

1 **The role of epistasis in amikacin, kanamycin, bedaquiline, and clofazimine resistance in**
2 ***Mycobacterium tuberculosis* complex**

3

4 Roger Vargas Jr^{1,2*}, Luca Freschi², Andrea Spitaleri³, Sabira Tahseen⁴, Ivan Barilar^{5,6}, Stefan
5 Niemann^{5,6}, Paolo Miotto³, Daniella Maria Cirillo³, Claudio U. Köser⁷, and Maha R. Farhat^{2,8*}

6

7 ¹ Department of Systems Biology, Harvard Medical School, Boston, USA.

8 ² Department of Biomedical Informatics, Harvard Medical School, Boston, USA.

9 ³ Emerging Bacterial Pathogens Unit, IRCCS San Raffaele Scientific Institute, Milan, Italy.

10 ⁴ National TB Reference laboratory, National TB Control Program, Islamabad, Pakistan.

11 ⁵ German Center for Infection Research, Partner site Hamburg-Lübeck-Borstel-Riems, Borstel, Germany.

12 ⁶ Molecular and Experimental Mycobacteriology, Research Center Borstel, Borstel, Germany.

13 ⁷ Department of Genetics, University of Cambridge, Cambridge, UK.

14 ⁸ Pulmonary and Critical Care Medicine, Massachusetts General Hospital, Boston, USA.

15 *Corresponding authors: roger_vargas@g.harvard.edu, Maha_Farhat@hms.harvard.edu

16 **ABSTRACT**

17 Antibiotic resistance among bacterial pathogens poses a major global health threat. *M. tuberculosis*
18 complex (MTBC) is estimated to have the highest resistance rates of any pathogen globally. Given
19 the slow growth rate and the need for a biosafety level 3 laboratory, the only realistic avenue to
20 scale up drug-susceptibility testing (DST) for this pathogen is to rely on genotypic techniques.
21 This raises the fundamental question of whether a mutation is a reliable surrogate for phenotypic
22 resistance or whether the presence of a second mutation can completely counteract its effect,
23 resulting in major diagnostic errors (i.e. systematic false resistance results). To date, such epistatic
24 interactions have only been reported for streptomycin that is now rarely used. By analyzing more
25 than 31,000 MTBC genomes, we demonstrated that *eis* C-14T promoter mutation, which is
26 interrogated by several genotypic DST assays endorsed by the World Health Organization, cannot
27 confer resistance to amikacin and kanamycin if it coincides with loss-of-function (LoF) mutations
28 in the coding region of *eis*. To our knowledge, this represents the first definitive example of
29 antibiotic reversion in MTBC. Moreover, we raise the possibility that *mmpR* (*Rv0678*) mutations
30 are not valid markers of resistance to bedaquiline and clofazimine if these coincide with LoF
31 mutation in the efflux pump encoded by *mmpS5* (*Rv0677c*) and *mmpL5* (*Rv0676c*).

32 INTRODUCTION

33 Tuberculosis (TB) and its causative pathogen *Mycobacterium tuberculosis* complex (MTBC) is a
34 major public health threat causing an estimated 10 million new cases of disease per year (World
35 Health Organization, 2020). Antibiotic resistance in particular poses a problem to controlling the
36 TB epidemic (World Health Organization, 2020). Owing to the inherently slow growth rate of
37 MTBC, genotypic drug-susceptibility testing (DST) represents the only realistic option to inform
38 the initial selection of the most appropriate treatment regimen (Mohamed et al., 2021). This raises
39 the fundamental question of whether the effect and clinical interpretation of a marker for resistance
40 depends on the presence of another mutation (i.e. epistasis) or whether the effect is universal.

41 Although it is known that the level of resistance conferred by resistance mutations in some
42 genes can differ, the only well-understood epistatic mechanism that completely counteracts the
43 effect of another mutation involves the *whiB7* (*Rv3197A*) regulatory gene (Ajileye et al., 2017;
44 Castro et al., 2020; Gagneux, 2018). Specifically, the over-expression of the *whiB7* cannot confer
45 streptomycin resistance in the vast majority of lineage 2 isolates because these have a loss-of-
46 function (LoF) mutation in the *tap* (*Rv1258c*) efflux pump (Köser et al., 2013; Merker et al., 2020).
47 Yet, because the use of streptomycin has been downgraded in the most recent treatment guidelines
48 by the World Health Organization (WHO), the clinical relevance of this example is limited (Viney
49 et al., 2021).

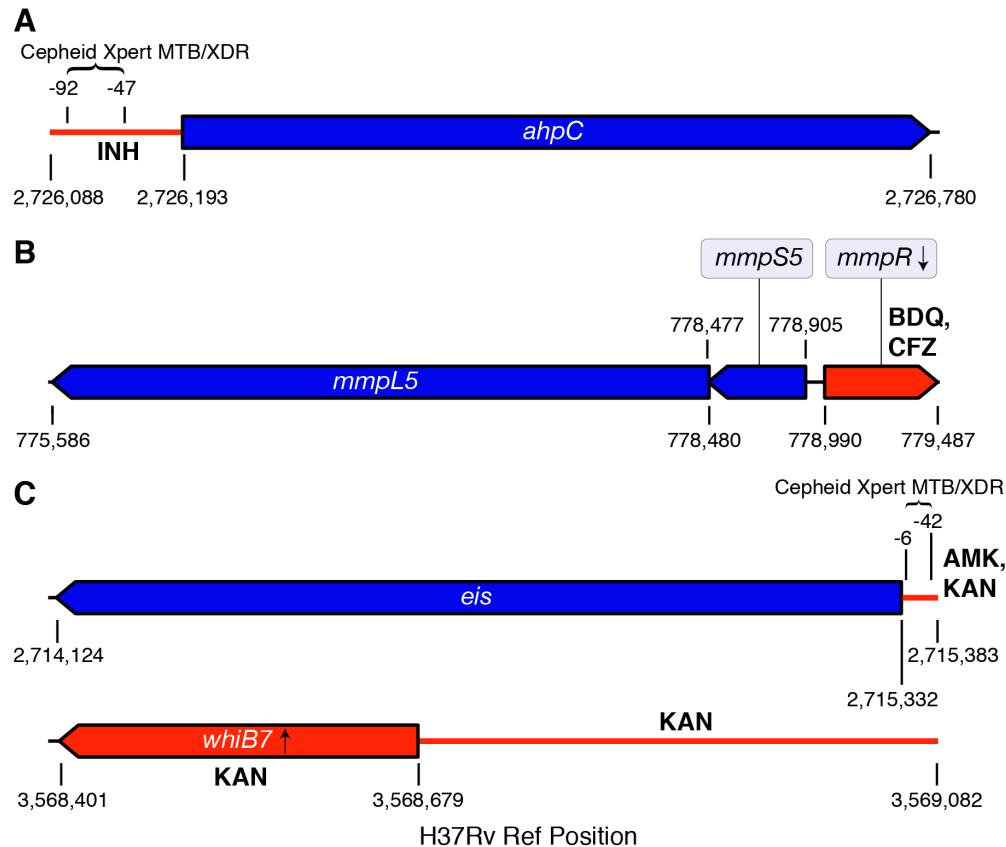
50 Using whole-genome sequencing (WGS) data for 31,440 isolates, we set out to survey
51 systematically whether other markers of resistance to more important antibiotics may be affected
52 by epistasis if they involve the over-expression of a non-essential drug resistance gene. First, we
53 analyzed the alkyl-hydroperoxidase *ahpC* (*Rv2428*), the function of which is not fully elucidated
54 but may act as a compensatory mechanism for isoniazid (INH) resistance caused by *katG*

55 mutations (i.e. the soon-to-be WHO-endorsed Cepheid Xpert MTB/XDR assay interrogates *ahpC*
56 promoter/upstream mutations) (World Health Organization, In press). Second, LoF mutations in
57 the transcriptional repressor *mmpR* (*Rv0678*), which is sequenced by several commercial targeted
58 next-generation sequencing (tNGS) assays being evaluated by WHO, confer cross-resistance to
59 bedaquiline (BDQ) and clofazimine (CFZ) via the over-expression of the non-essential efflux
60 pump encoded by *mmpS5* (*Rv0677c*) and *mmpL5* (*Rv0676c*) (Kadura et al., 2020; Mohamed et al.,
61 2021; Viljoen et al., 2017; Yamamoto et al., 2021). Third, four promoter/upstream mutations for
62 the *eis* (*Rv2416c*) acetyltransferase are responsible for kanamycin (KAN) resistance, of which the
63 C-14T mutation is due to be recognized by WHO as conferring cross-resistance to amikacin
64 (AMK), the only aminoglycoside (AG) now recommended for the treatment of TB (World Health
65 Organization, 2018, In press). In fact, the Xpert MTB/XDR already interprets this *eis* mutation
66 accordingly, whereas the WHO-endorsed Hain GenoType MTBDRsl VER 2.0 (SL-LPA) assay
67 will have to be updated accordingly. Finally, we included *whiB7* as it also regulates *eis* and,
68 therefore, could theoretically confer cross-resistance to both AGs rather than just to KAN (World
69 Health Organization, 2018, In press).

70 **RESULTS**

71 **INH: *ahpC* upstream mutations in combination with *ahpC* LoF mutations**

72 We observed 57 unique single nucleotide polymorphisms (SNPs) in the upstream region of *ahpC*
73 (**Fig 1A**), of which 18 were homoplasic and occurred in at least five isolates, consistent with
74 parallel evolution and known selection on this gene (**Table 1, Supplementary Table 1**). We
75 screened for frameshift indels, nonsense mutations, and mutations that abolish the start codon of
76 *ahpC* given that these are the most likely types of mutations to confer a LoF phenotype. This
77 yielded seven unique variants in eight isolates, of which just *ahpC* 323delC co-occurred with an
78 upstream mutation in a single isolate (**Table 2, Supplementary Table 1**). This particular upstream
79 mutation (i.e. G-88A **Table 1, Supplementary Table 1**) is a marker for MTBC lineage 3 and
80 correlates with only a 3-fold increase in the expression of *ahpC*, potentially by creating a new
81 Pribnow box (Chiner-Oms et al., 2019; Merker et al., 2020). As a result, this SNP is not considered
82 to be a marker of resistance (i.e. the Xpert MTB/XDR was designed not to detect it, unlike adjacent
83 mutations), which means that this is not an example of epistasis (World Health Organization, In
84 press). Indeed, this double mutant was phenotypically susceptible to INH at the critical
85 concentration (CC) of 0.1 mg/L in MGIT 960.



