

SMRT Genome Assembly Corrects Reference Errors, Resolving the Genetic Basis of Virulence in *Mycobacterium tuberculosis*

Afif Elghraoui *¹, Samuel J Modlin *¹, and Faramarz Valafar †¹

¹ Laboratory for Pathogenesis of Clinical Drug Resistance and Persistence, Biological and Medical Informatics Research Center, San Diego State University, San Diego, CA

July 19, 2016

Abstract

The genetic basis of virulence in *Mycobacterium tuberculosis* has been investigated through genome comparisons of its virulent (H37Rv) and attenuated (H37Ra) sister strains. Such analysis, however, relies heavily on the accuracy of the sequences. While the H37Rv reference genome has had several corrections to date, that of H37Ra is unmodified since its original publication. Here, we report the assembly and finishing of the H37Ra genome from single-molecule, real-time (SMRT) sequencing. Our assembly reveals that the number of H37Ra-specific variants is less than half of what the Sanger-based H37Ra reference sequence indicates, undermining and, in some cases, invalidating the conclusions of several studies. PE_PPE family genes, which are intractable to commonly-used sequencing platforms because of their repetitive and GC-rich nature, are overrepresented in the set of genes in which all reported H37Ra-specific variants are contradicted. We discuss how our results change the picture of virulence attenuation and the power of SMRT sequencing for producing high-quality reference genomes.

*Equal contribution

†Corresponding author: faramarz@sdsu.edu

1 Tuberculosis is a serious and pervasive public health problem [1]. It is a disease
2 caused by infection of bacteria from the *Mycobacterium tuberculosis* complex
3 (MTBC). The reference strain, *Mycobacterium tuberculosis* H37Rv, has an at-
4 tenuated counterpart known as H37Ra that is available for studies where facil-
5 ities to handle virulent samples are lacking. H37Ra exhibits a distinct colony
6 morphology, an absence of cord formation, decreased resistance to stress and
7 hypoxia, and attenuated virulence in mammalian models [2–4]. The H37Ra
8 genome was assembled by Zheng and colleagues in 2008 and compared to H37Rv
9 for the purpose of identifying the genetic basis of virulence attenuation [5]. The
10 resulting sequence has been used as the primary avirulent reference genome for
11 *M. tuberculosis* since its publication in 2008.

12 As genome sequencing technology has significantly improved [6], we sought to
13 assess the ability of single-molecule, real-time (SMRT) sequencing for finish-
14 ing mycobacterial genomes. In addition to a high overall GC-content, these
15 genomes have GC-rich repetitive sequences, a source of systematic error for
16 many sequencing protocols. Even sample preparation methods commonly used
17 for shotgun Sanger sequencing are prone to such bias [7]. Sequencing errors
18 in the H37Rv reference have been sought out, with some corrected, others re-
19 maining to be discovered, and still others discovered and remaining to be cor-
20 rected [8,9]. The Pacific Biosciences RS II platform has been shown to produce
21 finished-grade assemblies of microbial genomes exceeding the quality of Sanger
22 sequencing [10–12].

23 In this study, we sequenced and assembled the genome of *M. tuberculosis* H37Ra
24 and compared it to the reference sequence. We further compared both sequences
25 against the reference sequence for *M. tuberculosis* H37Rv and re-evaluated the
26 conclusions of Zheng and colleagues with respect to the genetic basis of virulence
27 attenuation.

28 **Results**

29 **Genome Assembly and Methylation Motif Detection**

30 With the data from one sequencing run (SMRTCell), the genome assembled
31 with 103x average coverage into a single contig containing 4426109 base pairs
32 after circularization and polishing. Performing the assembly with data from two
33 SMRTCells (217x average coverage) resulted in an identical sequence.

34 In the raw assembly, circularization was impeded by discrepancies in the edges
35 of the contig, where an IS6110 insertion was present in only one of the two
36 edges. It appears heterogeneously in our sample, as aligning our reads against
37 our assembly shows that a minority of reads have interrupted mapping to this
38 segment while the majority do not. With regard to base modifications, N6-
39 methyladenine was detected in 99.67% of the instances of the partner sequence
40 motifs CTGGAG and CTCCAG. The methylation of these motifs in both H37Ra
41 and H37Rv was previously reported by Zhu and colleagues in H37Ra as part of
42 their study of mycobacterial methylomes [13].

43 **Direct Comparison with the Hitherto H37Ra Reference** 44 **Genome**

45 Comparison of our assembly with the H37Ra reference sequence (NC_009525.1,
46 hereafter referred to as H37RaJH, for Johns Hopkins) showed significant varia-
47 tion. We found 33 single nucleotide polymorphisms (SNPs), and 77 insertions
48 and deletions in our assembly with respect to H37RaJH (Supplementary Data
49 1).

50 Structural Variations

51 Two of the insertions with respect to H37RaJH were substantial structural vari-
52 ations: one was an insertion of IS6110 into the gene corresponding to Rv1764
53 and the other was an in-frame insertion of 3456bp into the PPE54 gene.

54 The insertion of IS6110 into Rv1764 (an IS6110 transposase) is unsurprising,
55 as IS6110 insert frequently into that general region of the genome, as well as
56 within their transposase [14, 15]. This insertion was the heterogeneous inser-
57 tion responsible for the discrepant contig ends in our raw genome assembly.
58 Such heterogeneity implies either a lack of selection pressure on the insertion in
59 culture, a recent emergence of the insertion, or both.

60 The 3456bp insertion in *ppe54* with respect to H37RaJH incidentally corre-
61 sponds to a tandem duplication of a 1728bp sequence at the same site in H37Rv
62 with 100% identity. The complete absence of this tandem repeat at this site in
63 H37RaJH, however, is not necessarily an assembly error, as this is also observed
64 in several clinical isolates (unpublished data). This, along with the 100% iden-
65 tity between each 1728bp duplicate of the tandem repeat with respect to H37Rv,
66 lead us to believe that both the duplication in our sequence and the deletion ob-
67 served in H37RaJH are instances of *in vitro* evolution, following the divergence
68 of the lineages from which H37RaJH and our assembly were drawn.

69 These two structural variations, or, at least, very similar structural variations,
70 have been observed previously in virulent strains of *M. tuberculosis*, and there-
71 fore likely do not contribute to virulence attenuation in H37Rv (unpublished
72 data) [14, 16], but shed light on *in vitro* evolution of this strain [8, 17].

73 **Analysis of Motif Variants in H37Ra and H37Rv**

74 With the knowledge that the CTGGAG/CTCCAG motifs are methylated in
75 both H37Ra and H37Rv [13], we determined the motif variants, or sequence
76 polymorphisms that create or destroy motifs, between H37Rv and H37Ra. By
77 first comparing H37RaJH to H37Rv, we see that all but two motif variants were
78 due to structural variations. Both of these variants instantiate the CTGGAG
79 motif in H37Ra where it is absent in the H37Rv reference sequence. The first
80 is due to the $G \rightarrow T$ polymorphism at H37Rv position 2043284 (upstream of
81 PPE30) in H37RaJH, but this variant is contradicted by our H37Ra assembly.
82 The second is due to the $T \rightarrow G$ polymorphism at H37Rv position 2718852
83 (upstream of *nadD*) and confirmed by our H37Ra assembly, yet also appears
84 in CDC1551 and is a previously reported sequencing error in H37Rv [8] that
85 has not been applied to the current reference. Based on these results, DNA
86 methylation and motif variants do not play a role in the attenuation of virulence
87 in H37Ra.

88 **Status of Previously Reported “H37Ra-specific” Polymor-** 89 **phisms**

90 With our assembly, we aimed to replicate the study performed by Zheng and
91 colleagues when they first assembled the H37Ra genome [5]. In their study, they
92 compared their assembly with H37Rv, then filtered out variants also present in
93 CDC1551 (NC_002755.2) to find mutations likely specific to H37Ra [5]. Zheng
94 and colleagues identified a set of mutations in H37Ra unique with respect to
95 H37Rv and CDC1551 as “H37Ra-specific”. These mutations fall within or ad-
96 jacent to (which we term “affecting”) 56 genes in H37Rv, which we refer to as
97 the high-confidence (HC) gene set. While comparing the variants, Zheng and

98 colleagues also discovered sequencing errors in the H37Rv reference sequence [5],
99 a number of which were corrected in NC_000962.3 [9], the version used in our
100 study.

101 To see how well the HC genes are supported by our assembly of H37Ra, we
102 determined variants with respect to H37Rv for our assembly and H37RaJH and
103 performed set comparisons after excluding mutations shared with CDC1551
104 (Supplementary Data 1-2). We then categorized the HC genes as follows. We
105 labeled a gene “unsupported” if all mutations affecting it were observed only
106 in H37RaJH. We labeled a gene “supported” if all mutations affecting it were
107 observed in both H37Ra assemblies. Otherwise, we labeled a gene “adjusted”
108 if it had a different variant profile between H37RaJH and our assembly in a
109 manner distinct from the two categories defined above. Figure 1 shows example
110 classifications based on these criteria.

