

1 **Direct whole genome sequencing of sputum accurately identifies drug resistant**

2 ***Mycobacterium tuberculosis* faster than MGIT culture sequencing**

3

4 Ronan M. Doyle<sup>1,2\*</sup>, Carrie Burgess<sup>1</sup>, Rachel Williams<sup>1</sup>, Rebecca Gorton<sup>3</sup>, Helen Booth<sup>4,5</sup>,

5 James Brown<sup>6,12</sup>, Josephine M. Bryant<sup>7</sup>, Jackie Chan<sup>8</sup>, Dean Creer<sup>6</sup>, Jolyon Holdstock<sup>8</sup>,

6 Heinke Kunst<sup>9</sup>, Stefan Lozewicz<sup>10</sup>, Gareth Platt<sup>3</sup>, Erika Yara Romero<sup>1</sup>, Graham Speight<sup>8</sup>,

7 Simon Tiberi<sup>9</sup>, Ibrahim Abubakar<sup>11</sup>, Marc Lipman<sup>6,12</sup>, Timothy D. McHugh<sup>3</sup>, Judith Breuer<sup>1</sup>

8

9 1. Division of Infection and Immunity, University College London, London, UK

10 2. Microbiology, Virology and Infection Control, Great Ormond Street Hospital NHS  
11 Foundation Trust, London, UK

12 3. Centre for Clinical Microbiology, Division of Infection and Immunity, Royal Free  
13 Campus, UCL, London, UK

14 4. University College London Hospitals NHS Foundation Trust, London, UK

15 5. North Central London TB Service-South Hub, Whittington Hospital NHS Trust,  
16 London, UK

17 6. Royal Free London NHS Foundation Trust, London UK

18 7. Molecular Immunity Unit, MRC Laboratory of Molecular Biology, Department of  
19 Medicine, University of Cambridge, Cambridge, UK

20 8. Oxford Gene Technology, Oxford Begbroke Science Park, Begbroke, Oxfordshire,  
21 UK

22 9. Blizard Institute, Barts and The London School of Medicine and Dentistry, Queen  
23 Mary University of London, UK

24 10. North Middlesex University Hospital NHS Trust, London, UK

25 11. UCL Institute for Global Health, University College London, London, UK

26 12. UCL Respiratory, Division of Medicine, University College London, London, UK

27

28 **Running title:** WGS of drug resistant TB directly from sputum

29 **\*Corresponding Author:** Ronan M. Doyle, [ronan.doyle@gosh.nhs.uk](mailto:ronan.doyle@gosh.nhs.uk)

30

### 31 **Abstract**

32 The current methods available to diagnose antimicrobial resistant *Mycobacterium*  
33 *tuberculosis* infections require positive culture or only test a limited number of resistance-  
34 associated mutations. Rapid, accurate identification of antimicrobial resistance enables  
35 prompt initiation of effective treatment. Here, we determine the utility of whole-genome  
36 sequencing (WGS) *M. tuberculosis* directly from routinely obtained diagnostic sputum  
37 samples to provide a comprehensive resistance profile compared to Mycobacterial Growth  
38 Indicator Tube (MGIT) WGS. We sequenced *M. tuberculosis* from 43 sputum samples by  
39 targeted DNA enrichment using the Agilent SureSelectXT kit, and 43 MGIT positive samples  
40 from each participant. Thirty two (74%) sputum samples and 43 (100%) MGIT samples  
41 generated whole genomes. Time to antimicrobial resistance profile and concordance was  
42 compared with Xpert MTB/RIF and phenotypic resistance testing from culture of the same  
43 samples. Antibiotic susceptibility could be predicted from WGS of sputum within 5 days of  
44 sample receipt and up to 24 days earlier than WGS from MGIT culture and up to 31 days  
45 earlier than phenotypic testing. Direct sputum results could be reduced to 3 days with faster  
46 hybridisation and if only regions encoding drug resistance are sequenced. We show that  
47 direct sputum sequencing has the potential to provide comprehensive resistance detection

48 significantly faster than MGIT whole genome sequencing or phenotypic testing of resistance  
49 from culture in a clinical setting. This improved turnaround time enables prompt, appropriate  
50 treatment with associated patient and health service benefit. Improvements in sample  
51 preparation are necessary to ensure comparable sensitivity and complete resistance profile  
52 predictions in all cases.

53

54 **Keywords:** *Mycobacterium tuberculosis*; Whole-genome sequencing; Pathogen DNA  
55 enrichment; Antimicrobial resistance.

56

## 57 **Introduction**

58 Tuberculosis (TB) infection is a global emergency associated with an increasing burden of  
59 drug resistant *Mycobacterium tuberculosis* complex infections (1). Phenotypic testing for  
60 antimicrobial resistance detection is slow with results typically a month to six weeks after  
61 initial culture confirmation - leading to the potential for prolonged, suboptimal antibiotic  
62 treatment. Molecular assays such as the Xpert MTB/RIF (Cepheid), MTBDRplus and  
63 MTBDRs/ (Hain Lifescience) can rapidly detect a limited number of first and second line  
64 drug resistance mutations (2). However, none are currently able to identify the full range of  
65 antibiotic resistance mutations needed for appropriately targeted therapy in people with  
66 multi-drug resistant (MDR) TB. Further, these assays recognise only a fixed number of target  
67 mutations, missing less common resistance mutations; (3) whilst Xpert MTB/RIF can only  
68 detect DNA mutations and not predict amino acid changes, resulting in potential false  
69 positives (4).

70 Whole-genome sequencing (WGS) of *M. tuberculosis* allows comprehensive identification of  
71 all known drug resistant mutations for all classes of TB drugs and also can provide valuable  
72 contact tracing information (5). Recently the sequencing of organisms cultured in  
73 Mycobacterial Growth Indicator Tubes (MGIT) has been shown to be both an accurate  
74 method for detecting first and second line resistance mutations across the genome and  
75 cheaper than present routine diagnostic workflows (6). Although it is being rolled out in  
76 England, (7) it relies on bacterial culture which can delay the time to result by several weeks.