86

87 **Fig. 1. Genomic regions interrogated.** The non-essential genes conferring resistance or
88 compensating for resistance are shown in blue, whereas the corresponding regulatory regions or
89 non-essential regulators are shown in red along with the relevant antibiotic(s). For each of the four
90 regions, we screened for any type of mutation in the upstream regions and likely LoF mutations in
91 the coding regions (i.e. frameshift indels, nonsense mutations, and mutations that abolish the start
92 codon, including synonymous mutations at the start codons of *eis*, *mmpR*, and *whiB7* as these
93 genes start with a valine). Unlike Cepheid, Hain has not disclosed the precise *eis* promoter region
94 interrogated by its WHO-endorsed SL-LPA, which is why this information could not be included
95 (Hain Lifescience, 2017).

Position	Variant Name	Gene Position	Mutation Type	Codon Position	# Isolates	# Sublineages
776210	<i>mmpL5</i> C2271A	2271	N	Y757*	2	2
777499	<i>mmpL5</i> C982T	982	N	R328*	2	2
777581	<i>mmpL5</i> C900A	900	N	Y300*	293	2
778086	<i>mmpL5</i> 395insC	395	ins	132	6	2
779127	<i>mmpR</i> 138insG	138	ins	46	5	4
779181	<i>mmpR</i> 192-198delG	192	del	64	20	4
779181	<i>mmpR</i> 192-198insG	192	ins	64	86	2
779407	<i>mmpR</i> 418insG	418	ins	140	6	1
2714753	<i>eis</i> C580T	580	N	Q194*	10	2
2715287	<i>eis</i> 46insC	46	ins	16	3	2
2715305	<i>eis</i> G28T	28	N	E10*	6	2
2715330	<i>eis</i> G3A	3	S	V1V	6	2
2715342	<i>eis</i> G-10A	-10 <i>eis</i>	I		293	18
2715344	<i>eis</i> C-12T	-12 <i>eis</i>	I		332	17
2715346	<i>eis</i> C-14T	-14 <i>eis</i>	I		181	19
2715369	<i>eis</i> G-37T	-37 <i>eis</i>	I		285	8
2726105	<i>ahpC</i> G-88A	-88 <i>ahpC</i>	I		3350	12
2726141	<i>ahpC</i> C-52A	-52 <i>ahpC</i>	I		91	10
2726141	<i>ahpC</i> C-52T	-52 <i>ahpC</i>	I		92	25
2726145	<i>ahpC</i> G-48A	-48 <i>ahpC</i>	I		85	20
3568487	<i>whiB7</i> 193insG	193	ins	65	3	3
3568488	<i>whiB7</i> 192delC	192	del	64	573	3
3568547	<i>whiB7</i> 133delCA	133	del	45	2	2
3568626	<i>whiB7</i> 54delA	54	del	18	61	2
3568779	<i>whiB7</i> T-100C	-100 <i>whiB7</i>	I		256	2
3568857	<i>whiB7</i> C-178T	-178 <i>whiB7</i>	I		73	2
3568921	<i>whiB7</i> G-242C	-242 <i>whiB7</i>	I		117	2
3569029	<i>whiB7</i> A-350G	-350 <i>whiB7</i>	I		249	3

96

97 **Table 1. Mutations detected in a global sample of MTBC clinical isolates.** Mutations that occur

98 in our sample of 31,440 clinical isolates within the *mmpL5*, *mmpS5*, *mmpR*, *ahpC*, *eis*, *whiB7*
99 coding sequences and *oxyR-ahpC*, *eis-Rv2417c*, *whiB7-uvrD2* intergenic regions (**Figure 1**).

100 Mutations in **regulator** regions (*mmpR*, *oxyR-ahpC*, *eis-Rv2417c*, and *whiB7-uvrD2*) reported in
101 this table were among the four most commonly detected variants in each region. Mutations in

102 **regulated** regions (*mmpL5*, *mmpS5*, *ahpC*, *eis*, and *whiB7*) reported in this table were present in

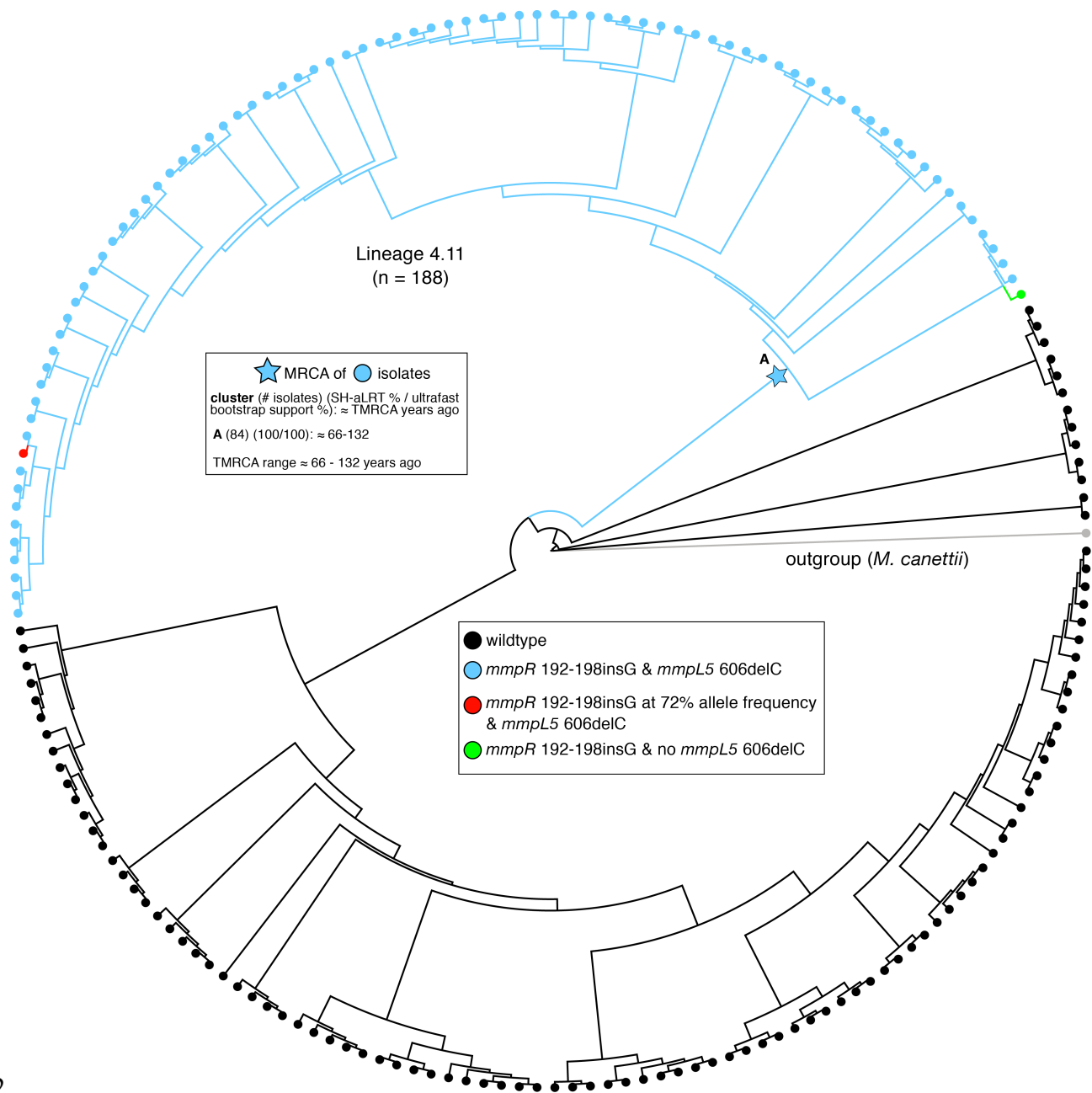
103 at least two MTBC sub-lineages. The full set of mutations detected within these genomic regions

104 is reported in **Supplementary Table 1** (S=Synonymous, N=Non-Synonymous, I=Intergenic).

105

106 **BDQ/CFZ: *mmpR* LoF mutations in combination with *mmpS5-mmpL5* LoF mutations**

107 We detected 91 fixed LoF variants in *mmpL5-mmpS5-mmpR*, of which 35 occurred in at least two
108 isolates (**Fig 1B, Supplementary Table 1**). Frameshifts were most common (39/68) in *mmpL5*,
109 followed by *mmpR* (21/68), and *mmpS5* (8/68). Each gene harbored frameshifts in isolates from at
110 least three MTBC major lineages, indicating parallel evolution (**Supplementary Table 1**). The
111 nonsense SNP *mmpL5* Y300* was observed in 293 isolates and in two genetically distinct lineages,
112 and the insC at nt395 of *mmpL5* also occurred two in distant lineages (**Table 1, Supplementary**
113 **Table 1**). The *mmpR* delG in the homopolymer (HP) nt192-198 was observed in 20 isolates from
114 three major lineages, whereas insG in the same HP was observed in 86 isolates from two major
115 lineages (**Table 1, Supplementary Table 1**). Noting the frequency of frameshifts in the
116 homopolymer region of *mmpR*, we investigated non-fixed frameshift variants (*i.e.* that had within-
117 sample allele frequencies of 10-75%) and recorded which isolates had >100x coverage of *mmpR*,
118 *mmpS5*, and *mmpL5*. Frameshift variants at low to intermediate allele frequency were rare and
119 occurred in a total of six isolates (3/7435 in *mmpR*, 1/8949 in *mmpS5*, and 2/6217 in *mmpL5*). Two
120 of these isolates had the frameshift insG in the aforementioned *mmpR* HP at 66% and 71% allele
121 frequencies (**Fig. 2, Supplementary Table 2**).



122

123 **Fig. 2. Phylogeny of 188 sub-lineage 4.11 isolates.** Isolates with both the *mmpR* 192-198insG
124 and *mmpL5* 606delC variants are colored in blue (n=82), isolates carrying neither variant are
125 colored in black (n=104). One isolate had the *mmpR* 192-198insG frameshift but not the *mmpL5*
126 606delC frameshift (green). Another isolate had the *mmpL5* 606delC frameshift and 72% of reads
127 supporting the *mmpR* 192-198insG frameshift (red). Time to most recent common ancestor
128 (TMRCA) estimates for the group of isolates with the *mmpR* 192-198insG and *mmpL5* 606delC
129 variants is given in the upper left.