111 We first noted that two of the HC variants reported by Zheng and colleagues,
112 those affecting *nadD* (Rv2421c) and *nrdH* (Rv3053), were included erroneously
113 (Table 1d). These variants were a $T \rightarrow G$ mutation 44 bases upstream of *nadD*,
114 at H37Rv position 2718852, and a 14bp deletion in the promoter of *nrdH*. These
115 mutations, although confirmed by our assembly, also appear in CDC1551 and
116 thus cannot be considered H37Ra-specific.

117 Of the variants in the remaining 54 HC genes, our assembly contradicts 35 (Ta-
118 ble 1a), adjusts 5 (Table 1b), and confirms 14 (Table 1c). We then considered
119 how these results affect the picture of how the genotypic differences between
120 H37Rv and H37Ra give rise to the phenotypic differences observed between the
121 two strains, which are discussed below and depicted graphically in Figure 3. As
122 our analysis focused on the HC gene set reported by Zheng and colleagues [5],
123 we did not re-evaluate whether additional genes and variants should belong to
124 this grouping. We did, however, carefully consider all variants unique to our as-

125 sembly (Table 2) and their potential effect on the organism's phenotype.

126 Accuracy of the H37Rv Reference Sequence

127 Ioerger and colleagues listed 73 polymorphisms (excluding those in PE_PPE
128 genes) with respect to the H37Rv reference shared between six H37Rv strains
129 from different laboratories, but considered all but one of them as errors in the
130 reference sequence because they also appeared in the H37Ra reference [8]. The
131 remaining polymorphism was a $A \rightarrow C$ transversion at position 459399, a posi-
132 tion upstream of Rv0383c masked by a 55bp deletion in H37RaJH. Interestingly,
133 our assembly contradicts this 55bp deletion, but is in perfect concordance with
134 the transversion at position 459399. The revelation that H37Ra is in fact the
135 same as all H37Rv strains at this position invalidates the maximum parsimony
136 tree in figure 1 of their publication [8]. Thus, through our improved assembly
137 of the H37Ra genome, we have identified an additional error in H37Rv, the
138 standard reference genome of *M. tuberculosis*.

139 SNPs Previously Reported to Cause Expression Changes in H37Ra 140 are Contradicted by Our Assembly

141 Interestingly, SNPs in the putative promoter regions of two genes, *phoH2* and
142 *sigC*, found by Zheng and colleagues to be up-regulated *in vitro* and down-
143 regulated in macrophage in H37Ra relative to H37Rv, were contradicted by our
144 assembly [5]. Zheng and colleagues attributed this differential expression to
145 these (now contradicted) SNPs, but it appears there instead must be a distal
146 causative factor driving the observed expression changes of both genes. The
147 SNP affecting *sigC* has been cited as the cause of the differential expression
148 of SigC in macrophages relative to H37Rv [18, 19], illustrating how incorrect

149 sequences can propagate through the literature.

150 **SNPs Previously Thought to Affect Polyketide Synthesis in H37Ra**
151 **are Contradicted by Our Assembly**

152 Altered polyketide synthesis has been proposed as one of the primary mech-
153 anisms attenuating virulence in H37Ra, through disrupting phthiocerol dimy-
154 cocerosate (PDIM) production, which has shown to manifest deleteriously in
155 H37Ra [20,21]. Our assembly contradicts both reported SNPs in *pks12* (polyke-
156 tide synthase 12) of H37RaJH. This means that some factor other than disrup-
157 tion of *pks12* causes the observed lowered PDIM production in H37Ra. Thus,
158 it remains unclear which (epi)genomic factor(s) underlie the observed reduction
159 in PDIM synthesis in H37Ra, as supported variants (those in *phoP* and *nrp*)
160 once considered to cause this reduction [22] have been shown not to [23, 24].
161 However, it is possible the decreased production of PDIMs is merely an artifact
162 of repeated subculturing *in vitro* [17].

163 **Variants in *phoP*, *mazG*, and *hadC* Account for Much of the Virulence**
164 **Attenuation in H37Ra**

165 Of all the HC genes, only variants in *phoP*, *mazG*, and *hadC* have been con-
166 nected strongly with virulence attenuation in H37Ra through wet-lab work, each
167 of which our assembly supports.

168 Of these, the most thoroughly studied is the nsSNP (S219L) in the DNA-binding
169 region of *phoP*, part of the two component *PhoPR* regulatory system. There
170 is an abundance of literature linking *phoP* to virulence attenuation in H37Ra,
171 through several mechanisms, including disrupted sulfolipid and trehalose synthe-
172 sis (Figure 2), diminished ESAT-6 secretion, and additional downstream effects

173 from altered expression of other genes under its regulon [5, 18, 23, 25–30]. How-
174 ever, several of these studies also show that *phoP* alone [23, 29] is not responsible
175 for virulence attenuation in H37Ra, but rather that the genomic cause behind
176 virulence attenuation in H37Ra is multifactorial.

177 The second gene, *mazG*, has a nsSNP (A219E) in a region coding for a highly
178 conserved alpha-helix residue in its protein product, a nucleoside triphosphate
179 (NTP) pyrophosphohydrolase [5]. MazG exhibits diminished hydrolysis activity
180 in H37Ra relative to both MazG in H37Rv and MazG of the fast-growing *M.*
181 *smegmatis*. Wild-type MazG hydrolyzes all NTPs, including those that are mu-
182 tagenic and appear more frequently with oxidative stress (Figure 3b), which is
183 experienced by the bacterium inside activated macrophages [31]. This decreased
184 ability to hydrolyze mutagenic NTPs contributes to virulence attenuation in
185 H37Ra [32].

186 In the third gene, *hadC*, there is a frameshift-inducing 1-bp insertion, which
187 creates a premature stop codon and results in truncation of HadC. *hadC* is a
188 member of the essential *hadA-hadB-hadC* gene cluster, which forms two hy-
189 dratases (HadAB and HadBC) of the *M. tuberculosis* fatty acid synthase II sys-
190 tem. Our assembly and H37RaJH both show a 5-bp insertion in *hadA* which,
191 along with *hadC*, are the only genes with variants in H37Ra [33] that encode
192 proteins known to be necessary for mycolic acid synthesis.

193 Recent complementation and knockout studies using *hadC* from H37Ra and
194 H37Rv showed that intact HadC is necessary for cord formation, and that the
195 truncated form *H37Ra/hadC* affects length and oxygenation of mycolic acids
196 (Figure 2b). Furthermore, when tested in murine lung and spleen, *H37Ra/hadC_{Rv}*
197 grew an intermediate amount of colony forming units, between that of H37Ra
198 and H37Rv, at a level commensurate with *H37RvΔhadC* which suggests that
199 the H37Ra *hadC* variant underlies some of its virulence attenuation [33].

200 Interestingly, while both our assembly and H37RaJH harbor a 5-bp insertion
201 in *hadA*, sequences obtained by Lee, Slama, and their respective colleagues do
202 not [29,33]. These two sequences were both derived from a culture from Institut
203 Pasteur, while ours and that of Zheng and colleagues [5] were acquired directly
204 from ATCC, which suggests that the two cultures diverged *in vitro* prior to
205 sequencing despite sharing the same ATCC identifier. We expect the deleteri-
206 ous effects of *hadC_{Ra}* shown by Slama and colleagues would be exacerbated by
207 the 5bp insertion in our assembly, as it results in an abnormal HadAB enzyme
208 which, when normal, has been posited to compensate for faulty HadBC [33].
209 However, the experiments discussed above indicate that the *hadC* variant alone
210 is sufficient to attenuate virulence, and is one of the primary sources of attenu-
211 ation in H37Ra.

212 **Copy Number Variation in *lpdA* Promoter**

213 The polymorphism reported in H37RaJH that affects *lpdA* (NAD(P)H quinone
214 reductase) is a third (as opposed to the two in H37Rv) 58bp repeat in its
215 promoter region. Our assembly reveals an additional two copies of this 58bp
216 region, resulting in a total of five copies of the repeat. LpdA has been shown to
217 protect bacilli from oxidative stress and improve *M. tuberculosis* survival in a
218 mouse model, which suggests that if this copy number variation disrupts typical
219 expression of LpdA, it may contribute to the phenotype of H37Ra [34]. This
220 may also affect the expression of *glpD2* (glycerol-3-phosphate dehydrogenase),
221 as it shares an operon with *lpdA* [5].

222 **Variants Affecting Uncharacterized Hypothetical Genes**

223 Several genes classified with unknown or hypothetical functions were among the
224 HC genes of H37RaJH (Table 1). Our assembly contradicts all variants in the
225 majority of these, leaving three which we supported in full.

226 Though none of these genes have an implicated role in virulence in the literature,
227 they may in reality. These genes should be investigated, as they are three of the
228 few supported HC genes yet to be explored. The value of exploring hypothetical
229 genes is evidenced by the recent discovery of a significant contribution of HadC
230 [33]—the function of which was unknown when H37RaJH was published—to
231 virulence attenuation in H37Ra (Figure 2).

232 **Significant Reduction of H37Ra-specific Variants in PE_PPE genes**

233 The PE_PPE family of genes is unique to mycobacteria but poorly characterized,
234 both functionally and genomically, in *M. tuberculosis*, the latter owing to the
235 family's high-GC content and repetitive nature [35]. Evidence for contribution
236 from PE_PPE family members to virulence has amassed support since 2008
237 [36–39], but this gene family was the most drastically altered by our assembly:
238 while PE_PPE genes comprise approximately 10% of the genome, they account
239 for nearly half (16/35) of the unsupported genes. It is likely that the majority
240 of these are errors in H37RaJH rather than manifestations of hypervariability,
241 as few PE_PPE genes fell into the adjusted or novel categories, as one would
242 expect if they were due to hypervariability.

