cluster belonging to lineage 4 and including 68 patients in the city of Bern  
332 (Switzerland), sampled between 1988 and 2011 (Figure S1, Supplementary Table S1; Genewein et  
333 al. 1993, Stucki et al. 2015). We first estimated a phylogenetic tree using the previously published  
334 short reads (NCBI bioproject PRJEB5925). A SNP alignment was created as previously described  
335 (Gygli et al. 2021), ignoring variants in known resistance genes and repetitive regions. The  
336 resulting alignment of 142 variable positions was used to estimate a phylogeny with raxml-ng v.  
337 1.2.1 (Kozlov et al. 2019), using the GTR substitution model and 100 bootstrap replicates. Branch  
338 lengths were rescaled to account for non-variable positions.

339 16 strains from the different subclusters were selected for sequencing. For two patients (P022,  
340 P028), different isolates than the ones used by Stucki et al. were accidentally selected. Since  
341 the aim of this study is not epidemiological inference, we still used them and unambiguously  
342 labelled the assemblies with the patient number and the isolate name separated by a dash.

343 For DNA extraction, strains were grown for two to three weeks on 7H11 plates inoculated  
344 from 7H9 liquid wake-up cultures. DNA was extracted from the plates using a CTAB method  
345 (Van Embden et al. 1993) that includes RNase treatment and a purification step using magnetic  
346 beads. SMRTbell libraries were prepared at the Genomics Core Facility jointly ran by ETH Zurich  
347 and the University of Basel and sent for sequencing on a Pacific Biosciences Sequel II System at  
348 the Lausanne Genomic Technologies Facility.



## 349 **Genome assembly and variant calling**

350 Read quality statistics were obtained with LongQC (Fukasawa et al. 2020). We used Flye v. 2.8.1  
351 (Kolmogorov et al. 2019) to assemble genomes from circular consensus sequences (CCS) and  
352 reoriented the assemblies to start with *dnaA* using Circlator v. 1.5.5 (Hunt et al. 2015). As *M.*  
353 *tuberculosis* is not known to contain plasmids, assembly is expected to result in a single circular  
354 chromosome.

355 A pangenome graph representation of the assembled single-contig genomes was build  
356 using PanGenome Graph Builder (pvgb) v. 077830d (Garrison, Guarracino, et al. 2023); for  
357 an evaluation of pvgb in bacteria see Yang et al. 2023). Percentage identity (-p) was set to 99,  
358 segment length (-s) to 5k, according to the high similarity expected in strains from a single  
359 transmission cluster and the maximum length of repeats in the genome. An arbitrary strain,  
360 P034-N1426, was used as a positional reference for outputting variants with *vg deconstruct* and  
361 for the annotation of repeats (see below).

362 To classify structural variants as insertions or deletions, we estimated the ancestral states  
363 of the variants by comparint them to a closely related strain that was also sequenced for the  
364 present study: N1015, a strain from Nepal assigned to sublineage 4.4.1.2, while the Bernese  
365 strains belong to sublineage 4.4.1.1. The minor allele in the outbreak sample was assumed to be  
366 the derived state if it differed from the outgroup allele; if the minor equalled the outgroup allele,  
367 the major allele was assumed as derived.

## 368 **Assembly validation and curation**

369 To evaluate the accuracy of our assemblies, we looked for inconsistencies between reads and  
370 assemblies by aligning the long reads back against the assemblies and calling variants. If the  
371 assembly is accurate, no variants should be found, while variants identified this way indicate  
372 errors during assembly, circularization, or the presence of true genetic variation in the culture. We  
373 used minimap2 2.24-r1122 (Li 2018) to align the reads and called variants with freebayes v. 1.3.4  
374 (Garrison and Marth 2012), setting ploidy to 1. Read-assembly inconsistencies were scrutinized  
375 with the Integrative Genomics Viewer (IGV, Robinson et al. 2011) and those unequivocally  
376 identified as assembly errors were curated using the pysam library of Biopython. The curated  
377 assemblies where validated through a second round of variant calling.