77 We have previously described a successful method for capturing *M. tuberculosis* DNA  
78 directly from sputum samples using biotinylated RNA baits (8). This protocol provides a  
79 possible faster alternative to sequencing *M. tuberculosis* whole genomes and could therefore  
80 offer quicker diagnosis of antibiotic resistance, leading to tailored treatment regimens with  
81 less use of antimicrobials and associated toxicity, fewer days in hospital, reduced cost and  
82 improved outcomes.

83 Mixed strain infections of *M. tuberculosis* are well-documented (9) and may lead to poor  
84 treatment outcomes and the possible emergence of minority drug resistant strains (10–12).  
85 Culture of *M. tuberculosis* is known to impact negatively on detection of mixtures and  
86 minority variant mutations (13), with short term MGIT culture being particularly poor at  
87 identifying mixed infections (14).

88 The aims of this study were: (1) to compare the utility of performing WGS directly from  
89 routinely-obtained diagnostic sputum with MGIT samples taken from the same participant  
90 (time to diagnosis plus their ability to predict antimicrobial resistance, AMR); (2) identify  
91 mixed infections and minority populations within samples.

92

## 93 **Materials and methods**

### 94 **Study recruitment**

95 Individuals aged 16 years or older attending a TB service with suspected pulmonary TB at  
96 seven clinics in London, UK were invited to take part in this study.

### 97 **DNA extraction**

98 DNA was extracted from 1ml clinical samples and MGIT cultures using mechanical ribolysis  
99 and automated DNA extraction workflow. Samples were centrifuged for 30 minutes at 16,200  
100 x g and the supernatant discarded. For MGIT cultures only, a saline prewash method was  
101 utilised to reduce the human nucleic acid component of the sample (15). 1ml of sterile saline  
102 was added to the pellet (0.9% w/v), the pellet was re-suspended and centrifuged for 15  
103 minutes at max speed (16,200 x g). The supernatant was discarded and the process was  
104 repeated. For sputum and MGIT cultures approximately 50ul of glass beads (425-600µm)  
105 were added to each sample pellet and ribolysis was performed on a FastPrep24 platform for  
106 45 seconds at 6.4 m/s. 240ul of extraction buffer 2 and 10ul of proteinase K was added to  
107 each sample, vortexed then incubated at 56°C for 10 minutes. DNA was extracted from  
108 samples lysates on the Diasorin IXT (Arrow) automated platform using DNA extraction  
109 cartridges eluting into 100ul.

### 110 **Quantification of extracted *Mycobacterium tuberculosis* DNA**

111 The Xpert MTB/RIF (Cepheid) assay was performed on sputum samples as per  
112 manufacturer's instructions; reporting the *M. tuberculosis* (MTB) quantity as either very low,  
113 low, medium or high alongside C<sub>T</sub> values. The Xpert MTB/RIF assay also reported rifampicin  
114 resistance as 'detected' or 'not detected'. Drug susceptibility testing was based on phenotypic  
115 culture for first-line drugs on solid media using the resistance ratio method and was carried

116 out by the National Mycobacterium Reference Service using their standard protocols. A  
117 second MTB specific qPCR targeting the 16S rRNA gene (*rrs*) was utilised to quantify the  
118 MTB DNA extracted from sputum samples and MGIT cultures. For MGIT culture extracts a  
119 1/1000 dilution was prepared prior to qPCR analysis. qPCR was performed using forward  
120 primer 5'-GTGATCTGCCCTGCACCTC-3' and reverse 5'-  
121 ATCCCACACCGCTAAAGCG-3' with a TaqMan probe ROX-  
122 AGGACCACGGGATGCATGTCTTGT-BHQ2 (16). The MTB specific qPCR reaction  
123 consisted of 12.5µl of Quantitect Multiplex NoROX mix (Qiagen), 0.2µM primers and probes  
124 and 5µl template per reaction in a total volume of 25µl. Reactions were performed in  
125 duplicate on a Rotorgene 8000 platform. PCR cycling conditions were as follows: 50°C for  
126 30 mins, 95°C for 15 mins and 40 cycles of 94°C for 45 secs and 60°C for 45 secs. Standards  
127 were prepared from commercially sourced MTB genomic DNA (Vircell), reconstituted as  
128 directed by the manufacturers.

### 129 **Sequencing library preparation and whole genome sequencing**

130 Total DNA was quantified in sputum and MGIT extracts using the Qubit High Sensitivity  
131 DNA assay (Life Technologies). Carrier human genomic DNA (Promega) was added where  
132 needed to obtain a total of 200 ng of DNA input for library preparation. All DNA samples  
133 were sheared using a Covaris S2 ultrasonicator for 150 seconds (PIP 175; duty factor 5; 200  
134 cycles per burst using frequency sweeping). Sputum samples were prepared using the  
135 SureSelectXT target enrichment system for the Illumina paired-end sequencing library  
136 protocol (Agilent Technologies). End repair, 3' addition of adenosine and ligation of adapters  
137 were all carried out according to Agilent's protocol. Prior to hybridisation, 12 cycles of  
138 precapture PCR were used using primers provided in the SureSelectXT kit. Hybridisation of  
139 MTB DNA to the streptavidin-coated beads was carried out using a MTB specific bait set  
140 described previously (8). Briefly, 120-mer RNA baits were designed to provide non-

141 redundant coverage of the entire length of the positive strand of the H37Rv reference  
142 genome, they were synthesized by Agilent Biotechnologies. The baits can be purchased from  
143 Agilent and the bait sequences are available upon request from authors. 18 cycles of  
144 postcapture PCR were performed with indexing primers provided in the SureSelectXT kit.  
145 All Agilent recommended quality control steps were carried out. In order to compare the  
146 effect of target enrichment on MGIT sequencing, the first 14 MGIT samples were underwent  
147 library preparation using SureSelectXT and all subsequent MGIT samples had DNA libraries  
148 prepared using the NEBNext Ultra II DNA Library Prep Kit (NEB) as per the manufacturer's  
149 protocol. The resulting DNA libraries were run on either a MiSeq or NextSeq sequencer  
150 (Illumina) using either a V2 500-cycle or 500/550 Mid Output 300-cycle kit, respectively.