243 Consequently, some extant work examining polymorphic PE_PPE genes between
244 H37RaJH and H37Rv is invalidated by our assembly. For example, our assembly
245 contradicts or changes the variant profile of all four PE_PPE genes reported to
246 be positively selected for in H37Ra in an evolutionary genomics study by Zhang

247 and colleagues [38] using H37RaJH.

248 Another study affected profoundly by our results is that of Kohli and colleagues
249 [36], which used H37RaJH and H37Rv in an *in silico* genomic and proteomic
250 comparison of PE_PPE family genes. Though our assembly renders much of the
251 results from their analyses invalid, applying their methodology to our updated
252 assembly would yield interesting results.

253 Our assembly contains polymorphisms in 6 of 22 genes that encode PE_PPE
254 family members reported as unique to H37RaJH (Table 1, Figure 3b). Of the
255 three PE_PPE family members fully corroborated by our sequence, one was
256 the duplication of *ppe38*, which McEvoy and colleagues have also identified in
257 3 different samples of H37Rv, suggesting this duplication likely plays no role
258 in virulence [40]. All 3 of the adjusted PE_PPE family members, as well as
259 the supported *Wag22* and *PPE13*, belong to PE_PPE sublineage V. Sublineage
260 V members comprise the majority of PE_PPE proteins that interact with the
261 host, and are overrepresented in proteomic studies of *in vivo* infection [35]. This
262 enrichment of subfamily V PE_PPE family members in the set of supported or
263 adjusted genes suggests they may be more integral to virulence attenuation
264 in H37Ra than other PE_PPE family members. The role of PE_PPE family
265 members in virulence should become better understood as more genomes are
266 sequenced using third-generation platforms.

267 In addition to the differences due to sequence alterations in PE_PPE family
268 genes, the corroborated polymorphism in *phoP* may confer altered expression
269 of many PE_PPE family proteins, as at least 13 are under its regulon [35], which
270 could mediate some virulence attenuation.

271 The precise roles of PE_PPE family members have yet to be elucidated in full.
272 It is difficult to evaluate rigorously the effect of each PE_PPE variant, as their

273 function in wild-type *M. tuberculosis* is poorly characterized [35]. Moreover,
274 their contribution to virulence may well require complexities of the native host
275 environment beyond what can be replicated *in vitro* or *ex vivo* with current
276 technology. Thus, the role the polymorphisms in this family play in the phe-
277 notype of H37Ra compel further research, which our reduction of variants has
278 made more tractable.

279 Discussion

280 Since its publication in 2008 [5], several studies have used the whole genome
281 [8, 36, 41–46] of H37RaJH, or the reported differences between H37RaJH and
282 H37Rv [47] in their analyses. Our improved assembly changes the implications
283 of several of these *in silico* analyses. Additionally, several studies have used the
284 set of genes with variants in H37RaJH with respect to H37Rv to guide wet-lab
285 experiments [48, 49]. Re-examining these studies with our assembly of H37Ra
286 may yield novel insights, as unsupported variants can serve as a retroactive
287 control.

288 Our *de novo* assembly using single-molecule sequencing has reduced the set
289 of genes polymorphic to H37Rv by more than half, clarifying which genomic
290 factors most likely give rise to virulence attenuation and other H37Ra-specific
291 phenotypes. For an expanded discussion of genes affected and their ties to vir-
292 ulence, see the supplementary note. Supported variants affecting PhoP, MazG,
293 and HadC have been experimentally affirmed [23, 32, 33], gaining insight into
294 how they manifest in the phenotype of H37Ra, but basic mechanisms for their
295 contributions are not fully elucidated. A few other supported or adjusted genes
296 (*lpdA*, *pabB*, and *nrp*) have been indirectly connected to the avirulence of H37Ra
297 through experiments on other mycobacterial species [22] or H37Ra complemen-

298 tation studies measuring proxies of virulence [48].

299 It is clear that the nsSNP in *phoP* remains a potent mediator of virulence of
300 *M. tuberculosis* through affecting SL and ATHL activity (Figure 2), while the
301 truncation of HadC enfeebles the mycomembrane (Figure 2b). Polymorphisms
302 in *mazG* and *lpdA* may each confer compromised stress response mechanisms in
303 H37Ra (Figure 3), which are critical to enduring the intramacrophage environ-
304 ment of the host [32, 34]. Variants affecting genes with regulatory functions—
305 *phoP* and others with roles in regulation not yet known—may also cause down-
306 stream effects on H37Ra phenotype, which may prove difficult to characterize.
307 The variants in genes of the PE_PPE family and hypothetical genes (Rv0010c,
308 Rv0039c, and Rv1006) potentially contribute to virulence attenuation through
309 mechanisms not yet identified. Thus, with the greater accuracy of our assem-
310 bly, wet-lab studies can focus on the true differences between the H37Ra and
311 H37Rv, and computational studies will be in greater concordance with reality,
312 yielding more useful results.

313 The advantages of single-molecule sequencing are readily apparent in our re-
314 sults. The random error profile of this technology allows for consensus accuracy
315 to increase as a function of sequencing depth [10]. Performing the assembly
316 with a doubled sequencing depth resulted in an identical sequence, indicating
317 that we were able to maximize the sequence’s accuracy with a single sequencing
318 run. The long reads produced by this technology allowed us to easily and un-
319 ambiguously capture known structural variants in H37Ra, as well as two novel
320 to the strain. We were also able to fully resolve the GC-rich and repetitive
321 PE_PPE genes, sequences which compound the weaknesses of most sequencing
322 technologies. As a result, our assembly demonstrates that H37Ra is significantly
323 more similar to H37Rv than indicated by H37Ra’s Sanger-based reference se-
324 quence, with contradicted variants overrepresented in the difficult sequences of

325 the PE_PPE genes. While *in vitro* evolution may underlie some of the differences
326 between our assembly and H37RaJH, we believe that most of the contradicted
327 variants (Table 1a) reflect sequencing errors in H37RaJH due to the disparity
328 in sequencing quality. Regardless, the contradicted variants should not be con-
329 sidered as characteristic of H37Ra or its attenuated virulence. These sites were
330 concordant with H37Rv and we did not find additional polymorphic PE_PPE
331 genes with respect to H37Rv (Table 2), indicating a disparity in sequencing
332 accuracy even among the Sanger-based references. On the other hand, the fact
333 that we have not resequenced H37Rv and CDC1551 is a limitation of our study,
334 where we have relied on their Sanger-based reference sequences for determining
335 H37Ra-specific variants. We believe that the impact of the latter is minimal
336 and that the former is the dominant factor, considering the level of concordance
337 with H37Rv and CDC1551 in the instances where our sequence disagreed with
338 H37RaJH.

339 Studies that have relied on the accuracy of PE_PPE sequences in the H37Ra ref-
340 erence sequence were the most severely impacted by our study. We consequently
341 advise caution when analyzing GC-rich and repetitive sequences among refer-
342 ence genomes, not to mention draft genomes. As *de novo* assembly can be rou-
343 tinely performed for microbes using single-molecule sequencing, we strongly rec-
344 ommend this for mycobacteria, especially because of their PE_PPE genes.

345 **Competing interests**

346 The authors declare that they have no competing interests.

347 **Author's contributions**

348 F.V., A.E., and S.J.M. designed the study. A.E. performed the *de novo* as-
349 sembly, methylation analysis, and comparative genomics analyses. S.J.M. per-
350 formed the literature review, interpreted the results, and wrote the supplemen-
351 tary note. A.E. and S.J.M. prepared the manuscript, which was reviewed and
352 approved by all authors.

353 **Acknowledgements**

354 We would like to thank Jason Chin and Richard Hall from Pacific Biosciences for
355 discussions on *de novo* assembly methodology and quality assessment. Logan
356 Fink also provided some assistance with quality assessment of the assembly. We
357 would also like to thank Antonino Catanzaro, Timothy Rodwell, and their staff
358 for bacterial culture and DNA extraction. Jonas Korlach, Anthony Baughn,
359 Sarah Ramirez-Busby, Ragavi Shanmugam, Carmela Chan, Amy Goodmanson,
360 Daeheon Oh, and Logan Fink reviewed and provided helpful feedback on drafts
361 of the manuscript. This work was funded by a grant from National Institute
362 of Allergy and Infectious Diseases (NIAID Grant No. R01AI105185). A.E.,
363 S.J.M., and F.V. were supported by this grant. S.J.M. was also supported by
364 scholarships from a National Science Foundation Grant (no. 0966391).