130

131 Three different LoF mutations in *mmpR* coincided with a LoF mutation in *mmpL5* (**Fig. 2**).

132 Of those, insG in the HP nt192-198, which had been repeatedly demonstrated to confer BDQ and

133 CFZ resistance during *in vitro* selection experiments and patient treatment, occurred in 82 isolates,

134 whereas the other two double mutations were observed in only a single isolate, respectively

135 (Andres et al., 2020; de Vos et al., 2019; Ghodousi et al., 2019; Peretokina et al., 2020; Sonnenkalb

136 et al., 2021; Zhang et al., 2015). All of the former 82 double-LoF mutants belonged to a

137 monophyletic group within sub-lineage 4.11 that was mostly multi-drug resistant (53 of 59 with

138 known phenotypic data). Most double-LoF mutants were isolated in Lima, Peru, between 1997

139 and 2012 and represented 43% (82/188) of the isolates from the sub-lineage 4.11 in our dataset

140 (**Fig. 2, Table 2, Supplementary Table 3-4**). Among the 84 isolates with co-occurrence of *mmpR*

141 and *mmpL5* LoF, there were no SNPs in the other BDQ resistance locus, *atpE*.

A: Mutation in Regulator				B: Mutation in Regulated Gene				isolate info with co-occurring mutations	
A type	A variant-name	A codon -pos	A num-isolates	B type	B variant-name	B codon -pos	B num-isolates	# isolates both mutations	isolate sublineages
INDEL	<i>mmpR</i> 192-198delG	64	20	INDEL	<i>mmpL5</i> 2028insA	676	1	1	2.2.1.1.1(1)
INDEL	<i>mmpR</i> 192-198insG	64	86	INDEL	<i>mmpL5</i> 606delC	202	83	82	4.11(82)
INDEL	<i>mmpR</i> 207insA	69	1	INDEL	<i>mmpL5</i> 1160insCGATG	387	1	1	2.2.1(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 627insC	209	1	1	2.2.1.1.1.i3(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 486insCT	162	2	2	4.1.i1.1.1(2)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 473insT	158	1	1	2.2.1.1.1.i3(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 448delA	150	7	7	2.2.1.1.1.i3(7)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 400insG	134	1	1	2.2.1.1.1(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 279delCGGCGATGCGT	93	1	1	2.2.1.1.1.i3(1)
SNP	<i>eis</i> C-14T		181	SNP	<i>eis</i> G39A	W13*	1	1	2.2.1.1.1(1)
SNP	<i>eis</i> C-14T		181	SNP	<i>eis</i> G38A	W13*	1	1	2.2.1.1.1(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 16insC	6	1	1	2.2.1.1.1(1)
SNP	<i>eis</i> C-14T		181	INDEL	<i>eis</i> 15insC	5	1	1	2.2.1.1.1(1)
SNP	<i>eis</i> C-14T		181	SNP	<i>eis</i> G3A	V1V	6	6	1.1.1.1.1(1) & 2.2.1.1.1.i3(5)
SNP	<i>ahpC</i> G-88A		3350	INDEL	<i>ahpC</i> 323delC	108	1	1	3.1.1(1)
SNP	<i>whiB7</i> T-147C		1	INDEL	<i>whiB7</i> 192delC	64	573	1	1.2.1.1.1(1)
INDEL	<i>whiB7</i> -214delG		9	INDEL	<i>whiB7</i> 192delC	64	573	1	1.2.1.1.1(1)
INDEL	<i>whiB7</i> -316insC		2	INDEL	<i>whiB7</i> 192delC	64	573	1	1.2.1.1.1(1)

142

143 **Table 2. Co-occurrence of regulator resistance mutations and regulon LoF mutations.** A list
 144 of antibiotic resistance mutations in **regulator** regions (*mmpR*, *oxyR-ahpC*, *eis-Rv2417c*, *whiB7-*
 145 *uvrD2*) that co-occur with LoF mutations in corresponding **regulated** regions (*mmpL5*, *mmpS5*,
 146 *mmpR*, *ahpC*, *eis*, *whiB7*) within our sample of 31,440 clinical isolates. A more detailed table can
 147 be found in **Supplementary Table 3**.

148

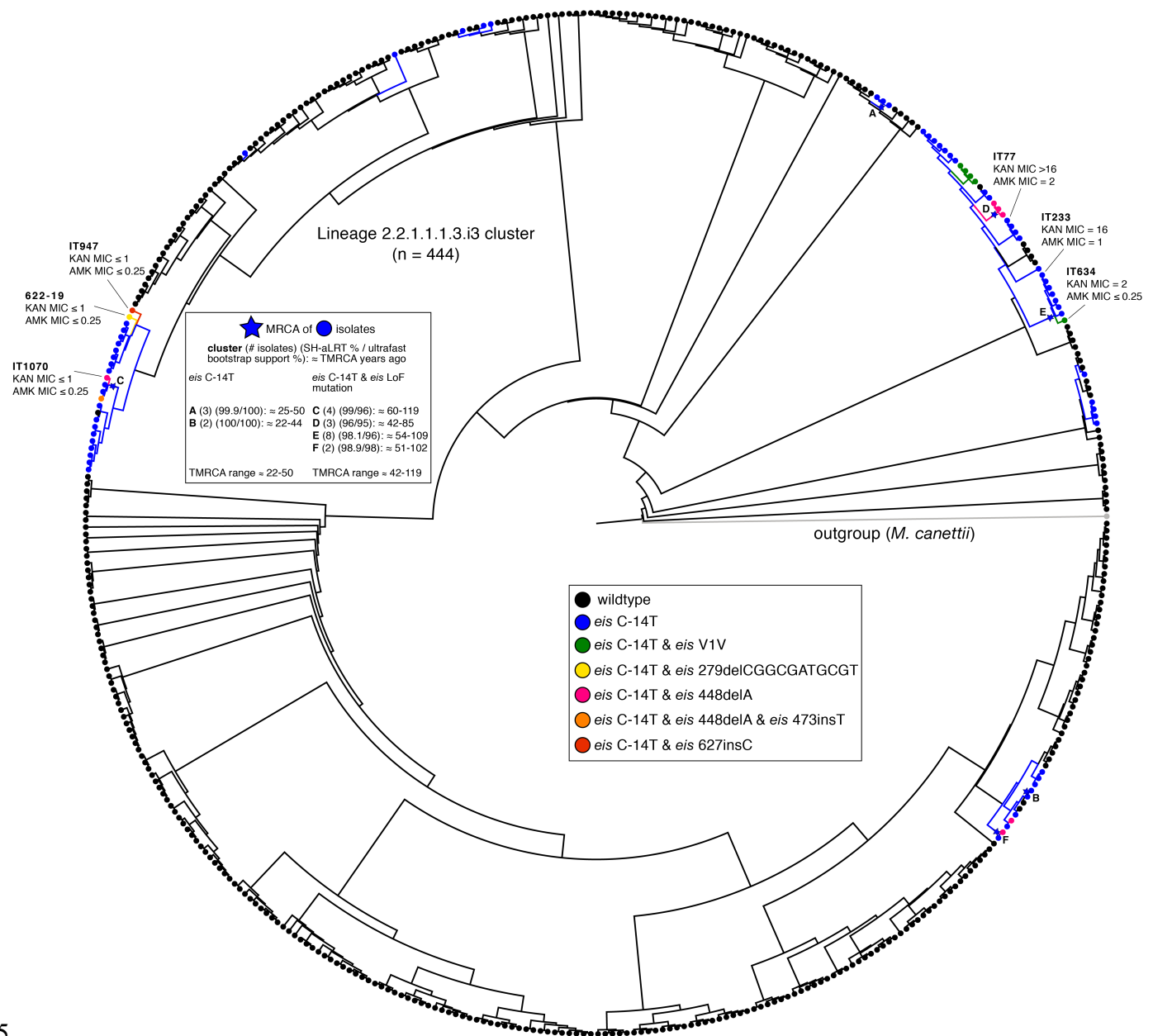
149 We constructed a phylogeny of all 188 MTBC sub-lineage 4.11 isolates and to study how
 150 the LoF mutations in *mmpR* and *mmpS5-mmpL5* evolved (**Fig. 2**). The majority of isolates with
 151 the *mmpR* or *mmpL5* frameshifts harbored both (82/84) but, based on the topology of the tree, we
 152 were unable to determine which of the two frameshifts arose first. Consequently, we could only date
 153 the common most recent common ancestor (MRCA). We approximated the age of the MRCA at
 154 66-132 years prior to sampling *i.e.* well before the use of BDQ or CFZ in treatment regimens for
 155 TB in the pre-antibiotic era and likely before the introduction of thioacetazone, which is also
 156 exported by *mmpS5-mmpL5* (Halloum et al., 2017; Ma et al., 2010).

157

158 **KAN: *eis* upstream mutations in combination with *eis* LoF mutations**

159 We observed 23 unique LoF mutations upstream of *eis* (**Fig. 1C**), of which ten were homoplasic
160 and occurred in at least five isolates (**Supplementary Table 1**). As expected, the classical G-37T,
161 C-14T, C-12T, and G-10A mutations, which are known to confer KAN resistance based on allelic
162 exchange and/or complementation experiments, were most frequent (**Table 1**) (Pholwat et al.,
163 2016; World Health Organization, 2018; Zaunbrecher et al., 2009). Specifically, 881 isolates with
164 either *eis* G-37T, C-12T, or G-10A and 179 isolates with *eis* C-14T did not have any of the other
165 key AG resistance mutations in *rrs* (i.e. A1401G, C1402T, or G1484T, see **Supplementary Table**
166 **5**) (World Health Organization, In press).

167 We identified 30 unique LoF mutations in *eis*, of which five were homoplasic and occurred
168 in at least five isolates (**Table 1, Supplementary Table 1**). These LoF never coincided with *eis*
169 G-37T, C-12T, or G-10A, whereas this was the case for 21 *eis* C-14T mutants (i.e. 13 isolates with
170 indels, six with a G3A synonymous change that abolished the valine start codon, and two with
171 nonsense mutations) (**Table 2, Supplementary Table 3**). MIC data were available for five of these
172 *eis* double mutants, which confirmed that they were susceptible to KAN whereas seven *eis* C-14T
173 control isolates with a wild-type *eis* coding region were KAN resistant (**Fig. 3, Table 3**). The
174 corresponding AMK MIC data mirrored the results for KAN.