#### 151 **Optimised sequencing method**

152 Three samples were prepared using SureSelectXT Fast Target Enrichment System (Agilent  
153 Technologies) as per manufacturer's protocol. A reduced bait set was designed to capture  
154 only genes associated with drug resistance and information for spoligotyping was used in  
155 hybridisation step. The reduced set of 120-mer RNA baits were synthesised by Agilent  
156 Technologies in the same way as the full set except that it only included baits that were  
157 complementary to the genes and regions in Table 1 of the H37Rv reference genome.

#### 158 **Bioinformatic analysis**

159 Sequencing reads were trimmed for adapter content and quality using Trim Galore, keeping  
160 reads longer than 100bp. Trimmed reads were deduplicated and mapped to the H37Rv  
161 (accession: NC\_000962) reference genome using BBmap allowing only successfully mapped  
162 paired reads at the 99% equivalent minimum identity across the entire read and a maximum  
163 insert size of 500bp. Duplicate mapped reads were removed using Picard tools and variants  
164 against the reference genome were called with freebayes keeping only variants with a

165 minimum of 10 supporting reads, greater than 2% frequency, mapping quality greater than 20  
166 and base quality score greater than 30, with reads present on both the forward and reverse  
167 strand, and on both the 5' or 3' end of reads. Variants found in and within 100bp of Pro-Glu  
168 (PE) and Pro-Pro-Glu (PPE) genes, mobile elements and repeat regions were discarded. For  
169 resistance calling single nucleotide variants (SNVs) were annotated using ANNOVAR (17).  
170 A maximum likelihood phylogeny was also inferred from 1113 core genome SNVs present in  
171 64 samples representing 32 participants using RAxML (v. 8.2.1) (18) with 99 bootstrap  
172 replicates. SNV distance between pairs was calculated using R package *seqinr*. The same  
173 filtering conditions were also applied to variants for minor variant analysis. The number of  
174 reads across each variant position was normalised between pairs of samples from the same  
175 patient to adjust for the effect of read depth on variant frequency. Minor variants were filtered  
176 from the dataset if found on reads with greater sequence identity to a different MTB complex  
177 species. To further control possible contamination of paired samples with low frequency  
178 variants, MGIT and sputum samples from each pair were extracted on separate days,  
179 prepared in separate sequencing libraries and sequenced on different runs.

#### 180 **Ethics approval and consent to participate**

181 Samples were collected with informed consent from patients attending a TB clinic setting at  
182 the participating hospitals. Approval for the study was granted by the NRES Committee East  
183 Midlands – Nottingham 1 (REC reference: 15/EM/0091). All samples were pseudo-  
184 anonymised and allocated a unique identification number.

#### 185 **Data availability**

186 All sequence data associated with this study has been deposited in the European Nucleotide  
187 Archive under study accession number PRJEB21685.



188

**189 Results****190 Genomic coverage**

191 Sixty three participants were prospectively enrolled. A paired sputum and MGIT sample was  
192 sequenced from 43 patients. This is due to ten participants not having a MGIT sample  
193 collected for sequencing, another eight samples being smear and Xpert MTB/RIF negative,  
194 and two where the volume of sputum was insufficient for DNA extraction. Samples  
195 sequenced from MGIT culture had a higher reference genome coverage as compared with  
196 those obtained directly from sputum (Fig. 1) and this was correlated to the increased *M.*  
197 *tuberculosis* DNA available from the former (Fig S1). Enrichment of MGIT culture samples  
198 using the *M. tuberculosis* probes also enhanced the quality and depth of sequence (Fig S1).

199 We next evaluated whether bacterial load, as measured by smear and Xpert MTB/RIF, could  
200 be used to predict the success of whole genome sequencing. From 43 patients, 32 sputum  
201 samples (74.4%) and 43 MGIT samples (100%) generated whole genomes (>85% coverage  
202 against reference genome) (Fig. 2A). Sputum sequencing success was linked to estimated  
203 input pathogen copy number. We stratified participants into 16 with high (3+ smear result,  
204 Xpert MTB/RIF High), 18 with medium (2+ smear result, Xpert MTB/RIF Medium) and 9  
205 low bacterial load (scanty or 1+ smear result, Xpert MTB/RIF Low) and found 87.5% of  
206 sputum samples with high bacterial load samples generated complete genomes as compared  
207 to 72.2% with medium and 55.5% with low bacterial load (Fig. 2B & 2C). We were also able  
208 to recover partial genomes for two sputum samples that were reported as negative by smear  
209 microscopy but Xpert MTB/RIF positive.

**210 MGIT and sputum sequence variation**

211 Comparison of the 32 patients with both complete sputum and MGIT genomes available,  
212 showed no unique consensus sequence variation between the pairs. The identity between  
213 sputum and MGIT consensus sequences is shown in a heat map (Fig. 3) and phylogenetic tree  
214 (Fig. S2). Twenty three MGIT and sputum pairs showed no SNV between them at the  
215 consensus level, while nine patients' sample pairs differed by one or two single nucleotides.  
216 In all nine patients the consensus polymorphism was present as a minority variant in the  
217 matched sputum or MGIT sample (Table S1).