365 **References**

- 366 [1] *Global tuberculosis report 2015* (World Health Organization).
- 367 [2] Middlebrook, G., Dubos, R. J. & Pierce, C. VIRULENCE AND MOR-
368 PHOLOGICAL CHARACTERISTICS OF MAMMALIAN TUBERCLE

- 369 BACILLI. *The Journal of Experimental Medicine* **86**, 175–184 (1947).
370 URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2135722/>.
- 371 [3] Heplar, J. Q., Clifton, C. E., Raffel, S. & Futrelle, C. M. Virulence of the
372 Tubercle Bacillus: I. Effect of Oxygen Tension upon Respiration of Virulent
373 and Avirulent Bacilli. *The Journal of Infectious Diseases* **94**, 90–98 (1954).
374 URL <http://www.jstor.org/stable/30089826>.
- 375 [4] Alsaadi, A.-I. & Smith, D. W. The Fate of Virulent and Attenuated
376 Mycobacteria in Guinea Pigs Infected by the Respiratory Route. *Amer-*
377 *ican Review of Respiratory Disease* **107**, 1041–1046 (1973). URL <http://www.atsjournals.org/doi/abs/10.1164/arrd.1973.107.6.1041>.
378
- 379 [5] Zheng, H. *et al.* Genetic Basis of Virulence Attenuation Revealed by Com-
380 parative Genomic Analysis of Mycobacterium tuberculosis Strain H37ra
381 versus H37rv. *PLoS ONE* **3**, e2375 (2008). URL [http://dx.plos.org/](http://dx.plos.org/10.1371/journal.pone.0002375)
382 [10.1371/journal.pone.0002375](http://dx.plos.org/10.1371/journal.pone.0002375).
- 383 [6] Koren, S. & Phillippy, A. M. One chromosome, one contig: complete mi-
384 crobial genomes from long-read sequencing and assembly. *Current Opinion*
385 *in Microbiology* **23**, 110–120 (2015). URL [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S1369527414001817)
386 [com/science/article/pii/S1369527414001817](http://www.sciencedirect.com/science/article/pii/S1369527414001817).
- 387 [7] Ross, M. G. *et al.* Characterizing and measuring bias in sequence data.
388 *Genome Biol* **14**, R51 (2013). URL [http://www.biomedcentral.com/](http://www.biomedcentral.com/content/pdf/gb-2013-14-5-r51.pdf)
389 [content/pdf/gb-2013-14-5-r51.pdf](http://www.biomedcentral.com/content/pdf/gb-2013-14-5-r51.pdf).
- 390 [8] Ioerger, T. R. *et al.* Variation among Genome Sequences of H37rv Strains of
391 Mycobacterium tuberculosis from Multiple Laboratories. *Journal of Bac-*
392 *teriology* **192**, 3645–3653 (2010). URL [http://www.ncbi.nlm.nih.gov/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2897344/)
393 [pmc/articles/PMC2897344/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2897344/).

- 394 [9] Köser, C. U., Niemann, S., Summers, D. K. & Archer, J. A. C. Overview of
395 errors in the reference sequence and annotation of *Mycobacterium tuber-*
396 *culosis* H37rv, and variation amongst its isolates. *Infection, Genetics and*
397 *Evolution* **12**, 807–810 (2012). URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S1567134811002243)
398 [science/article/pii/S1567134811002243](http://www.sciencedirect.com/science/article/pii/S1567134811002243).
- 399 [10] Roberts, R. J., Carneiro, M. O. & Schatz, M. C. The advantages of
400 SMRT sequencing. *Genome Biol* **14**, 405 (2013). URL [http://www.](http://www.biomedcentral.com/content/pdf/gb-2013-14-7-405.pdf)
401 [biomedcentral.com/content/pdf/gb-2013-14-7-405.pdf](http://www.biomedcentral.com/content/pdf/gb-2013-14-7-405.pdf).
- 402 [11] Koren, S. *et al.* Reducing assembly complexity of microbial genomes with
403 single-molecule sequencing. *Genome Biology* **14**, 1–16 (2013). URL [http:](http://dx.doi.org/10.1186/gb-2013-14-9-r101)
404 [//dx.doi.org/10.1186/gb-2013-14-9-r101](http://dx.doi.org/10.1186/gb-2013-14-9-r101).
- 405 [12] Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies
406 from long-read SMRT sequencing data. *Nature Methods* **10**, 563–
407 569 (2013). URL [http://www.nature.com.libproxy.sdsu.edu/nmeth/](http://www.nature.com.libproxy.sdsu.edu/nmeth/journal/v10/n6/abs/nmeth.2474.html)
408 [journal/v10/n6/abs/nmeth.2474.html](http://www.nature.com.libproxy.sdsu.edu/nmeth/journal/v10/n6/abs/nmeth.2474.html).
- 409 [13] Zhu, L. *et al.* Precision methylome characterization of *Mycobacterium*
410 *tuberculosis* complex (MTBC) using PacBio single-molecule real-time
411 (SMRT) technology. *Nucleic Acids Research* **44**, 730–743 (2016). URL
412 <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkv1498>.
- 413 [14] Roychowdhury, T., Mandal, S. & Bhattacharya, A. Analysis of IS6110 in-
414 sersion sites provide a glimpse into genome evolution of *Mycobacterium tu-*
415 *berculosis*. *Scientific Reports* **5**, 12567 (2015). URL [http://www.nature.](http://www.nature.com/articles/srep12567)
416 [com/articles/srep12567](http://www.nature.com/articles/srep12567).
- 417 [15] Alonso, H., Samper, S., Martín, C. & Otal, I. Mapping is6110 in high-
418 copy number mycobacterium tuberculosis strains shows specific insertion

- 419 points in the beijing genotype. *BMC Genomics* **14**, 1–11 (2013). URL
420 <http://dx.doi.org/10.1186/1471-2164-14-422>.
- 421 [16] Lari, N., Rindi, L. & Garzelli, C. Identification of one insertion site of
422 IS6110 in Mycobacterium tuberculosis H37ra and analysis of the RvD2
423 deletion in M. tuberculosis clinical isolates. *Journal of Medical Microbiology*
424 **50**, 805–811 (2001).
- 425 [17] Andreu, N. & Gibert, I. Cell population heterogeneity in Mycobacterium
426 tuberculosis H37rv. *Tuberculosis (Edinburgh, Scotland)* **88**, 553–559 (2008).
- 427 [18] Dokladda, K., Billamas, P. & Palittapongarnpim, P. Different behaviours of
428 promoters in Mycobacterium tuberculosis H37rv and H37ra. *World Journal*
429 *of Microbiology and Biotechnology* **31**, 407–413 (2015). URL [http://link.](http://link.springer.com/article/10.1007/s11274-014-1794-x)
430 [springer.com/article/10.1007/s11274-014-1794-x](http://link.springer.com/article/10.1007/s11274-014-1794-x).
- 431 [19] Malhotra, V., Tyagi, J. S. & Clark-Curtiss, J. E. DevR-mediated adap-
432 tive response in Mycobacterium tuberculosis H37ra: links to asparagine
433 metabolism. *Tuberculosis (Edinburgh, Scotland)* **89**, 169–174 (2009). URL
434 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2693488/>.
- 435 [20] Daffé, M., Lacave, C., Lanéelle, M. A., Gillois, M. & Lanéelle, G. Polyph-
436 thienoyl trehalose, glycolipids specific for virulent strains of the tubercle
437 bacillus. *European journal of biochemistry / FEBS* **172**, 579–584 (1988).
- 438 [21] Middlebrook, G., Coleman, C. M. & Schaefer, W. B. SULFOLIPID FROM
439 VIRULENT TUBERCLE BACILLI*. *Proceedings of the National Academy*
440 *of Sciences of the United States of America* **45**, 1801–1804 (1959). URL
441 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC222804/>.
- 442 [22] Hotter, G. S. *et al.* Transposon Mutagenesis of Mb0100 at the ppe1-
443 nrp Locus in Mycobacterium bovis Disrupts Phthiocerol Dimycocerosate

- 444 (PDIM) and Glycosylphenol-PDIM Biosynthesis, Producing an Aviru-
445 lent Strain with Vaccine Properties At Least Equal to Those of *M. bo-*
446 *vis* BCG. *Journal of Bacteriology* **187**, 2267–2277 (2005). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1065232/>.
- 448 [23] Chesne-Seck, M.-L. *et al.* A Point Mutation in the Two-Component Regu-
449 lator PhoP-PhoR Accounts for the Absence of Polyketide-Derived Acyltre-
450 haloses but Not That of Phthiocerol Dimycocerosates in *Mycobacterium*
451 *tuberculosis* H37ra. *Journal of Bacteriology* **190**, 1329–1334 (2008). URL
452 <http://jbs.asm.org/content/190/4/1329>.
- 453 [24] Hotter, G. S. & Collins, D. M. *Mycobacterium bovis* lipids: Virulence
454 and vaccines. *Veterinary Microbiology* **151**, 91–98 (2011). URL <http://www.sciencedirect.com/science/article/pii/S0378113511001118>.
- 456 [25] Li, A. H. *et al.* Contrasting Transcriptional Responses of a Viru-
457 lent and an Attenuated Strain of *Mycobacterium tuberculosis* Infecting
458 Macrophages. *PLoS ONE* **5** (2010). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2883559/>.
- 460 [26] Asensio, J. G. *et al.* The Virulence-associated Two-component PhoP-
461 PhoR System Controls the Biosynthesis of Polyketide-derived Lipids in
462 *Mycobacterium tuberculosis*. *Journal of Biological Chemistry* **281**, 1313–
463 1316 (2006). URL <http://www.jbc.org/content/281/3/1313>.
- 464 [27] Frigui, W. *et al.* Control of *M. tuberculosis* ESAT-6 Secretion and Specific
465 T Cell Recognition by PhoP. *PLoS Pathogens* **4** (2008). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2242835/>.
- 467 [28] Walters, S. B. *et al.* The *Mycobacterium tuberculosis* PhoPR
468 two-component system regulates genes essential for virulence and
469 complex lipid biosynthesis. *Molecular Microbiology* **60**, 312–330

