1    Article type: Research

2

3    Title: Whole genome sequencing *Mycobacterium tuberculosis* directly from sputum

4    identifies more genetic diversity than sequencing from culture

5

6    Authors:

Camus Nimmo[1,2]*

Liam P. Shaw[3,4]

Ronan Doyle[1]

Rachel Williams[1]

Kayleen Brien[2]

Carrie Burgess[1]

Judith Breuer[1]

Francois Balloux[3]

Alexander S. Pym[2]

7

8       1.  Division of Infection and Immunity, University College London, London WC1E 6BT,

9            UK

10      2.  Africa Health Research Institute, Durban, South Africa

11      3.  UCL Genetics Institute, University College London, London WC1E 6BT, UK

12      4.  Nuffield Department of Clinical Medicine, Oxford University, Oxford OX3 7BN, UK

13      *Corresponding author

14

15   Corresponding author email: c.nimmo.04@cantab.net

16 ## Abstract

17

18 **Background**

19 Repeated culture reduces within-sample *Mycobacterium tuberculosis* genetic diversity due

20 to selection of clones suited to growth in culture and/or random loss of lineages, but it is

21 not known to what extent omitting the culture step altogether alters genetic diversity. We

22 compared *M. tuberculosis* whole genome sequences generated from 33 paired clinical

23 samples using two methods. In one method DNA was extracted directly from sputum then

24 enriched with custom-designed SureSelect (Agilent) oligonucleotide baits and in the other it

25 was extracted from mycobacterial growth indicator tube (MGIT) culture.

26

27 **Results**

28 DNA directly sequenced from sputum showed significantly more within-sample diversity

29 than that from MGIT culture (median 5.0 vs 4.5 heterozygous alleles per sample, p=0.04).

30 Resistance associated variants present as HAs occurred in four patients, and in two cases

31 may provide a genotypic explanation for phenotypic resistance.

32

33 **Conclusions**

34 Culture-free *M. tuberculosis* whole genome sequencing detects more within-sample

35 diversity than a leading culture-based method and may allow detection of mycobacteria

36 that are not actively replicating.

37

38 Key words: Mycobacterium tuberculosis; drug-resistant tuberculosis; whole genome

39 sequencing; sputum; within-patient diversity; heteroresistance

2

40    # Background

41

42    International efforts to reduce tuberculosis (TB) infections and mortality over the last two

43    decades have only been partially successful. In 2017, 10 million people developed TB and it

44    has overtaken HIV as the infectious disease responsible for the most deaths worldwide(1,

45    2). Drug resistance is a major concern with a steady rise in the number of reported cases

46    globally and rapid increases in some areas(1). Patients with *Mycobacterium tuberculosis*

47    resistant to the first line drugs rifampicin and isoniazid are classed as having multidrug-

48    resistant (MDR) TB and usually treated with a standardised second line drug regimen for at

49    least nine months, which is also used for rifampicin monoresistance(3, 4). With the

50    emergence of resistance to fluoroquinolones and aminoglycosides (extensively drug-

51    resistant [XDR] TB) there is an increasing need for individualised therapy based on drug

52    susceptibility testing (DST). Individualised therapy ensures patients are treated with

53    sufficient active drugs which can prevent selection of additional resistance, improve

54    treatment outcomes and reduce duration of infectiousness(5-8).

55

56    Traditionally, phenotypic culture-based DST was used to identify drug resistance but this is

57    being replaced by rapid genetic tests that detect specific drug resistance-conferring

58    mutations. Next generation whole genome sequencing (WGS) of *M. tuberculosis* is being

59    increasingly used in research and clinical settings to comprehensively identify all drug

60    resistance associated mutations(9). *M. tuberculosis* has a conserved genome with little

61    genetic diversity between strains and no evidence of horizontal gene transfer(10), but more

62    detailed analysis of individual patient samples with WGS has identified genetically separate

63    bacterial subpopulations in sequential sputum samples(11-16) and across different

3

64     anatomical sites(17). This within-patient diversity can occur as a result of mixed infection

65     with genetically distinct strains or within-host evolution of a single infecting strain(18).

66

67     Bacterial subpopulations can be detected in clinical samples after sequencing reads are

68     mapped to a reference genome where multiple base calls are detected at a single genomic

69     site. These heterozygous alleles (HAs) at sites associated with drug resistance (resistance

70     associated variants, RAVs) may reflect heteroresistance, where a fraction of the total

71     bacterial population is drug susceptible while the remainder is resistant(19). Identification

72     of genetic diversity within clinical samples may improve detection of RAVs over currently

73     available rapid genetic tests(19) and can be achieved with freely available WGS analysis

74     toolkits(20-22). Identifying RAVs could improve individualised therapy, prevent acquired

75     resistance(12), and give insight into bacterial adaptation to the host.