175

176 **Fig. 3. Phylogeny of 444 sub-lineage 2.2.1.1.1.3.i3 isolates.** Isolates with the *eis* C-14T promoter
 177 SNP and no LoF variants in *eis* are colored in blue (n=61), whereas isolates carrying both the *eis*
 178 C-14T promoter SNP and a LoF mutation in *eis* are colored according to the legend (n=14).
 179 TMRCA estimates for groups of isolates with the *eis* C-14T promoter SNP are given in the upper
 180 left. The MICs (mg/L) for isolates from this sub-lineage are also included (MICs for isolates from
 181 other sub-lineages can be found in **Table 3**).

182

isolate ID	2.2.1.1.1.3.i3 cluster	<i>rrs</i> A1401G	<i>rrs</i> C1402T	<i>rrs</i> G1484T	<i>eis</i> C-14T	<i>eis</i> LOF mutation	KAN	AMK
IT123	no	A	C	G	yes	no	16	0.5
IT184	no	A	C	G	yes	no	8	0.5
IT952	no	A	C	G	yes	no	16	2
IT524	no	A	C	G	yes	no	16	2
655-19	no	A	C	G	yes	no	16	1
IT233	yes	A	C	G	yes	no	16	1
IT77	yes	A	C	G	yes	no	>16	2
622-19	yes	A	C	G	yes	yes (279delCGGCGATGCGT)	≤1	≤0.25
IT1070	yes	A	C	G	yes	yes (448delA)	≤1	≤0.25
IT947	yes	A	C	G	yes	yes (627insC)	≤1	≤0.25
168-19	no	A	C	G	yes	yes (400insG)	≤1	≤0.25
IT634	yes	A	C	G	yes	yes (V1V)	2	≤0.25
SAMN02419559	yes	A	C	G	yes	yes (V1V)		
SAMN02419535	yes	A	C	G	yes	yes (V1V)		
SAMN02419543	yes	A	C	G	yes	yes (V1V)		
SAMN02419586	yes	A	C	G	yes	yes (V1V)		
SAMN07236283	no	A	C	G	yes	yes (V1V)		
SAMEA1016073	yes	A	C	G	yes	yes (448delA)		
SAMEA1403685	yes	A	C	G	yes	yes (448delA)		
SAMEA1403638	yes	A	C	G	yes	yes (448delA)		
SAMN02584676	yes	A	C	G	yes	yes (448delA)		
SAMN04633319	yes	A	C	G	yes	yes (448delA)		
SAMN08376196	no	A	C	G	yes	yes (W13*)		
SAMN08709032	no	A	C	G	yes	yes (W13*)		
SAMN06210015	no	A	C	G	yes	yes (16insC)		
Peru2946	no	A	C	G	yes	yes (486insCT)		
Peru3354	no	A	C	G	yes	yes (486insCT)		
SAMN02584612	no	A	C	G	yes	yes (15insC)		
SAMN07956543	yes	A	C	G	yes	yes (448delA & 473insT)		

183

184 **Table 3. Overview of *eis* C-14T isolates with *eis* LoF mutations and corresponding MICs for**
185 **KAN and AMK.** All 29 double mutants lacked the three classical AG resistance mutations in *rrs*.
186 Isolates that are part of the 2.2.1.1.1.3.i3 cluster are shown in Fig 3. MICs (mg/L) were measured
187 using either the UKMYC5 or UKMYC6 broth microdilution plates by Thermo Fisher (Fowler and
188 CRyPTIC Consortium, 2021). The provisional CCs for KAN and AMK for these plates are 4 and
189 1 mg/L, respectively. Unlike for KAN, *eis* C-14T only has a modest effect on the MIC of AMK
190 (i.e. the MIC distribution of this mutation spans the CC when the efflux pump is active, which is
191 in line with data from other media) (Gygli et al., 2019; Pholwat et al., 2016; World Health
192 Organization, 2018; Zaunbrecher et al., 2009). Consequently, KAN is a more sensitive agent to
193 detect an inactive efflux pump than AMK. More details can be found in **Supplementary Table 6.**

194

195 The most common MTBC sub-lineage with the *eis* double mutants was 2.2.1.1.i3 (14/22
196 isolates). Of the 31,440 isolates, 444 were belonged to this sub-lineage and clustered closely based
197 on their pairwise SNP distance. The phylogeny of these 444 isolates showed that *eis* C-14T arose
198 more than nine times independently (**Fig. 3**). We approximated the MRCA for the six group of
199 isolates that had high bootstrap support. For the two groups of isolates that only harbored the *eis*
200 C-14T mutation the MRCA was dated to 22-50 years ago, in line with KAN introduction into
201 clinical use in approximately 1958 (Ektefaie et al., 2021). Of the 75 isolates with the *eis* C-14T
202 promoter mutation, 14 also harbored an *eis* LoF mutation that arose nine times independently (**Fig.**
203 **3, Table 2**). In each instance, the LoF variant emerged within a clade of *eis* C-14T mutants
204 suggesting that it appeared later in time. We compared the MRCA of the clades with double
205 mutants to those with *eis* C-14T only. We found the MRCA of the double mutants to be older on
206 average, suggesting that time and possibly fluctuating evolutionary pressures are needed for LoF
207 to develop in an *eis* C-14T background.

208

209 **KAN: *whiB7* upstream mutations in combination with *whiB7* or *eis* LoF mutations**

210 We found 116 unique SNPs upstream of *whiB7* (**Fig. 1C**), of which 8 were homoplasic and
211 occurred in at least five isolates (**Table 1, Supplementary Table 1**). We identified 10 unique LoF
212 mutations in *whiB7*, two of which (nt193insG & nt133delCA) evolved repeatedly, across 657
213 isolates. The most frequent mutation (nt192delC) occurred in 573 sub-lineage 1.2.1.1 isolates,
214 which was in agreement with earlier findings (Merker et al., 2020). This was the only LoF mutation
215 in *whiB7* to coincide with an upstream mutation (i.e. in three isolates in total, each with a different
216 upstream mutation) (**Table 2, Supplementary Table 3**). Because none of these upstream
217 mutations had been described in the literature, it was unclear whether these represented potential

218 examples of epistasis (Heyckendorf et al., 2018; Reeves et al., 2013). Finally, no LoF mutations
219 in *eis* were found in isolates harboring mutations upstream of *whiB7*.

220 DISCUSSION

221 Although our analysis did not yield any strong evidence for epistasis involving *ahpC* or *whiB7*,
222 our finding that epistasis is possible due to LoF mutations in *eis* is not only relevant for AGs but
223 has wider implications for the interpretation of sequencing data. First, with the exception of two
224 synonymous mutations in *aftA* (*Rv3792*) and *fabG1* (*Rv1483*) that confer ethambutol and
225 ethionamide/INH resistance respectively by creating alternative promoters, synonymous
226 mutations are typically excluded *a priori* from the analysis of WGS data (Ando et al., 2014; Safi
227 et al., 2013). We demonstrated that this assumption is not sound for start codons given that only
228 one of the four triplets encoding valine can act as a start codon (i.e. GTC). Second, evidence for
229 epistasis argues strongly that multivariate prediction approaches are needed for accurate resistance
230 prediction from sequencing data.

231 It is notable that *eis* LoF mutations coincided only with *eis* C-14T mutants, even though
232 isolates with *eis* G-37T, C-12T, and G-10A without any AG resistance mutations in *rrs* were
233 almost five times more frequent in our dataset (**Supplementary Table 5**). We hypothesize that
234 because *eis* C-14T leads to a greater up-regulation of *eis* than the other three mutations, this comes
235 at a fitness cost, unless selective antibiotic pressure is maintained (Pholwat et al., 2016; Sanz-
236 García et al., 2019; World Health Organization, 2018; Zaunbrecher et al., 2009). Indeed, molecular
237 dating and the topology of the in tree (**Fig. 3**) suggested that the LoF mutations arose independently
238 on multiple occasions after the acquisition of *eis* C-14T. To our knowledge, this represents the
239 strongest evidence to date for genotypic reversion from resistance to a susceptible phenotype for
240 MTBC (Richardson et al., 2009). We would like to stress, however, that even for AMK this is a
241 rare phenomenon, given that of the 179 isolates that harbored *eis* C-14T without any AG resistance
242 mutations in *rrs* only 12% had concomitant *eis* LoF mutations (**Table 2, Supplementary Table**

243 5). In other words, the cautious approach would be to still interpret *eis* C-14T as a marker for AMK
244 resistance to construct a relevant regimen, unless there is strong evidence that a particular isolate
245 is affected by epistasis (e.g. unlike the SL-LPA, the tNGS assays by ABL and Deeplex actually
246 interrogate part of the *eis* coding region) (Mohamed et al., 2021; World Health Organization,
247 2018). This initial treatment decision may then have to be adjusted based on the phenotypic DST
248 result.

249 Because we did not have BDQ or CFZ MICs for any of the *mmpR/mmpL5* double-LoF
250 mutants, it remains to be determined whether these are examples of epistasis (in the case of the
251 Peruvian cluster, this would be unrelated to antibiotic pressure, unlike for *eis*). We note, however,
252 that indirect evidence exists that is in line with this prediction. Villellas et al. reported the BDQ
253 MICs and *mmpR* sequence results for baseline isolates from the C208 phase 2b trial of BDQ, which
254 featured five isolates with the same *mmpR* frameshift that we observed in the Peruvian cluster (**Fig.**
255 **2**) (Diacon et al., 2014; Villellas et al., 2017). Three of the trial mutants were from South Africa
256 and had 7H11 MICs of 0.25-1 mg/L (i.e. \geq CC of 0.25 mg/L and, thus, consistent with a functional
257 efflux pump and resistant phenotype if an area of technical uncertainty is set at 0.25 mg/L, as
258 previously proposed) (Beckert et al., 2020; Nimmo et al., 2020; Villellas et al., 2017). By contrast,
259 even the lowest BDQ concentration tested (i.e. 0.008 mg/L) inhibited the growth of the remaining
260 two trial mutants that were isolated in 2009 in Lima, Peru (N. Lounis, personal communication).
261 Given that the Peruvian double-LoF cluster from this study was isolated in the same city during
262 the same period, it is possible that the latter two trial isolates were from this cluster, although this
263 remains to be confirmed using WGS data and, ideally, repeat MIC testing to exclude experimental
264 error.