## 378 **Gene and repeat annotation**

379 Genes were annotated *de novo* in all assemblies with bakta v.1.8.2 (Schwengers et al. 2021). To  
380 compare gene models, we also lifted over the H37Rv reference annotation (ASM19595v2) to  
381 P034-N1426 with liftoff v1.6.3 (Shumate and Salzberg 2021). To further characterize the repeat  
382 context of the variants, we annotated different types of repeats in the positional reference strain  
383 P034-N1426: insertion sequences with ISEScan v.1.7.2.3 (Xie and Tang 2017), short sequence  
384 repeats ( $\leq 9$  bp) with kmer-ssr v. 0.8 (Pickett et al. 2017), tandem repeats ( $> 9$  bp) with SPADE  
385 v. 1.0.0 (Mori et al. 2019), and homopolymers of at least 5 bp using our own script (see GitLab  
386 link above). To investigate these repeats in other assemblies, annotations were lifted over using  
387 *odgi position* on the pangenome graph. The variant positions resulting from *vg deconstruct* and  
388 the different annotations were intersected with bedtools v2.30.0 (Quinlan and Hall 2010).

389 To identify pairs of sequence homology across the genome, we used *nucmer* ( $-\text{maxmatch}$   
390  $-\text{nosimplify}$ ) and *show-coords* from the mummer4 tool (Marçais et al. 2018) and removed self-  
391 and overlapping hits. We further filtered out pairs with less than 90% sequence identity and an  
392 alignment length smaller than 50 bp, as well as pairs that overlapped with annotated tandem  
393 repeats. The locations of homology segments were intersected with the *de novo* gene and  
394 IS annotations, using bedtools v2.30.0 (Quinlan and Hall 2010), to identify genetic elements  
395 potentially involved in recombination.

## 396 **Identification of gene conversion tracts**

397 To test whether a cluster of variants identified in the same strain is due to gene conversion,  
398 we extracted the suspected conversion tracts from the respective genomes using *odgi* v0.8.3-  
399 57-gfdbdb4d23 (Garrison, Guarracino, et al. 2023) and *samtools* v.1.18 (Danecek et al. 2021). We  
400 then used *blastn* v. 2.12.0 (Camacho et al. 2009) to align the sequences against a) the genome in  
401 which the cluster was detected, b) a genome in the ancestral state. In case of gene conversion,  
402 we expect two matches in a, in the source and the target gene, and a single match in b, i.e. only  
403 in the source gene. Visualization of the gene conversion and duplication events was done using  
404 the R package *gviz* (Hahne and Ivanek 2016) and *dotter* (Sonnhammer and R Durbin 1995).

405 Epitopes in PPE18 were downloaded from the Immune Epitope Database ([iedb.org](http://iedb.org), accessed  
406 24.1.2024). *Miniprot* v. (Li 2023) was used to infer the location of the epitope peptides within the  
407 nucleotide sequences of the genes.

408

## Data access

409 Raw reads and assembled genomes were deposited on the European Nucleotide Archive  
410 (PRJEB73759). The assembly pipeline is available as a Snakemake workflow on <http://git.scicore.unibas.ch/TBRU/PacbioSnake>.

412

## Competing interest statement

413 The authors declare no competing interests.

414

## Acknowledgments

415 We wish to thank the members of the Gagneux group for the helpful discussions and the team  
416 of the sciCORE Center for Scientific Computing at the University of Basel for access to their  
417 computing cluster. This work was funded through grants from the European Research Council,  
418 grant number 883582, and the Swiss National Science Foundation, grant numbers 310030\_188888  
419 and CRSII5\_177163.