### 218 **Time to antibiotic resistance prediction**

219 Allowing for sample batching, which was only carried out for study purposes, direct  
220 sequencing of sputum using targeted enrichment reduced the time to antibiotic susceptibility  
221 prediction initially to five days, as compared with a mean of 11 (s.d. 6) days for MGIT  
222 sequencing (Fig. 4). This was reduced further by protocol optimisation to three days when a  
223 reduced bait set was used that captures only the regions with putative resistance mutations  
224 (Table 1 & Fig. 4). Hybridisation optimisation could also reduce the whole genome protocol  
225 to 4 days (Fig. 4). The reduced bait set targeted 35,960bp of the H37Rv reference genome in  
226 total and successfully re-sequenced three MDR-TB samples (noted in Fig 4.) from this study  
227 to a high average depth of coverage (>2,000X) over the captured regions. All eight genotypic  
228 resistant variants identified in the whole genome sequencing data were also identified after  
229 re-sequencing with the reduced bait set (Fig. S3). Overall, 36 sputum samples with complete  
230 genomes, including 77% of those with drug resistance mutations, would have been reported a  
231 mean of 9 days earlier than MGIT sequencing and a mean of 35 days earlier than phenotypic  
232 testing using the optimised three day protocol.

### 233 **Antibiotic resistance concordance**

234 We found complete concordance between resistance mutations identified in paired MGIT and  
235 sputum samples from nine participants when there was >85% single read coverage against  
236 reference genome. Four participants missed resistance mutations where sputum sequencing  
237 read coverage was too low to make a reliable call (Table 2). Xpert MTB/RIF, sputum WGS,  
238 MGIT WGS were concordant with phenotypic resistance testing in 21 out of 23 resistance  
239 mutations identified (Table 2). The exceptions were where a variant in participant RF015GT  
240 predicting an amino acid change Ser428Iso in *rpoB* (H37Rv codon numbering, *Escherichia*  
241 *coli rpoB* numbering S509I) reported as resistant by Xpert MTB/RIF, but the reference  
242 laboratory found it to be susceptible (Table 2). Previous publications have shown not all  
243 SNVs at this position are associated with resistance (19–21). A fixed mutation in *f* predicting  
244 an amino acid change Ser315Thr was confirmed to confer high levels of resistance to  
245 isoniazid in five samples, but patient BH052SA with the same mutation was found to be  
246 susceptible. This common polymorphism in *katG* has been previously been shown to confer  
247 consistently high levels of isoniazid resistance to *M. tuberculosis* (22–24). Whole genome  
248 sequencing from both MGIT and sputum samples also identified streptomycin (three patients)  
249 and Para-aminosalicylic acid (PAS, four patients) resistance mutations, neither of which is  
250 routinely tested within the phenotypic assay in the UK.

#### 251 **Mixed infections and minority variants**

252 No mixed infections were detected. Using data normalised for read depth, from 32 matched  
253 sputa-MGIT samples, minor frequency variation was low with only 88 minority bi-allelic  
254 sites meeting the quality criteria identified in all samples, representing 0.002% unique  
255 variable positions across the genomes. We undertook stringent procedures to exclude  
256 sequencing error, the presence of closely related *M. tuberculosis* complex species in sputum,  
257 contamination in and between sequencing runs, and lab contamination after sample collection  
258 as potential causes for the findings. None of the variant alleles were at positions known to be

259 associated with antimicrobial resistance. In 41% (13) of cases directly sequenced sputa had  
260 higher numbers of novel minority variants identified than the matched MGIT as compared to  
261 19% (6) of MGIT samples with more minority variants in the sputum (Fig. S4).

262 To control for the potential influence of the SureSelectXT step, we analysed the proportion of  
263 minority variants shared between sputum and eight SureSelectXT enriched MGIT samples  
264 (40%) as compared between sputum and 16 non-enriched samples (38%) and found no  
265 statistically significant difference ( $p=0.854$ ). Overall, 37.2% of minority variants were  
266 concordant between sputum and MGIT and the read frequencies with which they occurred  
267 were weakly correlated (Fig. S5). This correlation was skewed by one patient (WH044IL)  
268 whose MGIT and sputum samples were both more variable than other samples (Fig. S4) and  
269 in whom seven variants were at much higher frequency in the MGIT than in sputum (at 30-  
270 45% versus ~5% respectively, Fig. S5 & Fig. S6). However, there was no evidence of mixed  
271 genotypes. The five synonymous minor variants and two non-synonymous variants occurring  
272 in two genes of unknown function (Rv3529c and Rv3888c), were distributed across the  
273 genome, and were not shared across any other pairs of samples.

274

## 275 Discussion

276 We have shown that using target enrichment WGS methodology directly from diagnostic  
277 sputum samples generates resistance data, at most, up to 24 days earlier than MGIT culture  
278 WGS and up to 31 days faster than phenotypic testing of *Mycobacterium tuberculosis*.  
279 Sputum sequencing only achieved whole genome sequences suitable for predicting resistance  
280 mutations in 32/43 (74%), though this included a smear negative sputum sample. Our  
281 demonstration that the quality of sequence data is strongly correlated with the input level of  
282 TB DNA (Fig. S1), means that the success of sequencing can be predicted using semi-

283 quantitative methods such as smear microscopy and Xpert MTB/RIF (Fig. 2),  
284 notwithstanding their variable performance (25).

285 Our data compare well with a recent report describing WGS of sputum where contaminating  
286 human DNA had been depleted (26). While this method achieved a slightly faster turnaround  
287 time on diagnostic samples (2 versus 5 days), it may be less susceptible than targeted  
288 enrichment, as only 60% (24/40) of smear positive samples yielded sequence data suitable  
289 for resistance prediction. The study did not report bacterial load so a thorough comparison of  
290 sensitivity cannot be performed. The methods are, however, highly complementary and  
291 combining the two would likely improve genome copy input and increase direct WGS  
292 sensitivity. Our experience of enrichment methods (27, 28) also predicts that redesign of the  
293 first generation probe set would further improve detection of resistance mutations.