- 470 (2006). URL [http://onlinelibrary.wiley.com/doi/10.1111/j.](http://onlinelibrary.wiley.com/doi/10.1111/j.1365-2958.2006.05102.x/abstract)
471 [1365-2958.2006.05102.x/abstract](http://onlinelibrary.wiley.com/doi/10.1111/j.1365-2958.2006.05102.x/abstract).
- 472 [29] Lee, J. S. *et al.* Mutation in the Transcriptional Regulator PhoP Con-
473 tributes to Avirulence of Mycobacterium tuberculosis H37ra Strain. *Cell*
474 *Host & Microbe* **3**, 97–103 (2008). URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S1931312808000279)
475 [science/article/pii/S1931312808000279](http://www.sciencedirect.com/science/article/pii/S1931312808000279).
- 476 [30] Solans, L. *et al.* The PhoP-Dependent ncRNA Mcr7 Modulates the TAT Se-
477 cretion System in Mycobacterium tuberculosis. *PLoS Pathogens* **10** (2014).
478 URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4038636/>.
- 479 [31] Lyu, L.-D., Tang, B.-K., Fan, X.-Y., Ma, H. & Zhao, G.-P. Mycobacte-
480 rial MazG Safeguards Genetic Stability via Housecleaning of 5-OH-dCTP.
481 *PLoS Pathogens* **9** (2013). URL [http://www.ncbi.nlm.nih.gov/pmc/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3855555/)
482 [articles/PMC3855555/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3855555/).
- 483 [32] Lu, L.-d. *et al.* Mycobacterial MazG Is a Novel NTP Pyrophosphohydrolase
484 Involved in Oxidative Stress Response. *Journal of Biological Chemistry*
485 **285**, 28076–28085 (2010). URL [http://www.jbc.org/content/285/36/](http://www.jbc.org/content/285/36/28076)
486 [28076](http://www.jbc.org/content/285/36/28076).
- 487 [33] Slama, N. *et al.* The changes in mycolic acid structures caused by hadC mu-
488 tation have a dramatic effect on the virulence of Mycobacterium tuberculo-
489 sis. *Molecular Microbiology* n/a–n/a (2015). URL [http://onlinelibrary.](http://onlinelibrary.wiley.com/doi/10.1111/mmi.13266/abstract)
490 [wiley.com/doi/10.1111/mmi.13266/abstract](http://onlinelibrary.wiley.com/doi/10.1111/mmi.13266/abstract).
- 491 [34] Akhtar, P. *et al.* Variable-number tandem repeat 3690 polymorphism
492 in Indian clinical isolates of Mycobacterium tuberculosis and its in-
493 fluence on transcription. *Journal of Medical Microbiology* **58**, 798–
494 805 (2009). URL [http://jmm.microbiologyresearch.org/content/](http://jmm.microbiologyresearch.org/content/journal/jmm/10.1099/jmm.0.002550-0)
495 [journal/jmm/10.1099/jmm.0.002550-0](http://jmm.microbiologyresearch.org/content/journal/jmm/10.1099/jmm.0.002550-0).

- 496 [35] Fishbein, S., van Wyk, N., Warren, R. M. & Sampson, S. L. Phylogeny
497 to function: PE/PPE protein evolution and impact on Mycobacterium tu-
498 berculosis pathogenicity. *Molecular Microbiology* **96**, 901–916 (2015). URL
499 <http://onlinelibrary.wiley.com/doi/10.1111/mmi.12981/abstract>.
- 500 [36] Kohli, S. *et al.* Comparative genomic and proteomic analyses of PE/PPE
501 multigene family of Mycobacterium tuberculosis H37rv and H37ra reveal
502 novel and interesting differences with implications in virulence. *Nucleic*
503 *Acids Research* **40**, 7113–7122 (2012). URL [http://nar.oxfordjournals.](http://nar.oxfordjournals.org/content/40/15/7113)
504 [org/content/40/15/7113](http://nar.oxfordjournals.org/content/40/15/7113).
- 505 [37] Yu, G. *et al.* Integrative analysis of transcriptome and genome indi-
506 cates two potential genomic islands are associated with pathogenesis of
507 Mycobacterium tuberculosis. *Gene* **489**, 21–29 (2011). URL [http:](http://www.sciencedirect.com/science/article/pii/S0378111911004719)
508 [//www.sciencedirect.com/science/article/pii/S0378111911004719](http://www.sciencedirect.com/science/article/pii/S0378111911004719).
- 509 [38] Zhang, Y., Zhang, H., Zhou, T., Zhong, Y. & Jin, Q. Genes under posi-
510 tive selection in Mycobacterium tuberculosis. *Computational Biology and*
511 *Chemistry* **35**, 319–322 (2011). URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S1476927111000934)
512 [science/article/pii/S1476927111000934](http://www.sciencedirect.com/science/article/pii/S1476927111000934).
- 513 [39] Ahmed, A., Das, A. & Mukhopadhyay, S. Immunoregulatory functions and
514 expression patterns of PE/PPE family members: Roles in pathogenicity
515 and impact on anti-tuberculosis vaccine and drug design. *IUBMB Life* **67**,
516 414–427 (2015). URL [http://onlinelibrary.wiley.com/doi/10.1002/](http://onlinelibrary.wiley.com/doi/10.1002/iub.1387/abstract)
517 [iub.1387/abstract](http://onlinelibrary.wiley.com/doi/10.1002/iub.1387/abstract).
- 518 [40] McEvoy, C. R., Helden, P. D., Warren, R. M. & Pittius, N. C. G. Ev-
519 idence for a rapid rate of molecular evolution at the hypervariable and
520 immunogenic Mycobacterium tuberculosis PPE38 gene region. *BMC Evo-*

- 521 *lutionary Biology* **9**, 237 (2009). URL <http://dx.doi.org/10.1186/>
522 1471-2148-9-237.
- 523 [41] Liu, W. *et al.* Comparative genomic analyses of *Mycoplasma hyopneu-*
524 *moniae* pathogenic 168 strain and its high-passaged attenuated strain.
525 *BMC Genomics* **14**, 80 (2013). URL [http://www.ncbi.nlm.nih.gov/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3626624/)
526 [pmc/articles/PMC3626624/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3626624/).
- 527 [42] Zhang, W. *et al.* Genome Sequencing and Analysis of BCG Vaccine
528 Strains. *PLoS ONE* **8** (2013). URL [http://www.ncbi.nlm.nih.gov/pmc/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3747166/)
529 [articles/PMC3747166/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3747166/).
- 530 [43] Zhang, S. *et al.* Mutations in panD encoding aspartate decarboxylase are
531 associated with pyrazinamide resistance in *Mycobacterium tuberculosis*.
532 *Emerging Microbes & Infections* **2**, e34 (2013). URL [http://www.ncbi.](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3697303/)
533 [nlm.nih.gov/pmc/articles/PMC3697303/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3697303/).
- 534 [44] Song, T. *et al.* Fitness costs of rifampicin-resistance in *Mycobacterium*
535 *tuberculosis* are amplified under conditions of nutrient starvation and com-
536 pensated by mutation in the beta' subunit of RNA polymerase. *Molecu-*
537 *lar microbiology* **91**, 1106–1119 (2014). URL [http://www.ncbi.nlm.nih.](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3951610/)
538 [gov/pmc/articles/PMC3951610/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3951610/).
- 539 [45] Freidlin, P. J., Goldblatt, D., Kaidar-Shwartz, H. & Rorman, E. Poly-
540 morphic Exact Tandem Repeat A (PETRA): a Newly Defined Lineage of
541 *Mycobacterium tuberculosis* in Israel Originating Predominantly in Sub-
542 Saharan Africa. *Journal of Clinical Microbiology* **47**, 4006–4020 (2009).
543 URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2786624/>.
- 544 [46] Namouchi, A., Karboul, A., Fabre, M., Gutierrez, M. C. & Mardassi, H.
545 Evolution of Smooth Tubercle Bacilli PE and PE_{pgrs} Genes: Evidence
546 for a Prominent Role of Recombination and Imprint of Positive Selec-

547 tion. *PLoS ONE* **8** (2013). URL <http://www.ncbi.nlm.nih.gov/pmc/>
548 [articles/PMC3660525/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3660525/).

549 [47] Banerjee, R., Vats, P., Dahale, S., Kasibhatla, S. M. & Joshi, R. Com-
550 parative Genomics of Cell Envelope Components in Mycobacteria. *PLoS*
551 *ONE* **6** (2011). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/>
552 [PMC3089613/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3089613/).