420

## References

- 421 Abrahams JS, Weigand MR, Ring N, MacArthur I, Etty J, Peng S, Williams MM, Bready B,  
422 Catalano AP, Davis JR, et al. 2022. Towards comprehensive understanding of bacterial genetic  
423 diversity: large-scale amplifications in *Bordetella pertussis* and *Mycobacterium tuberculosis*.  
424 *Microbial Genomics*. **8**: 000761.
- 425 Achaz G. 2002. Origin and fate of repeats in bacteria. *Nucleic Acids Research*. **30**: 2987–2994.
- 426 Achtman M. 2012. Insights from genomic comparisons of genetically monomorphic bacterial  
427 pathogens. *Philosophical Transactions of the Royal Society B: Biological Sciences*. **367**: 860–867.
- 428 Armstrong J et al. 2020. Progressive Cactus is a multiple-genome aligner for the thousand-genome  
429 era. *Nature*. **587**: 246–251.
- 430 Brosch R, Gordon SV, Buchrieser C, Pym AS, Garnier T, and Cole ST. 2000. Comparative  
431 genomics uncovers large tandem chromosomal duplications in *Mycobacterium bovis* BCG  
432 Pasteur. *Yeast*. **1**: 111–123.

- 433 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, and Madden TL. 2009.  
434 BLAST+: architecture and applications. *BMC Bioinformatics*. **10**: 421.
- 435 Casola C and Hahn MW. 2009. Gene conversion among paralogs results in moderate false  
436 detection of positive selection using likelihood methods. *Journal of Molecular Evolution*. **68**:  
437 679–687.
- 438 Chiner-Oms Á, Sánchez-Busó L, Corander J, Gagneux S, Harris SR, Young D, González-Candelas  
439 F, and Comas I. 2019. Genomic determinants of speciation and spread of the *Mycobacterium*  
440 *tuberculosis* complex. *Science Advances*. **5**: eaaw3307.
- 441 Danecek P et al. 2021. Twelve years of SAMtools and BCFtools. *GigaScience*. **10**: giab008.
- 442 Darmon E and Leach DRF. 2014. Bacterial genome instability. *Microbiology and Molecular Biology*  
443 *Reviews*. **78**: 1–39.
- 444 De Maio N et al. 2019. Comparison of long-read sequencing technologies in the hybrid assembly  
445 of complex bacterial genomes. *Microbial Genomics*. **5**:
- 446 Delogu G, Cole ST, and Brosch R. 2008. The PE and PPE protein families of *Mycobacterium*  
447 *tuberculosis*. *Handbook of tuberculosis*. 131–150.
- 448 Domenech P, Kolly GS, Leon-Solis L, Fallow A, and Reed MB. 2010. Massive gene duplication  
449 event among clinical isolates of the *Mycobacterium tuberculosis* W/Beijing Family. *Journal of*  
450 *Bacteriology*. **192**: 4562–4570.
- 451 Fishbein S, Van Wyk N, Warren R, and Sampson S. 2015. Phylogeny to function: PE/PPE protein  
452 evolution and impact on *Mycobacterium tuberculosis* pathogenicity. *Molecular Microbiology*. **96**:  
453 901–916.
- 454 Fukasawa Y, Ermini L, Wang H, Carty K, and Cheung MS. 2020. LongQC: A quality control tool  
455 for third generation sequencing long read data. *G3 Genes | Genomes | Genetics*. **10**: 1193–1196.
- 456 Gagneux S. 2018. Ecology and evolution of *Mycobacterium tuberculosis*. *Nature Reviews Microbiology*.  
457 **16**: 202–213.
- 458 Garrison E, Guarracino A, Heumos S, Villani F, Bao Z, Tattini L, Haggmann J, Vorbrugg S,  
459 Marco-Sola S, Kubica C, et al. 2023. Building pangenome graphs, pp. 2023–04.
- 460 Garrison E and Marth G. 2012. Haplotype-based variant detection from short-read sequencing.  
461 *arXiv preprint arXiv:1207.3907*.