265 The possibility of epistasis underlines the need for comprehensive microbiological workup
266 of the ongoing clinical trials of BDQ. *mmpR* as well as *mmpS5-mmpL5* and the corresponding
267 inter-genic region have to be analyzed along with standardized MIC testing using an on-scale
268 quality control strain for both BDQ and CFZ (Kaniga et al., 2016; Schön et al., 2020, 2019). We
269 recommend that discordances between genotypic and phenotypic DST results are confirmed by
270 retesting and, where warranted, followed up with specialized testing. For example, the two
271 aforementioned Peruvian results may actually be hyper-susceptible to BDQ and CFZ (i.e. lower
272 concentrations would have to be tested to determine the MIC endpoint) (Merker et al., 2020). If
273 confirmed, this should also apply to *mmpS5-mmpL5* LoF mutants with wild-type *mmpR* (e.g. just
274 over half of the lineage 1.1.1.1 isolates in our dataset had a nonsense mutation in *mmpL5*) and may
275 have implications for the ongoing trials of TBAJ-876, TBAJ-587, TBI-166, and OPC-167832 as
276 these agents are also exported by this pump (Hariguchi et al., 2020; Xu et al., 2021, 2019).
277

278 **MATERIALS AND METHODS**

279 **Sequence Data**

280 We initially downloaded raw Illumina sequence data for 33,873 clinical isolates from NCBI
281 (Benson et al., 2008). We identified the BioSample for each isolate and downloaded all of the
282 associated Illumina sequencing runs. Isolates had to meet the following quality control measures
283 for inclusion in our study: (i) at least 90% of the reads had to be taxonomically classified as
284 belonging to MTBC after running the trimmed FASTQ files through Kraken (Wood and Salzberg,
285 2014) and (ii) at least 95% of bases had to have coverage of at least 10x after mapping the
286 processed reads to the H37Rv reference genome (Genbank accession: NC_000962).

287

288 **Illumina Sequencing FastQ Processing and Mapping to H37Rv**

289 The raw sequence reads from all sequenced isolates were trimmed with version 0.20.4 Prinseq
290 (settings: -min_qual_mean 20) (Schmieder and Edwards, 2011) and then aligned to H37Rv with
291 version 0.7.15 of the BWA mem algorithm using the -M settings (Li and Durbin, 2009). The
292 resulting SAM files were then sorted (settings: SORT_ORDER = coordinate), converted to BAM
293 format, and processed for duplicate removal with version 2.8.0 of Picard
294 (<http://broadinstitute.github.io/picard/>) (settings: REMOVE_DUPLICATES = true,
295 ASSUME_SORT_ORDER = coordinate). The processed BAM files were then indexed with
296 Samtools (Li et al., 2009). We used Pilon (settings: --variant) on the resulting BAM files to
297 generate VCF files that contained calls for all reference positions corresponding to H37Rv from
298 pileup (Walker et al., 2014).

299

300 **Empirical Score for Difficult-to-Call Regions**

301 We assessed the congruence in variant calls between short-read Illumina data and long-read
302 PacBio data for a set of isolates that underwent sequencing with both technologies. Using 31
303 isolates for which both Illumina and a complete PacBio assembly were available, we evaluated the
304 empirical base-pair recall (EBR) of all base-pair positions of the H37rv reference genome. For
305 each sample, the alignments of each high confidence genome assembly to the H37Rv genome were
306 used to infer the true nucleotide identity of each base pair position. To calculate the empirical base-
307 pair recall, we calculated what percentage of the time our Illumina based variant calling pipeline,
308 across 31 samples, confidently called the true nucleotide identity at a given genomic position. If
309 Pilon variant calls did not produce a confident base call (*Pass*) for the position, it did not count as
310 a correct base call. This yields a metric ranging from 0.0–1.0 for the consistency by which each
311 base-pair is both confidently and correctly sequenced by our Illumina WGS based variant calling
312 pipeline for each position on the H37Rv reference genome. An H37Rv position with an EBR score
313 of x% indicates that the base calls made from Illumina sequencing and mapping to H37Rv agreed
314 with the base calls made from the PacBio *de novo* assemblies in x% of the Illumina-PacBio pairs.
315 We masked difficult-to-call regions by dropping H37Rv positions with an EBR score below 0.9
316 as part of our variant calling procedure. Full details on the data and methodology can be found
317 elsewhere (Vargas et al., 2021).

318 **Variant Calling**

320 SNP Calling: To prune out low-quality base calls that may have arisen due to sequencing or
321 mapping error, we dropped any base calls that did not meet any of the following criteria: (i) the
322 call was flagged as *Pass* by Pilon, (ii) the mean base quality at the locus was >20, (iii) the mean
323 mapping quality at the locus was >30, (iv) none of the reads aligning to the locus supported an
324 insertion/deletion (indel), (v) a minimum coverage of 20 reads at the position, and (vi) at least 75%

325 of the reads aligning to that position supported 1 allele (using the *INFO.QP* field which gives the
326 proportion of reads supporting each base weighted by the base and mapping quality of the reads,
327 *BQ* and *MQ* respectively, at the specific position). A base call that did not meet all filters (i) – (vi)
328 was inferred to be low-quality/missing.

329 Indel Calling: To prune out low-quality indel variant calls, we dropped any indel that did not meet
330 any of the following criteria: (i) the call was flagged as *Pass* by Pilon, (ii) the maximum length of
331 the variant was 10bp, (iii) the mean mapping quality at the locus was >30, (iv) a minimum coverage
332 of 20 reads at the position, and (v) at least 75% of the reads aligning to that position supported the
333 indel allele (determined by calculating the proportion of total reads *TD* aligning to that position
334 that supported the insertion or deletion, *IC* and *DC* respectively). A variant call that met filters (i),
335 (iii), and (iv) but not (ii) or (v) was inferred as a high-quality call that did not support the indel
336 allele. Any variant call that did not meet all filters (i), (iii), and (iv) was inferred as low-
337 quality/missing.

338 Intermediate Allele Frequency Indel Calling: To call indel variants in which the indel allele was
339 detected at an intermediate frequency, we made the following modification to the *indel Calling*
340 filters outlined above. Filter (v) above is replaced with the following two filters: (v-i) at least 10%
341 but less than 75% of the reads aligning to that position supported the indel allele and (v-ii) at least
342 10 reads support the indel allele. The *mmpR* analysis was restricted to isolates with 100x coverage
343 across $\geq 99\%$ of the gene.

344

345 **SNP Genotypes Matrix**

346 We detected SNP sites at 899,035 H37Rv reference positions (of which 64,950 SNPs were not
347 biallelic) among our global sample of 33,873 isolates. We constructed a 899,035x33,873

348 genotypes matrix (coded as 0:A, 1:C, 2:G, 3:T, 9:Missing) and filled in the matrix for the allele
349 supported at each SNP site (row) for each isolate, according to the *SNP Calling* filters outlined
350 above. If a base call at a specific reference position for an isolate did not meet the filter criteria
351 that allele was coded as *Missing*. We excluded 20,360 SNP sites that had an EBR score <0.90,
352 another 9,137 SNP sites located within mobile genetic element regions (e.g. transposases,
353 intergrases, phages, or insertion sequences) (Comas et al., 2010; Vargas et al., 2021), then 31,215
354 SNP sites with missing calls in >10% of isolates, and 2,344 SNP sites located in overlapping genes
355 (coding sequences). These filtering steps yielded a genotypes matrix with dimensions
356 835,979x33,873. Next, we excluded 1,663 isolates with missing calls in >10% of SNP sites
357 yielding a genotypes matrix with dimensions 835,979x32,210. We used an expanded 96-SNP
358 barcode to type the global lineage of each isolate in our sample (Freschi et al., 2020). We further
359 excluded 325 isolates that either did not get assigned a global lineage, assigned to more than one
360 global lineage, or were typed as lineage 7. We then excluded 41,760 SNP sites from the filtered
361 genotypes matrix in which the minor allele count = 0 which resulted in a 794,219x31,885 matrix.
362 To provide further MTBC lineage resolution on the lineage 4 isolates, we required an MTBC sub-
363 lineage call for each lineage 4 isolate. We excluded 457 isolates typed as global lineage 4 but had
364 no further sub-lineage calls and then again excluded 11,654 SNP sites from the filtered genotypes
365 matrix in which the minor allele count=0. The genotypes matrix used for downstream analysis had
366 dimensions 782,565x31,428, representing 782,565 SNP sites across 31,428 isolates. The global
367 lineage (L) breakdown of the 31,428 isolates was: L1=2,815, L2=8,090, L3=3,398, L4=16,931,
368 L5=98, L6=96.

369
370 **Indel Genotypes Matrix**

371 We detected 53,167 unique indel variants within 50,576 H37Rv reference positions among our
372 global sample of 33,873 isolates. We constructed a 53,167x33,873 genotypes matrix (coded as
373 1:high quality call for the indel allele, 0:high quality call not for the indel allele, 9:Missing) and
374 filled in the matrix according to whether the indel allele was supported for each indel variant (row)
375 for each isolate, according to the *Indel Calling* filters outlined above. If a variant call at the
376 reference position for an indel variant did not meet the filter criteria that call was coded as *Missing*.
377 We excluded 2,006 indel variants that had an EBR score <0.90, another 694 indel variants located
378 within mobile genetic element regions, then 207 indel variants located in overlapping genes
379 (coding sequences). These filtering steps yielded a genotypes matrix with dimensions
380 50,260x33,873. Next, we excluded any isolate that was dropped while constructing the SNP
381 genotypes matrix to retain the same 31,428 isolates as described above. The genotypes matrix used
382 for downstream analysis had dimensions 50,260x31,428.

383

384 **Mixed Allele Frequency Indel Genotypes Matrix**

385 After following the same filtering steps outlined above under *Indel Genotypes Matrix*, we detected
386 7,731 unique indel variants in our filtered sample of 31,428 isolates in which at least one isolate
387 supported each indel variant at an intermediate allele frequency ($10\% \leq AR < 75\%$). We constructed
388 a 7,731x31,428 genotypes matrix (coded as 0:high quality call not for the indel allele, -9:Missing,
389 or 10-74:the % of reads supporting the indel allele) and filled the matrix according to whether the
390 indel allele was supported at an intermediate allele frequency for each indel variant (row) for each
391 isolate, according to the *Intermediate Allele Frequency Indel Calling* filters outlined above. To
392 determine the limit of detection for indels that might be present at lower allele frequencies, we
393 calculated the number of isolates in our sample that have 100x coverage in $\geq 99\%$ of the locus for

394 *mmpR* (7,435), *mmpS5* (8,949), and *mmpL5* (6,217) (**Supplementary Table 2**). We retained only
395 frameshift indels yielding a genotypes matrix with dimensions 5,925x31,428 and interrogated only
396 the *mmpR* - *mmpS5* - *mmpL5* chromosomal region for the presence of mixed indels
397 (**Supplementary Table 2**).