294 Direct sequencing of sputum is currently slower than rapid methods such as the Xpert  
295 MTB/RIF and the Hain MTBDRplus and MTBDRsl assays for detecting resistance. However  
296 there are major advantages using WGS. First, unlike existing rapid methods, it can accurately  
297 identify the precise nucleotide change causing resistance. In our study Xpert MTB/RIF  
298 reported resistance in a susceptible organism where there was a nucleotide change at position  
299 428 in the *rpoB* gene not associated with resistance (19–21). Where discordant resistance  
300 results were found between molecular methods and the routine phenotypic testing, we could  
301 not repeat the phenotypic testing as this was carried out historically by a centralised reference  
302 laboratory. Therefore these discrepancies cannot be confirmed and this represents a limitation  
303 of this type of study. Second, WGS, unlike rapid methods which target specific mutations, is  
304 able to detect resistance mutations for a wider range of second and third line drugs, and also  
305 new drugs where current rapid tests would require costly redesign. This has already helped us  
306 personalise treatment in a case of drug resistant *M. tuberculosis* (29). Third, the data from

307 direct WGS of sputum can report evolutionary relationships between samples (Fig. S2) -  
308 providing the most detailed information on transmission dynamics available.

309 An important objective of our study was to evaluate the potential for direct sequencing from  
310 sputum to detect mixed *M. tuberculosis* infections which are sub-optimally identified by  
311 MGIT and solid culture (13, 14). Mixed infections are important in the pathology of TB and  
312 the ability to detect resistance variants that are not at consensus level, although not  
313 necessarily common (30), can affect antibiotic stewardship (31). Mixed infections and MDR-  
314 TB are more prevalent in countries with a much higher burden of TB than the UK and a  
315 greater prevalence of drug resistance, whilst high levels of HIV amplify the problem (32, 33).

316 We were able to detect significantly more minority single nucleotide variants (SNVs) in  
317 sputum compared to the matched MGIT sequence (Fig. S5), despite the mostly clonal  
318 populations in this study and the greater read depths achieved from MGIT sequencing. The  
319 origin of this heterogeneity remains unconfirmed, although we rigorously excluded  
320 contamination and methodological error. SNVs could be due to the presence in sputum of  
321 non-tuberculous mycobacteria and other species which are known to have sequence  
322 homology with *M. tuberculosis* and may theoretically be detected by targeted enrichment.

323 However, our use of highly stringent sequence mapping and the fact that the SNVs were  
324 detected across the genome and not concentrated in regions generally associated with cross  
325 hybridisation, suggests that they are real. In case WH044IL, seven SNVs present in sputum  
326 increased in frequency in MGIT culture, possibly reflecting a selective growth advantage for  
327 this haplotype, particularly as one non-synonymous SNV occurred in the Rv3888c gene  
328 which has been shown to be essential for mycobacterial *in vitro* growth (34). This result  
329 confirms suggestions that diversity is lost and that culture-related selection of some variants  
330 can occur even during limited MGIT culture. Thus MGIT culture may not be representative  
331 of the original sample, and could potentially reduce the likelihood of identifying low level

332 resistance mutations and mixed infections that may act as a reservoir for resistance  
333 development.

334 The standard diagnostic workflows for *M. tuberculosis* are costly and time-consuming.  
335 Ground breaking work within Public Health England has demonstrated that sequencing *M.*  
336 *tuberculosis* whole genomes from positive MGIT cultures is faster and cheaper (6) and where  
337 sputum pathogen DNA concentration is low, early MGIT sequencing could still be the best  
338 possible workaround (15). Direct sequencing of *M. tuberculosis* from sputum has the  
339 potential to reduce the time to antimicrobial resistance detection within a clinically relevant  
340 timeframe (26). We show here that its success is critically dependent on the input genome  
341 copies of pathogen DNA. While enrichment increases the cost of pathogen sequencing, this  
342 could be offset, as demonstrated in our study, by only enriching areas of interest on the  
343 genome. It is important to note that the infrastructure and expertise for rapid, high  
344 throughput, targeted enrichment sequencing directly from clinical material of *M. tuberculosis*  
345 and other pathogens already exists in genomic centres where cancer and genetic disease  
346 sequencing uses this methodology. We believe, therefore, that effective scale-up and  
347 implementation of this rapid and accurate technology is relatively easy, once the technique  
348 has been optimised.

349

#### 350 **Acknowledgements**

351 We would like to thank all participants who agreed to take part of this study. We would also  
352 like to thank Lusha Kellgren, Maria McEwan, Narinder Boparai, Jacqui White, Neil Jones,  
353 Stephen Morris-Jones, Surjo De and Trupti Patel, for collection of samples at The  
354 Whittington Health NHS Trust, University College London Hospitals and the TB Service  
355 North Central London. Nirmala Ghimire and Devan Vaghela, for collection of samples at

356 Bart's Health NHS Trust. Josephine Silles, Pragna Patel, Marinos Christofi, Thomas Rendal,  
357 Yemi Martins and Esmeralda Sicat for the collection of samples at North Middlesex  
358 University Hospital. This work was funded by the UCLH/UCL NIHR Biomedical Resource  
359 Centre (Grant number: BRC/176/III/JB/101350) and the PATHSEEK European Union's  
360 Seventh Programme for research, technological development and demonstration (grant  
361 number 304875). The authors declare that they have no competing interests.