- 462 Genewein A, Telenti A, Bernasconi C, Schopfer K, Bodmer T, Mordasini C, Weiss S, Maurer AM,  
463 and Rieder H. 1993. Molecular approach to identifying route of transmission of tuberculosis  
464 in the community. *The Lancet*. **342**: 841–844.
- 465 Gey Van Pittius NC, Sampson SL, Lee H, Kim Y, Van Helden PD, and Warren RM. 2006.  
466 Evolution and expansion of the *Mycobacterium tuberculosis* PE and PPE multigene families  
467 and their association with the duplication of the ESAT-6 (esx) gene cluster regions. *BMC*  
468 *Evolutionary Biology*. **6**: 95.
- 469 Godfroid M, Dagan T, and Kupczok A. 2018. Recombination signal in *Mycobacterium tuberculosis*  
470 stems from reference-guided assemblies and alignment artefacts. *Genome Biology and Evolution*.  
471 **10**: 1920–1926.
- 472 Green AG, Vargas R, Marin MG, Freschi L, Xie J, and Farhat MR. 2023. Analysis of genome-  
473 wide mutational dependence in naturally evolving *Mycobacterium tuberculosis* populations.  
474 *Molecular Biology and Evolution*. **40**: msad131.
- 475 Guo F, Wei J, Song Y, Li B, Wang X, Wang H, and Xu T. 2023. Immunological effects of  
476 the PE/PPE family proteins of *Mycobacterium tuberculosis* and related vaccines. *Frontiers in*  
477 *Immunology*. **14**: 1255920.
- 478 Gupta R, Barkan D, Redelman-Sidi G, Shuman S, and Glickman MS. 2011. Mycobacteria exploit  
479 three genetically distinct DNA double-strand break repair pathways. *Molecular Microbiology*.  
480 **79**: 316–330.
- 481 Gygli SM et al. 2021. Prisons as ecological drivers of fitness-compensated multidrug-resistant  
482 *Mycobacterium tuberculosis*. *Nature Medicine*. **27**: 1171–1177.
- 483 Hahne F and Ivanek R. 2016. Visualizing genomic data using Gviz and bioconductor. *Statistical*  
484 *genomics: methods and protocols*. 335–351.
- 485 Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA, and Harris SR. 2015. Circlator: automated  
486 circularization of genome assemblies using long sequencing reads. *Genome Biology*. **16**: 1–10.
- 487 Karboul A et al. 2006. Insights into the evolutionary history of tubercle bacilli as disclosed by  
488 genetic rearrangements within a PE\_PGRS duplicated gene pair. *BMC Evolutionary Biology*. **6**:  
489 107.
- 490 Kolmogorov M, Yuan J, Lin Y, and Pevzner PA. 2019. Assembly of long, error-prone reads using  
491 repeat graphs. *Nature Biotechnology*. **37**: 540–546.

- 492 Kozlov AM, Darriba D, Flouri T, Morel B, and Stamatakis A. 2019. RAxML-NG: a fast, scalable  
493 and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics*. **35**:  
494 4453–4455.
- 495 Kühnert D, Coscolla M, Brites D, Stucki D, Metcalfe J, Fenner L, Gagneux S, and Stadler T. 2018.  
496 Tuberculosis outbreak investigation using phylodynamic analysis. *Epidemics*. **25**: 47–53.
- 497 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. **34**: 3094–3100.
- 498 Li H. 2023. Protein-to-genome alignment with miniprot. *Bioinformatics*. **39**: btad014.
- 499 Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, and Zimin A. 2018. MUMmer4: A  
500 fast and versatile genome alignment system. *PLoS Computational Biology*. **14**: e1005944.
- 501 Marin M, Vargas Jr R, Harris M, Jeffrey B, Epperson LE, Durbin D, Strong M, Salfinger M,  
502 Iqbal Z, Akhundova I, et al. 2021. Genomic sequence characteristics and the empiric accuracy  
503 of short-read sequencing. *bioRxiv*. 2021–04.
- 504 McEvoy CR, Cloete R, Müller B, Schürch AC, Van Helden PD, Gagneux S, Warren RM, and  
505 Gey Van Pittius NC. 2012. Comparative analysis of *Mycobacterium tuberculosis* PE and PPE  
506 genes reveals high sequence variation and an apparent absence of selective constraints. *PLoS*  
507 *ONE*. **7**: e30593.
- 508 McEvoy CR, Falmer AA, Van Pittius NCG, Victor TC, Van Helden PD, and Warren RM. 2007.  
509 The role of IS6110 in the evolution of *Mycobacterium tuberculosis*. *Tuberculosis*. **87**: 393–404.
- 510 McEvoy CR, Van Helden PD, Warren RM, and Van Pittius N. 2009. Evidence for a rapid rate of  
511 molecular evolution at the hypervariable and immunogenic *Mycobacterium tuberculosis* PPE38  
512 gene region. *BMC Evolutionary Biology*. **9**: 237.
- 513 Modlin SJ, Robinhold C, Morrissey C, Mitchell SN, Ramirez-Busby SM, Shmaya T, and Valafar F.  
514 2021. Exact mapping of Illumina blind spots in the *Mycobacterium tuberculosis* genome  
515 reveals platform-wide and workflow-specific biases. *Microbial Genomics*. **7**:
- 516 Mori H, Evans-Yamamoto D, Ishiguro S, Tomita M, and Yachie N. 2019. Fast and global detection  
517 of periodic sequence repeats in large genomic resources. *Nucleic Acids Research*. **47**: e8–e8.
- 518 Nair S, Ramaswamy PA, Ghosh S, Joshi DC, Pathak N, Siddiqui I, Sharma P, Hasnain SE, Mande  
519 SC, and Mukhopadhyay S. 2009. The PPE18 of *Mycobacterium tuberculosis* interacts with TLR2  
520 and activates IL-10 induction in macrophage. *The Journal of Immunology*. **183**: 6269–6281.