398

399 **Inclusion and Processing of 12 *eis* C-14T mutants with AG MICs**

400 We added 12 clinical *eis* C-14T mutants to the dataset, for which we had KAN and AMK MICs
401 and some of which had a LoF mutation in *eis*. We processed the raw sequencing reads according
402 to the methods described above to generate VCF files. We genotyped SNPs for these isolates at
403 the 782,565 SNP sites and genotyped indels for the 50,260 indel variants previously identified
404 using the same filters described above to construct 782,565x12 and 50,260x12 matrices,
405 respectively.

406 During analysis, we observed that 3/12 isolates (IT947, 622-19 and 168-19) carried the *eis* C-14T
407 promoter resistance mutation and no observed LoF mutation in *eis* but were phenotypically
408 susceptible according to KAN MICs. Upon further inspection of the VCF files for these isolates,
409 we found that all three isolates had a LoF mutation in *eis* that we originally did not detect per
410 our variant calling methodology. We found that one isolate (622-19) had an 11bp deletion in *eis*
411 which was not represented in the 50,260 indel variants since we restricted our analysis to indels
412 ≤ 10 bp and consequently was excluded from our 50,260x12 matrix. Each of the other two strains,
413 IT947 and 168-19, had a different 1bp insertion in *eis* that was not identified in our original pool
414 of 31,428 isolates, so it also was also not represented in the 50,260x12 matrix. We updated our
415 variant call data by incorporating these newly identified variants (**Table 2, Supplementary 1,**
416 **Supplementary Table 3**).

417

418 **Targeted Chromosomal Regions**

419 We queried our SNP and indel matrices for the following types of mutations in the following
420 regions of the H37Rv Reference Genome: [1] *mmpR - mmpS5 - mmpL5*: the coding sequences
421 for *mmpR* (778990 - 779487), *mmpS5* (778477 - 778905), and *mmpL5* (775586 - 778480) for
422 nonsense SNVs (single nucleotide variant), frameshift indels, missense SNVs that abolish the start
423 codon, and synonymous SNVs that abolish the start codon for *mmpR* which starts with a valine
424 (we did not check for synonymous SNVs at the first codon for *mmpS5* or *mmpL5* because these
425 coding sequences start with a methionine). [2] upstream *ahpC - ahpC*: the intergenic region *oxyR-*
426 *ahpC* (2726088 - 2726192) for SNVs and indels, and the coding sequence for *ahpC* (2726193 -
427 2726780) for nonsense SNVs, frameshift indels, and missense SNVs that abolish the start codon.
428 We did not check for synonymous SNVs at the first codon for *ahpC* because the coding sequence
429 starts with a methionine (and also serves as the initiation site). [3] upstream *eis - eis*: the intergenic
430 region *eis-Rv2417c* (2715333 - 2715383) for SNVs and indels, and the coding sequence
431 for *eis* (2714124 - 2715332) for nonsense SNVs, frameshift indels, missense SNVs that abolish
432 the start codon, and synonymous SNVs that abolish the START codon. [4] upstream *whiB7 -*
433 *whiB7*: the intergenic region *whiB7-uvrD2* (3568680 - 3569082) for SNVs and indels, and the
434 coding sequence for *whiB7* (3568401 - 3568679) for nonsense SNVs, missense SNVs that abolish
435 the start codon, frameshift indels, and synonymous SNVs that abolish the start codon.

436

437 **Antibiotic Resistance Mutations in *rrs* and *atpE***

438 Resistance to aminoglycosides can occur as a result of mutations in the 1,400bp region of the 16S
439 rRNA (*rrs*), where *rrs* A1401G, C1402T, and G1484T mutations have all been implicated in

440 aminoglycoside resistance (Kambli et al., 2016; Reeves et al., 2013). To ensure that isolates were
441 not aminoglycoside resistant directly from harboring one of these *rrs* mutations, we genotyped
442 (with $\geq 75\%$ allele frequency) the 1401, 1402, and 1484 nucleotide coordinates in *rrs* for the set of
443 12 added isolates with *eis* C-14T promoter resistance mutations and 17 other isolates (from our
444 original set of 31428 isolates) with coinciding *eis* C-14T promoter resistance mutation and *eis* LoF
445 mutations (**Fig. 3, Table 3, Supplementary Table 3**). None of these 29 isolates harbored any of
446 the *rrs* A1401G, C1402T, or G1484T aminoglycoside resistance mutations (**Table 3**). Similarly,
447 single nucleotide variants in the gene *atpE*, which encodes the BDQ target, have been associated
448 with high-level BDQ resistance (Kadura et al., 2020). We interrogated the genotypes for 29 SNP
449 sites in *atpE* (SNPs that were present within our pool of 31,428 isolates) in the 84 isolates that
450 harbored both a frameshift in *mmpR* and frameshift in *mmpL5* (**Fig. 1, Table 2**) and found that
451 none of the isolates carried a mutant allele at any of these SNP sites.

452

453 **Phylogeny Construction and assessment of convergent evolution**

454 To generate the trees, we first merged the VCF files of the isolates in the sample (188 lineage 4.11
455 isolates & 444 lineage 2.2.1.1.1.3.i3 isolates) with bcftools (Li et al., 2009). We then removed
456 repetitive, antibiotic resistance and low coverage regions (Freschi et al., 2020). We generated a
457 multi-sequence FASTA alignment from the merged VCF file with vcf2phylip (version 1.5,
458 <https://doi.org/10.5281/zenodo.1257057>). We constructed the phylogenetic tree with IQ-TREE
459 (Nguyen et al., 2015). We used the *mset* option to restrict model selection to GTR models,
460 implemented the automatic model selection with ModelFinder Plus (Kalyaanamoorthy et al., 2017)
461 and computed the SH-aLRT test and bootstrap values with UFBoot (Minh et al., 2013) with 1000
462 bootstrap replicates.

463 To quantify the number of independent mutational events (SNPs & indels) in the original sample
464 of 31,428, we grouped isolates into eight groups based off of genetic similarity, five groups
465 corresponding to global lineages 1, 2, 3, 5, 6 and three groups for global lineage 4. We constructed
466 eight phylogenies from these groups, then used the genotypes in conjunction with the phylogenies
467 to assess the number of independent arisals for each mutation observed. We used an ancestral
468 reconstruction approach to quantify the number of times each SNV arose independently in the
469 phylogenies using SNPPar (Edwards et al., 2020). This yielded a *homoplasy score* or an estimate
470 for the number of independent arisals for each SNV (**Supplementary Table 1**). To quantify the
471 number of independent arisals for each indel, we developed a simple method to count the number
472 of times each indel allele “breaks” the phylogenies. If a given mutant allele is observed in two
473 separate parts of a phylogeny, then we can assume that this allele arose twice in pool of isolates
474 used to construct the tree. We calculated a *homoplasy score* by counting these topology disruptions
475 for both SNVs & indels. The results for the SNVs were congruent with the *homoplasy scores*
476 computed from the ancestral reconstructions, validating this approach for computing *homoplasy*
477 *scores* for indels.

478

479 **MRCA Dating Approximation**

480 To date the arisal of a specific mutation within a group of isolates on a phylogeny, we looked for
481 groups of isolates on the trees that carried the mutant allele of interest. We grouped isolates
482 according to the following principles: (1) a group of isolates had to be a sub-tree of 2 or more
483 monophyletic mutants, and (2) we identified the MRCA of all mutants in that sub-tree assuming
484 that reversion of mutations is impossible. For a given group, we checked that the MRCA of the
485 isolates had an SH-aLRT of $\geq 80\%$ and an ultrafast bootstrap support of $\geq 95\%$. If these conditions

486 were satisfied, indicating high confidence in the branch, we then calculated the median branch
487 length (SNPs/site) between the MRCA and the tips. We multiplied the median branch length
488 (SNPs/site) by the number of sites in the SNP concatenate used to construct the tree to get the
489 median branch length in SNPs/genome. Molecular clock estimates for MTBC range from 0.3-0.6
490 SNPs/genome/year, we divided the branch lengths in SNPs/genome by 0.3 SNPs/genome/year and
491 0.6 SNPs/genome/year to get upper and lower bound estimates for the MRCA age.

492

493 **Data Analysis and Variant Annotation**

494 Data analysis was performed using custom scripts run in Python and interfaced with iPython (Pérez
495 and Granger, 2007). Statistical tests were run with Statsmodels (Seabold and Perktold, 2010) and
496 Figures were plotted using Matplotlib (Hunter, 2007). Numpy (Van Der Walt et al., 2011),
497 Biopython (Cock et al., 2009) and Pandas (McKinney and others, 2010) were all used extensively
498 in data cleaning and manipulation. Functional annotation of SNPs was done in Biopython using
499 the H37Rv reference genome and the corresponding genome annotation. For every SNP variant
500 called, we used the H37Rv reference position provided by the Pilon (Walker et al., 2014) generated
501 VCF file to determine the nucleotide and codon positions if the SNP was located within a coding
502 sequence in H37Rv. We extracted any overlapping CDS region and annotated SNPs accordingly,
503 each overlapping CDS regions was then translated into its corresponding peptide sequence with
504 both the reference and alternate allele. SNPs in which the peptide sequences did not differ between
505 alleles were labeled synonymous, SNPs in which the peptide sequences did differ were labeled
506 non-synonymous and if there were no overlapping CDS regions for that reference position, then
507 the SNP was labeled intergenic. Functional annotation of indels was also done in Biopython using
508 the H37Rv reference genome and the corresponding genome annotation. For every indel variant

509 called, we used the H37Rv reference position provided by the Pilon generated VCF file to
510 determine the nucleotide and codon positions if the indel was located within a coding sequence in
511 H37Rv. An indel variant was classified as in-frame if the length of the indel allele was divisible
512 by three, otherwise it was classified as a frameshift.

513 **SUPPLEMENTARY TABLE DESCRIPTIONS**

514

515 **Supplementary Table 1. Mutations detected in a global sample of MTBC clinical isolates.** A
516 full list of mutations that occur within our sample of 31,440 clinical isolates within the *mmpL5*,
517 *mmpS5*, *mmpR*, *ahpC*, *eis*, *whiB7* coding sequences and *oxyR-ahpC*, *eis-Rv2417c*, *whiB7-uvrD2*
518 intergenic regions.

519

520 **Supplementary Table 2. Mixed indels in the *mmpR-mmpL5-mmpS5* chromosomal region.** A
521 list of frameshift indels that were detected at an intermediate allele frequencies between 10% and
522 75% in *mmpR*, *mmpS5*, or *mmpL5* within our sample of 31,428 isolates (excludes the set of 12
523 added isolates, see **Methods**).