362

### 363 **References**

- 364 1. WHO. 2017. Global tuberculosis report 2017. WHO.
- 365 2. Scott LE, McCarthy K, Gous N, Nduna M, Rie AV, Sanne I, Venter WF, Duse A,  
366 Stevens W. 2011. Comparison of Xpert MTB/RIF with Other Nucleic Acid  
367 Technologies for Diagnosing Pulmonary Tuberculosis in a High HIV Prevalence  
368 Setting: A Prospective Study. *PLOS Med* 8:e1001061.
- 369 3. Jenkins C, Claxton AP, Shorten RJ, McHugh TD, Gillespie SH. 2005. Rifampicin  
370 Resistance in Tuberculosis Outbreak, London, England. *Emerg Infect Dis* 11:931.
- 371 4. Ocheretina O, Byrt E, Mabou M-M, Royal-Mardi G, Merveille Y-M, Rouzier V,  
372 Fitzgerald DW, Pape JW. 2016. False-positive rifampin resistant results with Xpert  
373 MTB/RIF version 4 assay in clinical samples with a low bacterial load. *Diagn Microbiol*  
374 *Infect Dis* 85:53–55.
- 375 5. Witney AA, Gould KA, Arnold A, Coleman D, Delgado R, Dhillon J, Pond MJ, Pope  
376 CF, Planche TD, Stoker NG, Cosgrove CA, Butcher PD, Harrison TS, Hinds J. 2015.



- 377 Clinical Application of Whole-Genome Sequencing To Inform Treatment for  
378 Multidrug-Resistant Tuberculosis Cases. *J Clin Microbiol* 53:1473–1483.
- 379 6. Pankhurst LJ, del Ojo Elias C, Votintseva AA, Walker TM, Cole K, Davies J, Fermont  
380 JM, Gascoyne-Binzi DM, Kohl TA, Kong C, Lemaitre N, Niemann S, Paul J, Rogers  
381 TR, Roycroft E, Smith EG, Supply P, Tang P, Wilcox MH, Wordsworth S, Wyllie D,  
382 Xu L, Crook DW. 2016. Rapid, comprehensive, and affordable mycobacterial diagnosis  
383 with whole-genome sequencing: a prospective study. *Lancet Respir Med* 4:49–58.
- 384 7. Walker TM, Cruz ALG, Peto TE, Smith EG, Esmail H, Crook DW. 2017. Tuberculosis  
385 is changing. *Lancet Infect Dis* 17:359–361.
- 386 8. Brown AC, Bryant JM, Einer-Jensen K, Holdstock J, Houniet DT, Chan JZM, Depledge  
387 DP, Nikolayevskyy V, Broda A, Stone MJ, Christiansen MT, Williams R, McAndrew  
388 MB, Tutill H, Brown J, Melzer M, Rosmarin C, McHugh TD, Shorten RJ, Drobniewski  
389 F, Speight G, Breuer J. 2015. Rapid Whole-Genome Sequencing of Mycobacterium  
390 tuberculosis Isolates Directly from Clinical Samples. *J Clin Microbiol* 53:2230–2237.
- 391 9. Cohen T, van Helden PD, Wilson D, Colijn C, McLaughlin MM, Abubakar I, Warren  
392 RM. 2012. Mixed-Strain Mycobacterium tuberculosis Infections and the Implications  
393 for Tuberculosis Treatment and Control. *Clin Microbiol Rev* 25:708–719.
- 394 10. van Rie A, Victor TC, Richardson M, Johnson R, van der Spuy GD, Murray EJ, Beyers  
395 N, Gey van Pittius NC, van Helden PD, Warren RM. 2005. Reinfection and mixed  
396 infection cause changing Mycobacterium tuberculosis drug-resistance patterns. *Am J*  
397 *Respir Crit Care Med* 172:636–642.
- 398 11. Baldeviano-Vidalón GC, Quispe-Torres N, Bonilla-Asalde C, Gastiaburú-Rodríguez D,  
399 Pro-Cuba JE, Llanos-Zavalaga F. 2005. Multiple infection with resistant and sensitive

- 400 M. tuberculosis strains during treatment of pulmonary tuberculosis patients. *Int J Tuberc*  
401 *Lung Dis Off J Int Union Tuberc Lung Dis* 9:1155–1160.
- 402 12. Zetola NM, Shin SS, Tumedí KA, Moeti K, Ncube R, Nicol M, Collman RG, Klausner  
403 JD, Modongo C. 2014. Mixed Mycobacterium tuberculosis complex infections and  
404 false-negative results for rifampin resistance by GeneXpert MTB/RIF are associated  
405 with poor clinical outcomes. *J Clin Microbiol* 52:2422–2429.
- 406 13. Martín A, Herranz M, Ruiz Serrano MJ, Bouza E, García de Viedma D. 2010. The  
407 clonal composition of Mycobacterium tuberculosis in clinical specimens could be  
408 modified by culture. *Tuberc Edinb Scotl* 90:201–207.
- 409 14. Hanekom M, Streicher EM, Berg DV de, Cox H, McDermid C, Bosman M, Pittius NCG  
410 van, Victor TC, Kidd M, Soolingen D van, Helden PD van, Warren RM. 2013.  
411 Population Structure of Mixed Mycobacterium tuberculosis Infection Is Strain Genotype  
412 and Culture Medium Dependent. *PLOS ONE* 8:e70178.
- 413 15. Votintseva AA, Pankhurst LJ, Anson LW, Morgan MR, Gascoyne-Binzi D, Walker  
414 TM, Quan TP, Wyllie DH, Del Ojo Elias C, Wilcox M, Walker AS, Peto TEA, Crook  
415 DW. 2015. Mycobacterial DNA extraction for whole-genome sequencing from early  
416 positive liquid (MGIT) cultures. *J Clin Microbiol* 53:1137–1143.
- 417 16. Honeyborne I, McHugh TD, Phillips PPJ, Bannoo S, Bateson A, Carroll N, Perrin FM,  
418 Ronacher K, Wright L, van Helden PD, Walzl G, Gillespie SH. 2011. Molecular  
419 bacterial load assay, a culture-free biomarker for rapid and accurate quantification of  
420 sputum Mycobacterium tuberculosis bacillary load during treatment. *J Clin Microbiol*  
421 49:3905–3911.