- 521 Namouchi A, Karboul A, Fabre M, Gutierrez MC, and Mardassi H. 2013. Evolution of smooth  
522 tubercle bacilli PE and PE\_PGRS genes: evidence for a prominent role of recombination and  
523 imprint of positive selection. *PLoS ONE*. **8**: e64718.
- 524 Phelan JE et al. 2016. Recombination in pe/ppe genes contributes to genetic variation in  
525 *Mycobacterium tuberculosis* lineages. *BMC Genomics*. **17**: 151.
- 526 Pickett BD, Miller JB, and Ridge PG. 2017. Kmer-SSR: a fast and exhaustive SSR search algorithm.  
527 *Bioinformatics*. **33**: 3922–3928.
- 528 Quinlan AR and Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic  
529 features. *Bioinformatics*. **26**: 841–842.
- 530 Redelings BD and Suchard MA. 2007. Incorporating indel information into phylogeny estimation  
531 for rapidly emerging pathogens. *BMC Evolutionary Biology*. **7**: 40.
- 532 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, and Mesirov JP.  
533 2011. Integrative genomics viewer. *Nature Biotechnology*. **29**: 24–26.
- 534 Roychowdhury T, Mandal S, and Bhattacharya A. 2015. Analysis of IS 6110 insertion sites provide  
535 a glimpse into genome evolution of *Mycobacterium tuberculosis*. *Scientific Reports*. **5**: 12567.
- 536 Santoyo G and Romero D. 2005. Gene conversion and concerted evolution in bacterial genomes.  
537 *FEMS Microbiology Reviews*. **29**: 169–183.
- 538 Schmid M, Frei D, Patrignani A, Schlapbach R, Frey JE, Remus-Emsermann MN, and Ahrens CH.  
539 2018. Pushing the limits of de novo genome assembly for complex prokaryotic genomes  
540 harboring very long, near identical repeats. *Nucleic Acids Research*. **46**: 8953–8965.
- 541 Schwengers O, Jelonek L, Dieckmann MA, Beyvers S, Blom J, and Goesmann A. 2021. Bakta: rapid  
542 and standardized annotation of bacterial genomes via alignment-free sequence identification.  
543 *Microbial Genomics*. **7**:
- 544 Shen P and Huang HV. 1986. Homologous recombination in *Escherichia coli*: dependence on  
545 substrate length and homology. *Genetics*. **112**: 441–457.
- 546 Shumate A and Salzberg SL. 2021. Liftoff: accurate mapping of gene annotations. *Bioinformatics*.  
547 **37**: 1639–1643.
- 548 Smith TM, Youngblom MA, Kernien JF, Mohamed MA, Fry SS, Bohr LL, Mortimer TD, O'Neill  
549 MB, and Pepperell CS. 2022. Rapid adaptation of a complex trait during experimental  
550 evolution of *Mycobacterium tuberculosis*. *eLife*. **11**: e78454.