553 [48] Zhang, G. *et al.* Evaluation of mycobacterial virulence using rabbit skin
554 liquefaction model. *Virulence* **1**, 156–163 (2010). URL [http://dx.doi.](http://dx.doi.org/10.4161/viru.1.3.11748)
555 [org/10.4161/viru.1.3.11748](http://dx.doi.org/10.4161/viru.1.3.11748).

556 [49] Målen, H., De Souza, G. A., Pathak, S., Søfteland, T. & Wiker, H. G.
557 Comparison of membrane proteins of Mycobacterium tuberculosis H37rv
558 and H37ra strains. *BMC Microbiology* **11**, 18 (2011). URL [http://www.](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3033788/)
559 [ncbi.nlm.nih.gov/pmc/articles/PMC3033788/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3033788/).

560 [50] Groenewald, W., Baird, M. S., Verschoor, J. A., Minnikin, D. E. &
561 Croft, A. K. Differential spontaneous folding of mycolic acids from My-
562 cobacterium tuberculosis. *Chemistry and Physics of Lipids* **180**, 15–22
563 (2014). URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0009308413001667)
564 [S0009308413001667](http://www.sciencedirect.com/science/article/pii/S0009308413001667).

565 [51] Minnikin, D. E. *et al.* Pathophysiological Implications of Cell Envelope
566 Structure in Mycobacterium tuberculosis and Related Taxa. In Ribon,
567 W. (ed.) *Tuberculosis - Expanding Knowledge* (InTech, 2015). URL [http:](http://www.intechopen.com/books/tuberculosis-expanding-knowledge/pathophysiological-implications-of-cell-envelope-structure-in-mycobacterium-tuberculosis)
568 [//www.intechopen.com/books/tuberculosis-expanding-knowledge/](http://www.intechopen.com/books/tuberculosis-expanding-knowledge/pathophysiological-implications-of-cell-envelope-structure-in-mycobacterium-tuberculosis)
569 [pathophysiological-implications-of-cell-envelope-structure-in-mycobacterium-tuberculosis](http://www.intechopen.com/books/tuberculosis-expanding-knowledge/pathophysiological-implications-of-cell-envelope-structure-in-mycobacterium-tuberculosis)

570 [52] Hunt, M. *et al.* Circlator: automated circularization of genome assemblies
571 using long sequencing reads. *Genome Biology* **16** (2015). URL [http:](http://genomebiology.com/2015/16/1/294)
572 [//genomebiology.com/2015/16/1/294](http://genomebiology.com/2015/16/1/294).

- 573 [53] English, A. C., Salerno, W. J. & Reid, J. G. PBHoney: identifying ge-
574 nomic variants via long-read discordance and interrupted mapping. *BMC*
575 *bioinformatics* **15**, 180 (2014). URL [http://www.biomedcentral.com/](http://www.biomedcentral.com/1471-2105/15/180)
576 [1471-2105/15/180](http://www.biomedcentral.com/1471-2105/15/180).
- 577 [54] Myers, E. W. An $o(ND)$ difference algorithm and its variations. *Algorith-*
578 *mica* **1**, 251–266 (1986). URL [http://link.springer.com/article/10.](http://link.springer.com/article/10.1007/BF01840446)
579 [1007/BF01840446](http://link.springer.com/article/10.1007/BF01840446).
- 580 [55] Miller, W. & Myers, E. W. A file comparison program. *Software: Prac-*
581 *tice and Experience* **15**, 1025–1040 (1985). URL [http://onlinelibrary.](http://onlinelibrary.wiley.com/doi/10.1002/spe.4380151102/abstract)
582 [wiley.com/doi/10.1002/spe.4380151102/abstract](http://onlinelibrary.wiley.com/doi/10.1002/spe.4380151102/abstract).
- 583 [56] McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biology*
584 **17**, 122 (2016). URL <http://dx.doi.org/10.1186/s13059-016-0974-4>.

585 **Figures**

Figure 1: Example Classification of Genes Based on Variant Comparisons.

Considering the profile of H37Ra-specific variants (those with respect to H37Rv not also appearing in CDC1551), a given gene (blue arrow) is categorized as “supported”, “contradicted”, or “adjusted” by our H37Ra assembly as a result of comparison with the hitherto reference sequence NC_009525.1. The illustration shows examples of the different variant profiles a gene could have and their resulting classifications. Genes in the “supported” and “contradicted” categories are strictly those where our assembly either fully matches the H37Ra reference (supported) or the H37Rv reference (contradicted). Multiple factors may cause a gene to be classified as “adjusted”. Such genes may have variant profiles not fully meeting the criteria of “supported” or “contradicted”, or they may have novel H37Ra-specific variants observed only in our assembly.

Figure 2: Cell Wall Differences in H37Ra and H37Rv

A) State of knowledge following publication of H37RaJH. At this time it was known that the SNP in the DNA-binding site of *phoP* abrogated synthesis of sulfolipids (yellow) and acyltrehaloses (purple and red) of the mycomembrane outer leaflet, while two SNPs in *pks12*, both of which were refuted in our assembly, were thought to cause the observed lack of phthiocerol dimycocerosates (blue) in H37Ra. **B) Current state of knowledge.** Advances were made in understanding the inner leaflet. A single nucleotide, frameshift deletion in the now annotated *hadC* gene was shown by Slama and colleagues [33] to alter the mycolic acid profile in three distinct ways: i. Lower proportion of oxygenated mycolic acids (K-MA and Me-MA; green and blue carbon skeletons, respectively) to α -MAs (orange carbon skeleton). There are seven Me-MAs depicted in H37Rv compared to three in H37Ra, reflecting the proportions reported by Slama and colleagues [33]. ii. Extra degree of unsaturation (red circles) in H37Ra mycolic acids due to truncation of the HadC protein in H37Ra. iii. Shorter chain lengths of mycolic acids in H37Ra. Note that Me-MAs have larger loops in H37Rv than in H37Ra, and that the height of the α -MAs is shorter in H37Ra than H37Rv. Carbon chain lengths are based on results reported by Slama and colleagues. The folding geometry of the mycolic acids is depicted in panel B, as described by Groenewald and colleagues [50], and inspired by the illustration style of Minnikin and colleagues [51].

Figure 3: Visualization of the Reduced Set of H37Ra-specific Variants and Their Effect on Phenotype

Our assembly contradicts many variants previously thought to be H37Ra-specific, reducing the number of genes that may contribute to H37Ra's virulence attenuation. Several of these genes have been reassigned function since the first published assembly of the H37Ra genome [5], which is reflected in the figure. Blue stars indicate that the H37Ra-specific variant(s) in that gene has been shown to confer a phenotypic change in H37Ra relative to H37Rv in wet-lab studies. For these genes, the mechanisms affected by the H37Ra-specific variant are illustrated in detail (see Figure 2 for *hadC* and *phoP*). For other genes, their general function is described or briefly illustrated. a) The set of genes identified to carry H37Ra-specific polymorphisms in the original H37Ra genome publication [5] and their contribution to phenotype as understood at that time. 57 genes are affected, the majority of which were PE_PPE genes or were of unknown function. b) The set of genes with H37Ra-specific variants confirmed by our assembly is reduced markedly, particularly in PE_PPE genes, highlighting the strength of single-molecule sequencing in resolving GC-rich and repetitive stretches of DNA. Genes with functions not yet characterized were also reduced significantly. Though in a few instances this was because these genes' function was characterized between 2008 and now (indicated by an asterisk), most were due to our assembly showing that they matched that of H37Rv and, therefore, are not H37Ra-specific.

586 **Tables**

Table 1: Status of Genes Previously Reported as Affected by H37Ra-specific Mutations.

- (a) Genes with all High-Confidence Variants Unsupported by our Assembly
- (b) Genes with Different H37Ra-specific Variant Profiles in our Assembly
- (c) Genes with High-Confidence Variant Profiles Fully Confirmed by our Assembly
- (d) Genes with Variant Profiles Erroneously Declared as H37Ra-specific

Table 2: Variants in H37Ra Unique to our Assembly

587 **Supplementary Information**

588 **Supplementary Note — Expanded Discussion of Virulence**
589 **Attenuation Mechanisms in *M. tuberculosis* H37Ra**

590 **Supplementary Data 1 — Raw Variants**

591 Zip archive containing the following data in Variant Call Format (VCF):

592 **A6_7-H37Ra_NC009525_1.vcf** Variants in our H37Ra assembly with respect
593 to the H37Ra reference sequence.

594 **A6_7-H37Rv_NC000962_3.vcf** Variants in our H37Ra assembly with respect
595 to the H37Rv reference sequence.

596 **H37Ra_NC009525_1-H37Rv_NC000962_3.vcf** Variants in the H37Ra ref-
597 erence sequence with respect to the H37Rv reference sequence.

598 **Supplementary Data 2 — Annotated Variants with Respect**
599 **to H37Rv**

600 Spreadsheet containing annotated variants in our assembly and the H37Ra ref-
601 erence sequence with respect to the H37Rv reference sequence. The sheets
602 separate variants that are common to the two H37Ra assemblies and those that
603 are unique to each.