- 422 17. Yang H, Wang K. 2015. Genomic variant annotation and prioritization with  
423 ANNOVAR and wANNOVAR. *Nat Protoc* 10:1556–1566.
- 424 18. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-  
425 analysis of large phylogenies. *Bioinformatics* 30:1312–1313.
- 426 19. Williams DL, Waguespack C, Eisenach K, Crawford JT, Portaels F, Salfinger M, Nolan  
427 CM, Abe C, Sticht-Groh V, Gillis TP. 1994. Characterization of rifampin-resistance in  
428 pathogenic mycobacteria. *Antimicrob Agents Chemother* 38:2380–2386.
- 429 20. Sekiguchi J, Miyoshi-Akiyama T, Augustynowicz-Kopec E, Zwolska Z, Kirikae F,  
430 Toyota E, Kobayashi I, Morita K, Kudo K, Kato S, Kuratsuji T, Mori T, Kirikae T.  
431 2007. Detection of Multidrug Resistance in *Mycobacterium tuberculosis*. *J Clin*  
432 *Microbiol* 45:179–192.
- 433 21. Zenteno-Cuevas R, Zenteno JC, Cuellar A, Cuevas B, Sampieri CL, Riviera JE, Parissi  
434 A. 2009. Mutations in *rpoB* and *katG* genes in *Mycobacterium* isolates from the  
435 Southeast of Mexico. *Mem Inst Oswaldo Cruz* 104:468–472.
- 436 22. Wengenack NL, Uhl JR, St Amand AL, Tomlinson AJ, Benson LM, Naylor S, Kline  
437 BC, Cockerill FR, Rusnak F. 1997. Recombinant *Mycobacterium tuberculosis*  
438 *KatG*(S315T) is a competent catalase-peroxidase with reduced activity toward isoniazid.  
439 *J Infect Dis* 176:722–727.
- 440 23. Saint-Joanis B, Souchon H, Wilming M, Johnsson K, Alzari PM, Cole ST. 1999. Use of  
441 site-directed mutagenesis to probe the structure, function and isoniazid activation of the  
442 catalase/peroxidase, *KatG*, from *Mycobacterium tuberculosis*. *Biochem J* 338 ( Pt  
443 3):753–760.

- 444 24. Pym AS, Saint-Joanis B, Cole ST. 2002. Effect of katG mutations on the virulence of  
445 Mycobacterium tuberculosis and the implication for transmission in humans. Infect  
446 Immun 70:4955–4960.
- 447 25. Devonshire AS, O’Sullivan DM, Honeyborne I, Jones G, Karczmarczyk M, Pavšič J,  
448 Gutteridge A, Milavec M, Mendoza P, Schimmel H, Van Heuverswyn F, Gorton R,  
449 Cirillo DM, Borroni E, Harris K, Barnard M, Heydenrych A, Ndusilo N, Wallis CL,  
450 Pillay K, Barry T, Reddington K, Richter E, Mozioğlu E, Akyürek S, Yalçinkaya B,  
451 Akgoz M, Žel J, Foy CA, McHugh TD, Huggett JF. 2016. The use of digital PCR to  
452 improve the application of quantitative molecular diagnostic methods for tuberculosis.  
453 BMC Infect Dis 16:366.
- 454 26. Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K,  
455 Chatterjee A, Smith EG, Sanderson N, Walker TM, Morgan MR, Wyllie DH, Walker  
456 AS, Peto TEA, Crook DW, Iqbal Z. 2017. Same-day diagnostic and surveillance data  
457 for tuberculosis via whole genome sequencing of direct respiratory samples. J Clin  
458 Microbiol.
- 459 27. Depledge DP, Palser AL, Watson SJ, Lai IY-C, Gray ER, Grant P, Kanda RK, Leproust  
460 E, Kellam P, Breuer J. 2011. Specific capture and whole-genome sequencing of viruses  
461 from clinical samples. PloS One 6:e27805.
- 462 28. Christiansen MT, Brown AC, Kundu S, Tutill HJ, Williams R, Brown JR, Holdstock J,  
463 Holland MJ, Stevenson S, Dave J, Tong CYW, Einer-Jensen K, Depledge DP, Breuer J.  
464 2014. Whole-genome enrichment and sequencing of Chlamydia trachomatis directly  
465 from clinical samples. BMC Infect Dis 14:591.

- 466 29. Nimmo C, Doyle R, Burgess C, Williams R, Gorton R, McHugh TD, Brown M, Morris-  
467 Jones S, Booth H, Breuer J. 2017. Rapid identification of a *Mycobacterium tuberculosis*  
468 full genetic drug resistance profile through whole genome sequencing directly from  
469 sputum. *Int J Infect Dis* 62:44–46.
- 470 30. Witney AA, Bateson ALE, Jindani A, Phillips PPJ, Coleman D, Stoker NG, Butcher  
471 PD, McHugh TD, RIFAQUIN Study Team. 2017. Use of whole-genome sequencing to  
472 distinguish relapse from reinfection in a completed tuberculosis clinical trial. *BMC Med*  
473 15:71.
- 474 31. Cohen T, Chindelevitch L, Misra R, Kempner ME, Galea J, Moodley P, Wilson D.  
475 2016. Within-Host Heterogeneity of *Mycobacterium tuberculosis* Infection Is  
476 Associated With Poor Early Treatment Response: A Prospective Cohort Study. *J Infect*  
477 *Dis* 213:1796–1799.
- 478 32. Bifani PJ, Mathema B, Kurepina NE, Kreiswirth BN. 2002. Global dissemination of the  
479 *Mycobacterium tuberculosis* W-Beijing family strains. *Trends Microbiol* 10:45–52.
- 480 33. Eldholm V, Rieux A, Monteserin J, Lopez JM, Palmero D, Lopez B, Ritacco V, Didelot  
481 X, Balloux F. 2016. Impact of HIV co-infection on the evolution and transmission of  
482 multidrug-resistant tuberculosis. *eLife* 5:e16644.
- 483 34. Griffin JE, Gawronski JD, DeJesus MA, Ioerger TR, Akerley BJ, Sasseti CM. 2011.  
484 High-Resolution Phenotypic Profiling Defines Genes Essential for *Mycobacterial*  
485 Growth and Cholesterol Catabolism. *PLOS Pathog* 7:e1002251.