524

525 **Supplementary Table 3. Co-occurrence of regulator resistance mutations and regulon LoF**
526 **mutations.** A more detailed version of **Table 2**.

527

528 **Supplementary Table 4. Binary resistance phenotypes for MTBC sub-lineage 4.11 isolates.**
529 A table of binary resistance phenotype (STR, INH, RIF, EMB, PZA, AMK & KAN) data for a
530 subset isolates that belong to sub-lineage 4.11 (**Fig. 2**), curated from multiple studies (Groschel et
531 al., 2021).

532

533 **Supplementary Table 5. Count of isolates with *eis* promoter mutations and no coinciding *rrs***
534 **AG resistance mutations.** The count of isolates with *eis* promoter mutations (G-10A, C-12T, C-
535 14T, G-37T) that coincide with any AG resistance mutations in *rrs* (A1401G, C1402T, G1484T).

536

537 **Supplementary Table 6. KAN and AMK resistance details for strains with MICs and strains**
538 **with double *eis* promoter SNP & *eis* LoF mutations.** A more detailed version of **Table 3** with
539 binary resistance phenotype (STR, INH, RIF, EMB, PZA, AMK & KAN) data for a subset of
540 isolates (Groschel et al., 2021).

541

542 REFERENCES

- 543 Ajileye A, Alvarez N, Merker M, Walker TM, Akter S, Brown K, Moradigaravand D, Schön T,
544 Andres S, Schleusener V. 2017. Some synonymous and nonsynonymous gyrA mutations
545 in Mycobacterium tuberculosis lead to systematic false-positive fluoroquinolone
546 resistance results with the Hain GenoType MTBDRsl assays. *Antimicrob Agents*
547 *Chemother* **61**:e02169-16.
- 548 Ando H, Miyoshi-Akiyama T, Watanabe S, Kirikae T. 2014. A silent mutation in mabA confers
549 isoniazid resistance on Mycobacterium tuberculosis. *Mol Microbiol* **91**:538–547.
- 550 Andres S, Merker M, Heyckendorf J, Kalsdorf B, Rumetshofer R, Indra A, Hofmann-Thiel S,
551 Hoffmann H, Lange C, Niemann S. 2020. Bedaquiline-resistant Tuberculosis: Dark
552 Clouds on the Horizon. *Am J Respir Crit Care Med* **201**:1564–1568.
- 553 Beckert P, Sanchez-Padilla E, Merker M, Dreyer V, Kohl TA, Utpatel C, Köser CU, Barilar I,
554 Ismail N, Omar SV. 2020. MDR M. tuberculosis outbreak clone in Eswatini missed by
555 Xpert has elevated bedaquiline resistance dated to the pre-treatment era. *Genome Med*
556 **12**:1–11.
- 557 Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. 2008. GenBank. *Nucleic Acids*
558 *Res* **37**:D26–D31.
- 559 Castro RA, Ross A, Kamwela L, Reinhard M, Loiseau C, Feldmann J, Borrell S, Trauner A,
560 Gagneux S. 2020. The genetic background modulates the evolution of fluoroquinolone-
561 resistance in Mycobacterium tuberculosis. *Mol Biol Evol* **37**:195–207.
- 562 Chiner-Oms Á, Berney M, Boinett C, González-Candelas F, Young DB, Gagneux S, Jacobs WR,
563 Parkhill J, Cortes T, Comas I. 2019. Genome-wide mutational biases fuel transcriptional
564 diversity in the Mycobacterium tuberculosis complex. *Nat Commun* **10**:1–11.
- 565 Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T,
566 Kauff F, Wilczynski B, others. 2009. Biopython: freely available Python tools for
567 computational molecular biology and bioinformatics. *Bioinformatics* **25**:1422–1423.
- 568 Comas I, Chakravarti J, Small PM, Galagan J, Niemann S, Kremer K, Ernst JD, Gagneux S.
569 2010. Human T cell epitopes of Mycobacterium tuberculosis are evolutionarily
570 hyperconserved. *Nat Genet* **42**:498–498.
- 571 de Vos M, Ley SD, Wiggins KB, Derendinger B, Dippenaar A, Grobbelaar M, Reuter A, Dolby
572 T, Burns S, Schito M. 2019. Bedaquiline microheteroresistance after cessation of
573 tuberculosis treatment. *N Engl J Med* **380**:2178–2180.
- 574 Diacon AH, Pym A, Grobusch MP, de los Rios JM, Gotuzzo E, Vasilyeva I, Leimane V, Andries
575 K, Bakare N, De Marez T. 2014. Multidrug-resistant tuberculosis and culture conversion
576 with bedaquiline. *N Engl J Med* **371**:723–732.
- 577 Edwards DJ, Duchêne S, Pope B, Holt KE. 2020. SNPPar: identifying convergent evolution and
578 other homoplasies from microbial whole-genome alignments. *bioRxiv*.
- 579 Ektefaie Y, Dixit A, Freschi L, Farhat MR. 2021. Globally diverse Mycobacterium tuberculosis
580 resistance acquisition: a retrospective geographical and temporal analysis of whole
581 genome sequences. *Lancet Microbe* **2**:e96–e104.
- 582 Fowler PW, CRyPTIC Consortium. 2021. Epidemiological cutoff values for a 96-well broth
583 microdilution plate for high throughput research antibiotic susceptibility testing of M.
584 tuberculosis. *medRxiv*.
- 585 Freschi L, Vargas R, Hussain A, Kamal SM, Skrahina A, Tahseen S, Ismail N, Barbova A,
586 Niemann S, Cirillo DM. 2020. Population structure, biogeography and transmissibility of
587 Mycobacterium tuberculosis. *bioRxiv*.

- 588 Gagneux S. 2018. Ecology and evolution of Mycobacterium tuberculosis. *Nat Rev Microbiol*
589 **16**:202–202.
- 590 Ghodousi A, Rizvi AH, Baloch AQ, Ghafoor A, Khanzada FM, Qadir M, Borroni E, Trovato A,
591 Tahseen S, Cirillo DM. 2019. Acquisition of cross-resistance to bedaquiline and
592 clofazimine following treatment for tuberculosis in Pakistan. *Antimicrob Agents*
593 *Chemother* **63**:e00915-19.
- 594 Groschel MI, Owens M, Freschi L, Vargas R, Marin MG, Phelan J, Iqbal Z, Dixit A, Farhat MR.
595 2021. GenTB: A user-friendly genome-based predictor for tuberculosis resistance
596 powered by machine learning. *bioRxiv*.
- 597 Gygli SM, Keller PM, Ballif M, Blöchliger N, Hömke R, Reinhard M, Loiseau C, Ritter C,
598 Sander P, Borrell S. 2019. Whole-genome sequencing for drug resistance profile
599 prediction in Mycobacterium tuberculosis. *Antimicrob Agents Chemother* **63**.
- 600 Hain Lifescience. 2017. GenoType MTBDRsl VER 2.0 - Molecular Genetic Assay for
601 Identification of the M. tuberculosis Complex and its Resistance to Fluoroquinolones and
602 Aminoglycosides/Cyclic Peptides from Sputum Specimens or Cultivated Samples (No.
603 IFU-317A-04).
- 604 Halloum I, Viljoen A, Khanna V, Craig D, Bouchier C, Brosch R, Coxon G, Kremer L. 2017.
605 Resistance to thiacetazone derivatives active against Mycobacterium abscessus involves
606 mutations in the MmpL5 transcriptional repressor MAB_4384. *Antimicrob Agents*
607 *Chemother* **61**:e02509-16.
- 608 Hariguchi N, Chen X, Hayashi Y, Kawano Y, Fujiwara M, Matsuba M, Shimizu H, Ohba Y,
609 Nakamura I, Kitamoto R. 2020. OPC-167832, a novel carbostyryl derivative with potent
610 antituberculosis activity as a dpre1 inhibitor. *Antimicrob Agents Chemother* **64**.
- 611 Heyckendorf J, Andres S, Köser CU, Olaru ID, Schön T, Sturegård E, Beckert P, Schleusener V,
612 Kohl TA, Hillemann D. 2018. What is resistance? Impact of phenotypic versus molecular
613 drug resistance testing on therapy for multi-and extensively drug-resistant tuberculosis.
614 *Antimicrob Agents Chemother* **62**.
- 615 Hunter JD. 2007. Matplotlib: A 2D graphics environment. *Comput Sci Eng* **9**:90–95.
- 616 Kadura S, King N, Nakhoul M, Zhu H, Theron G, Köser CU, Farhat M. 2020. Systematic review
617 of mutations associated with resistance to the new and repurposed Mycobacterium
618 tuberculosis drugs bedaquiline, clofazimine, linezolid, delamanid and pretomanid. *J*
619 *Antimicrob Chemother*.
- 620 Kalyaanamoorthy S, Minh BQ, Wong TK, Von Haeseler A, Jermiin LS. 2017. ModelFinder: fast
621 model selection for accurate phylogenetic estimates. *Nat Methods* **14**:587–589.
- 622 Kambli P, Ajbani K, Nikam C, Sadani M, Shetty A, Udawadia Z, Georghiou SB, Rodwell TC,
623 Catanzaro A, Rodrigues C. 2016. Correlating rrs and eis promoter mutations in clinical
624 isolates of Mycobacterium tuberculosis with phenotypic susceptibility levels to the
625 second-line injectables. *Int J Mycobacteriology* **5**:1–6.
- 626 Kaniga K, Cirillo DM, Hoffner S, Ismail NA, Kaur D, Lounis N, Metchock B, Pfyffer GE,
627 Venter A. 2016. A multilaboratory, multicountry study to determine bedaquiline MIC
628 quality control ranges for phenotypic drug susceptibility testing. *J Clin Microbiol*
629 **54**:2956–2962.
- 630 Köser CU, Bryant JM, Parkhill J, Peacock SJ. 2013. Consequences of whiB7 (Rv3197A)
631 mutations in Beijing genotype isolates of the Mycobacterium tuberculosis complex.
632 *Antimicrob Agents Chemother* **57**:3461–3461.