- 551 Sonnhammer EL and Durbin R. 1995. A dot-matrix program with dynamic threshold control  
552 suited for genomic DNA and protein sequence analysis. *Gene*. **167**: GC1–GC10.
- 553 Stritt C and Gagneux S. 2023. How do monomorphic bacteria evolve? The *Mycobacterium*  
554 *tuberculosis* complex and the awkward population genetics of extreme clonality. *Peer Community*  
555 *Journal*. **3**: e92.
- 556 Stucki D et al. 2015. Tracking a tuberculosis outbreak over 21 years: strain-specific single-  
557 nucleotide polymorphism typing combined with targeted whole-genome sequencing. *The*  
558 *Journal of Infectious Diseases*. **211**: 1306–1316.
- 559 Treangen TJ, Abraham AL, Touchon M, and Rocha EP. 2009. Genesis, effects and fates of repeats  
560 in prokaryotic genomes. *FEMS Microbiology Reviews*. **33**: 539–571.
- 561 Tvedte ES et al. 2021. Comparison of long-read sequencing technologies in interrogating bacteria  
562 and fly genomes. *G3 Genes | Genomes | Genetics*. **11**: jkab083.
- 563 Uplekar S, Heym B, Friocourt V, Rougemont J, and Cole ST. 2011. Comparative genomics of  
564 *esx* genes from clinical isolates of *Mycobacterium tuberculosis* provides evidence for gene  
565 conversion and epitope variation. *Infection and Immunity*. **79**: 4042–4049.
- 566 Van Embden JD, Cave MD, Crawford JT, Dale JW, Eisenach KD, Gicquel B, Hermans P, Martin C,  
567 McAdam R, and Shinnick TM. 1993. Strain identification of *Mycobacterium tuberculosis* by  
568 DNA fingerprinting: recommendations for a standardized methodology. *Journal of Clinical*  
569 *Microbiology*. **31**: 406–409.
- 570 Wang L et al. 2022. Multiple genetic paths including massive gene amplification allow *Mycobacterium*  
571 *tuberculosis* to overcome loss of ESX-3 secretion system substrates. *Proceedings of the National*  
572 *Academy of Sciences*. **119**: e2112608119.
- 573 Weiner B et al. 2012. Independent large-scale duplications in multiple *M. tuberculosis* lineages  
574 overlapping the same genomic region. *PLoS ONE*. **7**: e26038.
- 575 WHO 2023. Global Tuberculosis Report 2023. World Health Organization.
- 576 Wick RR, Judd LM, and Holt KE. 2023. Assembling the perfect bacterial genome using Oxford  
577 Nanopore and Illumina sequencing. *PLOS Computational Biology*. **19**: e1010905.
- 578 Xie Z and Tang H. 2017. ISEScan: automated identification of insertion sequence elements in  
579 prokaryotic genomes. *Bioinformatics*. **33**: 3340–3347.



- 580 Xu Y, Yang E, Huang Q, Ni W, Kong C, Liu G, Li G, Su H, and Wang H. 2015. PPE57 induces  
581 activation of macrophages and drives Th1-type immune responses through TLR2. *Journal of*  
582 *Molecular Medicine*. **93**: 645–662.
- 583 Yang Z. 2014. Molecular evolution: a statistical approach. In. Oxford University Press.
- 584 Yang Z, Guarracino A, Biggs PJ, Black MA, Ismail N, Wold JR, Merriman TR, Prins P, Garrison  
585 E, and De Ligt J. 2023. Pangenome graphs in infectious disease: a comprehensive genetic  
586 variation analysis of *Neisseria meningitidis* leveraging Oxford Nanopore long reads. *Frontiers*  
587 *in Genetics*. **14**: 1225248.
- 588 Zhang Y, Zhang H, Zhou T, Zhong Y, and Jin Q. 2011. Genes under positive selection in  
589 *Mycobacterium tuberculosis*. *Computational Biology and Chemistry*. **35**: 319–322.