604 **Supplementary Data 3 — Computer Code used for Analy-**
605 **ses**

606 **Online Methods**

607 **Sample Preparation and Whole-Genome Sequencing**

608 *M. tuberculosis* H37Ra (ATCC25177) was obtained from ATCC and cultured
609 on Lowenstein-Jenson slants and Middlebrook 7H11 plates. Cultures were incu-
610 bated until growth of a full bacterial lawn. DNA was extracted using Genomic-
611 tips (Qiagen Inc.) following the manufacturer's sample preparation and lysis
612 protocol for bacteria with the following modifications. Culture was harvested
613 directly into buffer B1/RNase solution, homogenized by vigorous vortex mixing
614 and inactivated at 80°C for 15 minutes. Lysozyme was added and incubated
615 at 37°C for 30 minutes followed by the addition of proteinase K and further
616 incubation at 37°C for an additional 60 minutes. Buffer B2 was added and the
617 mixture was incubated overnight at 50°C. Wide-bore pipet tips were used to
618 optimize recovery of large DNA fragments. The remainder of the Genomic-tip
619 protocol was carried out exactly as described by the manufacturer. DNA pu-
620 rity and concentration was analyzed on a Nanodrop 1000 (Thermo Scientific).
621 The DNA was then sequenced using two SMRTCells on the Pacific Biosciences
622 RS II instrument with the P6-C4 chemistry and a 20kb insert library prepara-
623 tion.

624 **Genome Assembly and Methylation Determination**

625 The genome was assembled using Pacific Biosciences' Hierarchical Genome As-
626 sembly Process [12] (HGAP) as implemented in SMRTAnalysis 2.3.0. This ver-
627 sion of SMRTAnalysis provides two implementations of HGAP: HGAP.2 and
628 the newer HGAP.3. HGAP.3 differs from HGAP.2 by replacing the Celera
629 Assembler's assembly consensus step with Pacific Biosciences' speed-optimized

630 implementation. We used HGAP.2 because, in our experiments, we found that
631 HGAP.3 consistently produced spurious contigs while HGAP.2 did not.

632 The overlapping ends of the contig, an artifact of the assembly due to the cir-
633 cularity of the chromosome, were trimmed and joined using the minimus2 pro-
634 gram from AMOS (<http://amos.sourceforge.net>). Discrepancies between
635 the contig ends were resolved manually by selecting an authoritative sequence
636 and trimming the discrepant one. The circularization was also performed with
637 Circlator [52] to confirm the minimus2 results.

638 We validated the assembly structure using PBHoney [53], a structural varia-
639 tion detection tool, by using our assembled genome as input. Any structural
640 variations detected would indicate potential misassemblies.

641 Circularization was followed by three rounds of assembly polishing using Quiver
642 in SMRTAnalysis. Quiver was used with the maximum coverage parameter set
643 to 1000 and otherwise default settings.

644 The methylome was determined using the base modification and motif analysis
645 protocol in SMRTAnalysis.

646 **Comparative Genomics**

647 In all cases, variants were determined using GNU diff ([http://www.gnu.org/](http://www.gnu.org/software/diffutils)
648 [software/diffutils](http://www.gnu.org/software/diffutils)), an implementation of Myers' algorithm for solving the
649 longest-common-subsequence problem [54, 55] and converted to the Variant
650 Call Format for analysis. This process is implemented in our custom tool,
651 biodiff (<http://www.github.com/valafarlab/biodiff>). Because insertions
652 and deletions in repetitive regions can be represented equivalently in multi-
653 ple ways, we normalized the variants using the "norm" function of bcftools
654 (<http://samtools.github.io/bcftools>), giving every mutation a standard

655 representation to facilitate a proper comparison. Variants were then compared
656 using bcftools isec and annotated using the Ensembl Variant Effect Predic-
657 tor [56]. Motif variants were analyzed using *in villa* code.

658 **Literature Review**

659 In order to gain a holistic view of the research built off of and conclusions drawn
660 from the unique variants of H37RaJH with respect to H37Rv, we performed an
661 exhaustive literature review. Common names and Rv numbers were searched
662 using Google scholar within all publications which cited Zheng et al, 2008 [5]
663 as of March 14th, 2016, for all genes with H37RaJH specific variants within
664 CDS or potential promoter regions, according to Table 2 of [5]. All mentions of
665 these genes were compiled and evaluated to illustrate how our assembly alters
666 the picture of how the genomic differences between the reference strains con-
667 tribute to the observed virulence attenuation of H37Ra (Discussion). Genes are
668 discussed in the present study according to the H37Rv annotation (as opposed
669 to H37Ra's own annotation), as this convention relates to extant publications
670 most readily.

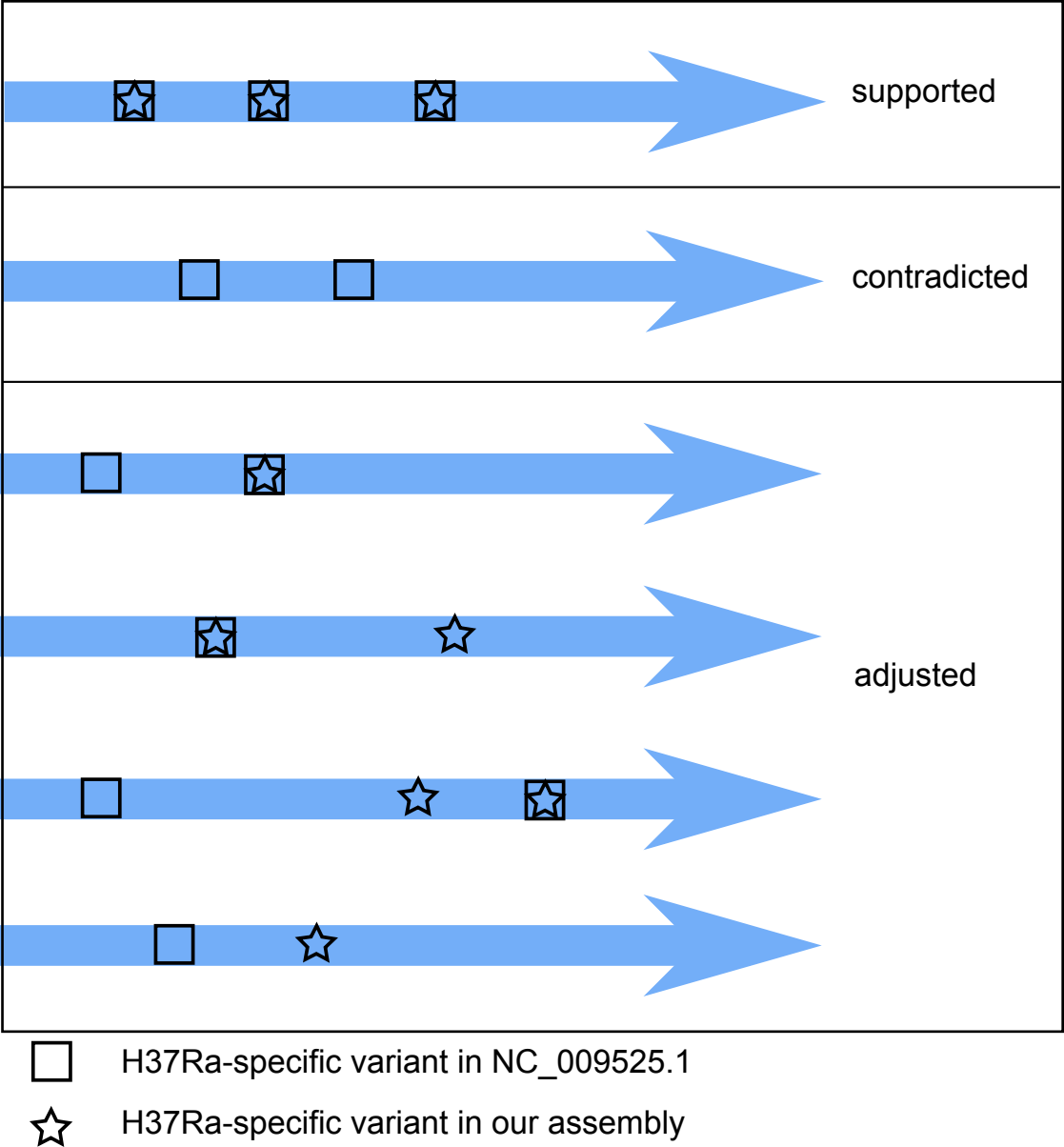
671 **Data Availability**

672 The assembled sequence and raw sequencing data for this project are available
673 through NCBI under Bioproject PRJNA329548.