- 633 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows--Wheeler transform.
634 *Bioinformatics* **25**:1754–1760.
- 635 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.
636 2009. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**:2078–2079.
- 637 Ma Z, Lienhardt C, McIlleron H, Nunn AJ, Wang X. 2010. Global tuberculosis drug
638 development pipeline: the need and the reality. *The Lancet* **375**:2100–2109.
- 639 McKinney W, others. 2010. Data structures for statistical computing in python Proceedings of the
640 9th Python in Science Conference. pp. 51–56.
- 641 Merker M, Kohl TA, Barilar I, Andres S, Fowler PW, Chryssanthou E, Ängeby K, Jureen P,
642 Moradigaravand D, Parkhill J. 2020. Phylogenetically informative mutations in genes
643 implicated in antibiotic resistance in Mycobacterium tuberculosis complex. *Genome Med*
644 **12**:1–8.
- 645 Minh BQ, Nguyen MAT, von Haeseler A. 2013. Ultrafast approximation for phylogenetic
646 bootstrap. *Mol Biol Evol* **30**:1188–1195.
- 647 Mohamed S, Köser CU, Salfinger M, Sougakoff W, Heysell SK. 2021. Targeted next-generation
648 sequencing: a Swiss army knife for mycobacterial diagnostics?
- 649 Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective
650 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*
651 **32**:268–274.
- 652 Nimmo C, Millard J, Brien K, Moodley S, van Dorp L, Lutchminarain K, Wolf A, Grant AD,
653 Balloux F, Pym AS. 2020. Bedaquiline resistance in drug-resistant tuberculosis HIV co-
654 infected patients. *Eur Respir J* **55**.
- 655 Peretokina IV, Krylova LY, Antonova OV, Kholina MS, Kulagina EV, Nosova EY, Safonova
656 SG, Borisov SE, Zimenkov DV. 2020. Reduced susceptibility and resistance to
657 bedaquiline in clinical M. tuberculosis isolates. *J Infect* **80**:527–535.
- 658 Pérez F, Granger BE. 2007. IPython: a system for interactive scientific computing. *Comput Sci*
659 *Eng* **9**.
- 660 Pholwat S, Stroup S, Heysell S, Ogarkov O, Zhdanova S, Ramakrishnan G, Houpt E. 2016. eis
661 promoter C14G and C15G mutations do not confer kanamycin resistance in
662 Mycobacterium tuberculosis. *Antimicrob Agents Chemother* **60**:7522–7523.
- 663 Reeves AZ, Campbell PJ, Sultana R, Malik S, Murray M, Plikaytis BB, Shinnick TM, Posey JE.
664 2013. Aminoglycoside cross-resistance in Mycobacterium tuberculosis due to mutations
665 in the 5' untranslated region of whiB7. *Antimicrob Agents Chemother* **57**.
- 666 Richardson E, Lin S, Pinsky B, Desmond E, Banaei N. 2009. First documentation of isoniazid
667 reversion in Mycobacterium tuberculosis. *Int J Tuberc Lung Dis* **13**:1347–1354.
- 668 Safi H, Lingaraju S, Amin A, Kim S, Jones M, Holmes M, McNeil M, Peterson SN, Chatterjee
669 D, Fleischmann R. 2013. Evolution of high-level ethambutol-resistant tuberculosis
670 through interacting mutations in decaprenylphosphoryl- β -D-arabinose biosynthetic and
671 utilization pathway genes. *Nat Genet* **45**:1190–1197.
- 672 Sanz-García F, Anoz-Carbonell E, Pérez-Herrán E, Martín C, Lucía A, Rodrigues L, Aínsa JA.
673 2019. Mycobacterial aminoglycoside acetyltransferases: a little of drug resistance, and a
674 lot of other roles. *Front Microbiol* **10**.
- 675 Schmieder R, Edwards R. 2011. Quality control and preprocessing of metagenomic datasets.
676 *Bioinformatics* **27**:863–864.
- 677 Schön T, Köser CU, Werngren J, Viveiros M, Georghiou S, Kahlmeter G, Giske C, Maurer F,
678 Lina G, Turnidge J. 2020. What is the role of the EUCAST reference method for MIC

- 679 testing of the Mycobacterium tuberculosis complex? *Clin Microbiol Infect Off Publ Eur*
680 *Soc Clin Microbiol Infect Dis* **26**:1453–1455.
- 681 Schön T, Matuschek E, Mohamed S, Utukuri M, Heysell S, Alffenaar J-W, Shin S, Martinez E,
682 Sintchenko V, Maurer F. 2019. Standards for MIC testing that apply to the majority of
683 bacterial pathogens should also be enforced for Mycobacterium tuberculosis complex.
684 *Clin Microbiol Infect* **25**:403–405.
- 685 Seabold S, Perktold J. 2010. Statsmodels: Econometric and statistical modeling with
686 pythonProceedings of the 9th Python in Science Conference. pp. 61–61.
- 687 Sonnenkalb L, Carter J, Spitaleri A, Iqbal Z, Hunt M, Malone K, Utpatel C, Cirillo DM,
688 Rodrigues C, Nilgiriwala KS. 2021. Deciphering Bedaquiline and Clofazimine
689 Resistance in Tuberculosis: An Evolutionary Medicine Approach. *bioRxiv*.
- 690 Van Der Walt S, Colbert SC, Varoquaux G. 2011. The NumPy array: a structure for efficient
691 numerical computation. *Comput Sci Eng* **13**:22–22.
- 692 Vargas R, Freschi L, Marin M, Epperson LE, Smith M, Oussenko I, Durbin D, Strong M,
693 Salfinger M, Farhat MR. 2021. In-host population dynamics of Mycobacterium
694 tuberculosis complex during active disease. *Elife* **10**:e61805.
- 695 Viljoen A, Dubois V, Girard-Misguich F, Blaise M, Herrmann J, Kremer L. 2017. The diverse
696 family of MmpL transporters in mycobacteria: from regulation to antimicrobial
697 developments. *Mol Microbiol* **104**:889–904.
- 698 Villellas C, Coeck N, Meehan CJ, Lounis N, de Jong B, Rigouts L, Andries K. 2017. Unexpected
699 high prevalence of resistance-associated Rv0678 variants in MDR-TB patients without
700 documented prior use of clofazimine or bedaquiline. *J Antimicrob Chemother* **72**:684–
701 690.
- 702 Viney K, Linh NN, Gegia M, Zignol M, Glaziou P, Ismail N, Kasaeva T, Mirzayev F. 2021.
703 New definitions of pre-extensively and extensively drug-resistant tuberculosis: update
704 from the World Health Organization.
- 705 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
706 Wortman J, Young SK, others. 2014. Pilon: an integrated tool for comprehensive
707 microbial variant detection and genome assembly improvement. *PLoS One* **9**:e112963–
708 e112963.
- 709 Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification using
710 exact alignments. *Genome Biol* **15**:R46–R46.
- 711 World Health Organization. 2020. Global Tuberculosis Report. World Health Organization.
- 712 World Health Organization. 2018. Technical report on critical concentrations for drug
713 susceptibility testing of medicines used in the treatment of drug-resistant tuberculosis.
714 World Health Organization.
- 715 World Health Organization. In press. Catalogue of mutations in Mycobacterium tuberculosis
716 complex associated with drug resistance phenotypes. World Health Organization.
- 717 Xu J, Converse PJ, Upton AM, Mdluli K, Fotouhi N, Nuermberger EL. 2021. Comparative
718 efficacy of the novel diarylquinoline TBAJ-587 and bedaquiline against a resistant
719 Rv0678 mutant in a mouse model of tuberculosis. *Antimicrob Agents Chemother* **65**.
- 720 Xu J, Wang B, Fu L, Zhu H, Guo S, Huang H, Yin D, Zhang Y, Lu Y. 2019. In vitro and in vivo
721 activities of the riminophenazine TBI-166 against Mycobacterium tuberculosis.
722 *Antimicrob Agents Chemother* **63**.

- 723 Yamamoto K, Nakata N, Mukai T, Kawagishi I, Ato M. 2021. Coexpression of MmpS5 and
724 MmpL5 Contributes to Both Efflux Transporter MmpL5 Trimerization and Drug
725 Resistance in *Mycobacterium tuberculosis*. *Msphere* **6**.
- 726 Zaunbrecher MA, Sikes RD, Metchock B, Shinnick TM, Posey JE. 2009. Overexpression of the
727 chromosomally encoded aminoglycoside acetyltransferase eis confers kanamycin
728 resistance in *Mycobacterium tuberculosis*. *Proc Natl Acad Sci* **106**:20004–20009.
- 729 Zhang S, Chen J, Cui P, Shi W, Zhang W, Zhang Y. 2015. Identification of novel mutations
730 associated with clofazimine resistance in *Mycobacterium tuberculosis*. *J Antimicrob*
731 *Chemother* **70**:2507–2510.
- 732

733 **ACKNOWLEDGEMENTS**

734 We thank Koné Kaniga and Nacer Lounis for sharing information about the C208 trial conducted
735 by Janssen and Thomas Schön for helpful discussions regarding aminoglycoside resistance. We
736 thank the members of the Farhat lab for helpful discussions and comments on the research project
737 and manuscript. R.V.J. was supported by the National Science Foundation Graduate Research
738 Fellowship under Grant No. DGE1745303. C.U.K. received an observership from the European
739 Society of Clinical Microbiology and Infectious Diseases. M.R.F. was supported by NIH NIAID
740 R01 AI55765. Portions of this research were conducted on the O2 High Performance Compute
741 Cluster, supported by the Research Computing Group, at Harvard Medical School.

742

743 **COMPETING INTERESTS**

744 C.U.K.'s work for Becton Dickinson involves a collaboration with Janssen and Thermo Fisher
745 Scientific. C.U.K. is a consultant for Becton Dickinson, the Foundation for Innovative New
746 Diagnostics, the Stop TB Partnership, and the TB Alliance. C.U.K. worked as a consultant for
747 QuantuMDx, the WHO Global TB Programme, and WHO Regional Office for Europe. C.U.K.
748 gave a paid educational talk for Oxford Immunotec. Hain Lifescience covered C.U.K.'s travel and
749 accommodation to present at a meeting. C.U.K. is an unpaid advisor to BioVersys and
750 GenoScreen.

751

752 **DATA AND MATERIALS AVAILABILITY**

753 Mtb sequencing data was collected from NCBI and is publicly available. WGS data for the set of
754 added 12 clinical *eis* C-14T mutants (**Materials and Methods**) will be uploaded to a public
755 sequence repository upon acceptance of this manuscript for publication. All packages and software

756 used in this study have been noted in the **Materials and Methods**. Custom scripts written in
757 python version 2.7.15 were used to conduct all analyses and interfaced via Jupyter Notebooks. All
758 scripts and notebooks will be uploaded to a GitHub repository upon acceptance of this manuscript
759 for publication.