486

487 **Tables**

488 **Table 1. Target genes/regions and reason for inclusion in the reduced *M. tuberculosis***

489 **bait set.**

Gene target	Region property
<i>gyrB/gyrA</i>	Fluoroquinolone resistance
<i>rpoB/rpoC</i>	Rifampicin resistance
<i>rpsL</i>	Streptomycin resistance
<i>Rrs</i>	Streptomycin, amikacin and kanamycin resistance
<i>gidB</i>	Streptomycin resistance
<i>mabA/fabG1</i> + promoter + <i>inhA</i>	Isoniazid and ethionamide resistance
<i>katG</i>	Isoniazid resistance
<i>kasA</i>	Isoniazid resistance
<i>aphC-oxvR</i>	Isoniazid resistance
<i>tlyA</i>	Capreomycin resistance
Promoter + <i>pncA</i>	Pyrazinamide resistance
<i>eis</i> + promoter	Kanamycin resistance
<i>thyA</i>	PAS resistance
<i>embC</i>	Ethambutol resistance
<i>embB</i>	Ethambutol resistance
<i>ethA</i>	Ethionamide resistance
Direct repeat locus	Spoligotyping

490 Table 2. Antimicrobial resistance profiles from 13 patients with evidence of resistance from direct sputum sequencing using whole  
491 genome bait set.

Patient	RIFAMPICIN								ISONIAZID					ETHAMBUTOL				PYRAZINAMIDE			STREPTOMYCIN					PAS								
	gene	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype	phenotype						
BH001MC	+	S	ND	S	ND	S	ND	ND	S	S	ND	S	ND	S	ND	S	ND	S	S	S	S	S	ND	S	ND	ND	R	ND	ND					
BH041OS	+	S	S	S	S	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	ND				
BH052SA	+++	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND				
BH056ESDS	+	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	R	R	ND	
RF002AH	+	S	S	S	S	S	S	ND	S	R	R	R	S	S	S	R	S	S	S	S	S	S	S	S	S	S	S	S	ND	S	S	ND		
RF009ZC	Neg	S	S	S	S	S	S	S	S	R	R	S	S	S	R	S	S	S	S	S	S	S	R	S	ND	S	S	S	S	S	ND			
RF015GT	+	S	S	S	S	R	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	S	S	ND		
RF016SW	++	S	S	S	S	S	S	S	S	R	R	S	S	S	R	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	S	S	ND	
WH006AW	++	R	R	S	S	S	R	R	R	R	S	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	R	ND	S	S	ND
WH017KL	+	S	S	S	S	S	S	S	S	S	S	S	S	S	R	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	R	R	ND
WH026NS	Neg	ND	S	ND	S	ND	S	ND	S	ND	S	ND	S	ND	S	ND	S	ND	S	S	S	S	ND	S	ND	S	ND	ND	R	ND	R	ND	ND	
WH036ES	Neg	S	S	S	S	S	S	S	S	S	R	R	S	S	R	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	S	S	ND
WH037PD	++	S	S	S	R	S	R	R	R	R	R	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	S	ND	S	S	ND	

492 "ND" = Sample or result not available. "." = Low sequence read coverage at position. "R" = Resistant, "S" = Susceptible

493 **Figure legends**

494 **Figure 1.** Comparison of the percentage of reference genome with at least one sequence read  
495 covering a position on the y axis by the median depth of coverage on the x axis for each  
496 individual sample. This is stratified by whether original sample material was a sputum or  
497 MGIT culture.

498 **Figure 2.** (A) Barplot showing the percentage reference genome coverage (A single read  
499 covering each genome position) for patients with both a sputum and MGIT sample  
500 sequenced. The plot is annotated with both Xpert MTB/RIF (GX) and smear microscopy  
501 (Smear) results. For GX: \*\*\*\* = high, \*\*\* = medium, \*\* = low, and \* = very low. For smear  
502 \*\*\* = 3+, \*\* = 2+, \* = 1+, S = scanty and - = negative. Where a result is missing the test was  
503 not carried out. (B & C) Boxplot showing median depth of coverage for sputum samples  
504 stratified by both the quantitative Xpert MTB/RIF measure (B) and semi-quantitative smear  
505 result (C).

506 **Figure 3.** Heatmap clustering samples by the pairwise number of single nucleotide  
507 differences between them. Sample names are formatted such that the patient identifier is at  
508 the start followed by whether the sample originated from a MGIT culture or sputum.

509 **Figure 4.** The time taken in days on the x axis from sample collection (day 0) to when the  
510 MGIT samples flag positive and sequence result becomes available is denoted by the grey  
511 bars. The time taken for a patient sputum result to become available is marked by three  
512 different identifiers depending on whether the sample underwent whole genome sequencing  
513 or partial genome sequencing and whether the 24 hour or one hour hybridisation protocol was  
514 used. Sputum samples that failed sequencing are marked by missing symbols. Patient  
515 identifiers marked in red had confirmed drug resistant TB infections.