674 **Code Availability**

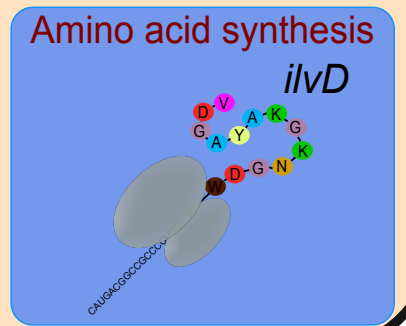
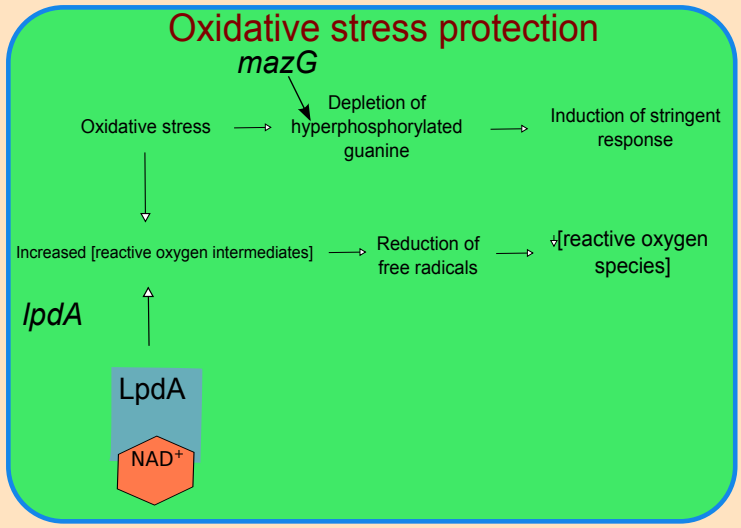
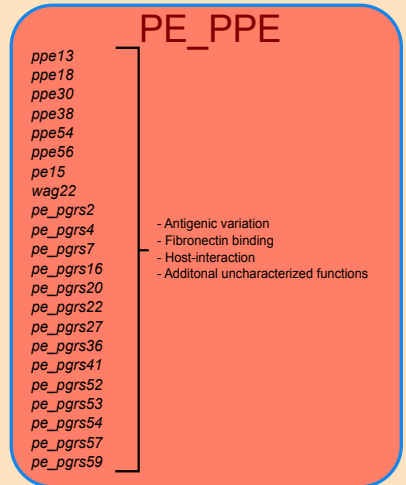
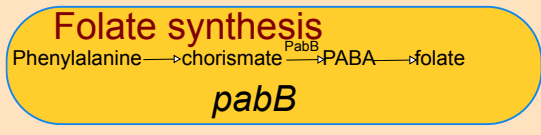
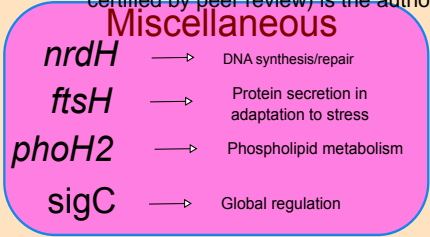
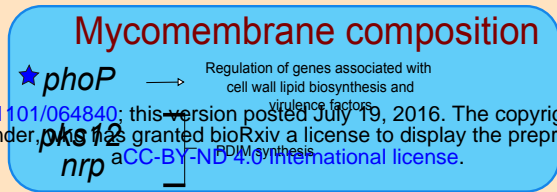
675 Our motif variants detection tool is available from [http://github.com/valafarlab/](http://github.com/valafarlab/motif-variants)
676 `motif-variants`. Analysis code for this study is provided as Supplemental Data

677 3.

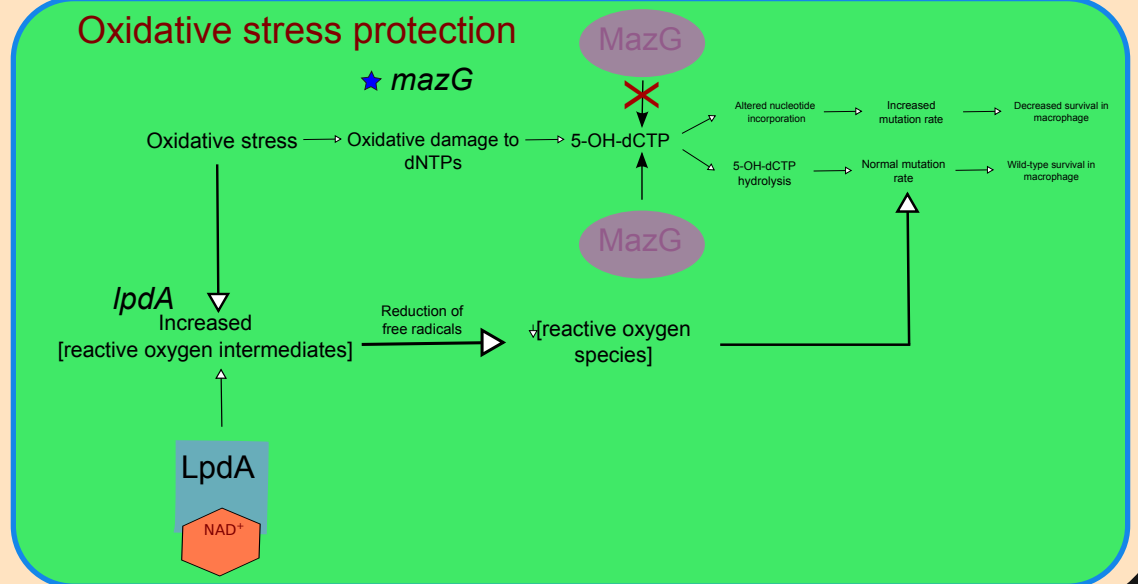
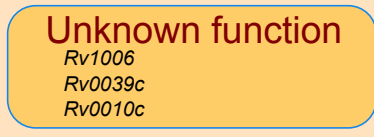
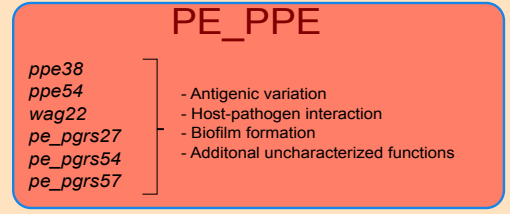
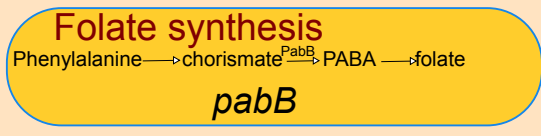
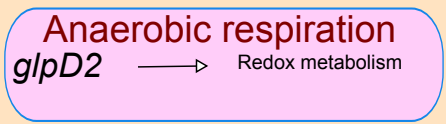
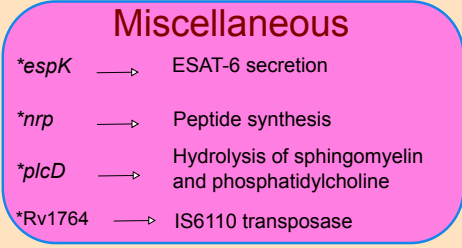
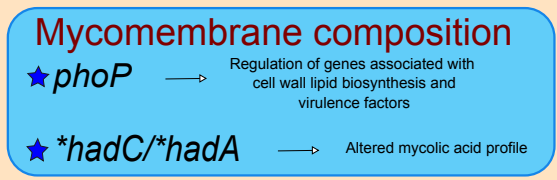


a

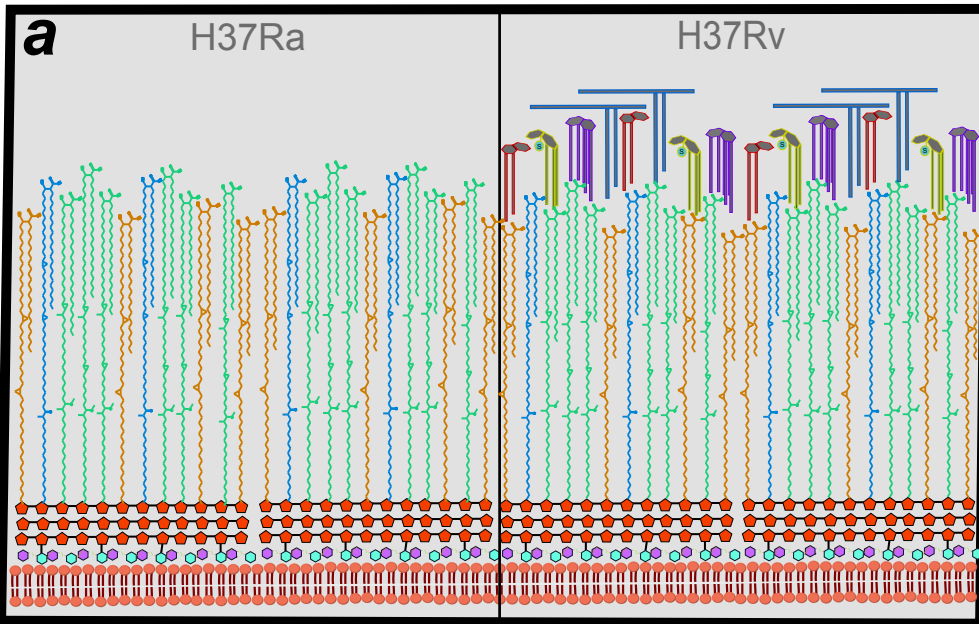
bioRxiv preprint doi: <https://doi.org/10.1101/064840>; this version posted July 19, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.







b



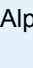


State of Knowledge after H37RaJH (2008)



Trehaloses & Dimycoserate (outer leaflet)

-  Polyacyltrehalose (PAT)
-  Sulfolipids (SL)
-  Diacyltrehalose (DAT)
-  Phthiocerol Dimycoserate (PDIM)

Mycolic acids (inner leaflet)

-  Alpha-mycolic acid (α-MA)
-  Methoxy-mycolic acid (Me-MA)
-  Keto-mycolic acid (K-MA)

State of Knowledge after H37RaSD (2016)