76

77     *M. tuberculosis* WGS is usually performed on fresh or stored frozen cultured isolates to

78     obtain sufficient purified mycobacterial DNA(23, 24). However, the culture process can

79     change the population structure from that of the original sample due to genetic drift

80     (random loss of lineages) and/or the selection of subpopulations more suited to growth in

81     culture(25-27). Repeated subculture leads to loss of genetic diversity and

82     heteroresistance(28). Additionally, in the normal course of *M. tuberculosis* infection, some

83     bacteria exist as viable non-culturable persister organisms that are hypothesised to cause

84     the high relapse rate seen following treatment of insufficient duration. Although these

85     organisms may be identified in sputum by techniques such as reporter phages or culture

86     with resuscitation promoting factors(29, 30) they are likely to be missed by any sequencing

87     method reliant on standard culture.

4

88

89    WGS directly from sputum without enrichment is challenging(23). It has recently been

90    improved by depleting human DNA during DNA extraction(31). We have previously reported

91    the use of oligonucleotide enrichment technology SureSelect (Agilent, CA, USA) to sequence

92    *M. tuberculosis* DNA directly from sputum(32) and demonstrated its utility in determining a

93    rapid genetic drug resistance profile(33, 34).

94

95    It remains unclear to what extent WGS of cultured *M. tuberculosis* samples underestimates

96    the genetic diversity of the population in sputum samples. One previous study of 16 patients

97    did not identify increased genetic diversity in *M. tuberculosis* DNA sequenced directly from

98    sputum compared to DNA from culture(31), whereas another study of mostly drug

99    susceptible patients showed sequencing directly from sputum identified a slight excess of

100   HAs relative to culture(33). Here we reanalyse heterozygous alleles (HAs) for the 12

101   available paired sequences with >60-fold mean genome coverage from that study(33) in

102   addition to 21 newly collected samples from patients with MDR-TB and further explore the

103   genomic location of the additional diversity identified.

104

105   Results

106

107   **Patient Characteristics and Drug Susceptibility Testing**

108

109   Whole genome sequences were obtained for 33 patients from both mycobacterial growth

110   indicator tube (MGIT) culture and direct sputum sequencing. The patients were

111   predominantly of black African ethnicity (83%) and 50% were HIV positive. First line

5

112     phenotypic drug susceptibility testing (DST) results identified 20 patients with MDR-TB and

113     one with rifampicin monoresistance. In addition there were two isoniazid monoresistant

114     patients and ethambutol resistance was detected in 7 patients. Second-line phenotypic DST

115     was performed for patients with rifampicin-resistant or MDR-TB and identified one case of

116     kanamycin resistance (Table 1).

117

118     All samples had mean genome coverage of 60x or above with at least 85% of the genome

119     covered at 20x (Supplementary Material: Table 1). We observed greater mean coverage

120     depth in sputum-derived sequences than MGIT sequences (median 173.7 vs 142.4, p=0.03,

121     Supplementary Material: Table 1), and so mapped reads were randomly downsampled to

122     give equal mean coverage depth in each pair. A genotypic susceptibility profile was

123     determined by evaluating MGIT WGS for consensus-level RAVs using a modified version of

124     publicly available lists(22, 35). Genotypic RAVs predicted all rifampicin phenotypic resistance

125     and >95% of isoniazid phenotypic resistance. Ethambutol genotypic RAVs were poorly

126     predictive of phenotypic resistance in line with findings from other studies(36) (Table 1).

127     The patient with kanamycin phenotypic resistance was correctly identified by an *rrs* a1401g

128     RAV. No full phenotypic fluoroquinolone phenotypic resistance was identified, but several

129     colonies from patient F1013 did grow in the presence of ofloxacin (although not enough to

130     be classified as resistant). The consensus sequences from this patient harboured a *gyrB*

131     E501D mutation which is believed to confer resistance to moxifloxacin but not other

132     fluoroquinolones, which may explain the borderline phenotypic DST result(37).

133

134     **Genetic Diversity**

135

6

136    To compare consensus sequences from sputum and MGIT, a WGS consensus sequence-level

137    maximum likelihood phylogenetic tree was constructed (Supplementary Material: Figure 1).

138    As expected, all paired sequences were closely related, with a median difference of 0.0

139    (range 0-1) single nucleotide polymorphisms (SNPs). Samples from patients F1066 and

140    F1067 were closely related with only one consensus-level SNP separating all four consensus

141    sequences. There was no obvious epidemiological link between these patients (although

142    this study was not designed to collect comprehensive epidemiological information) and they

143    lived 20km apart in Durban. However, both patients were admitted contemporaneously to

144    an MDR treatment facility and sampled on the same day. DNA extraction and sequencing

145    occurred on different runs. Therefore the close genetic linkage may represent direct

146    transmission within a hospital setting, a community transmission chain or an unlikely cross-

147    contamination during sample collection.

148

149    Having established congruence between sputum and MGIT sequences at the consensus

150    level we then compared genetic diversity by DNA source. We first defined a threshold for

151    calling variants present as heterozygous alleles (HAs) in our entire dataset by using a range

152    of minimum read count frequencies as described in the methods (Figure 1). Below a

153    minimum of three supporting reads there was an exponential increase in the number of HAs

154    identified, which may be indicative of the inclusion of sequencing errors. To reduce this risk,

155    we used a threshold of a minimum of four supporting reads.

156

157    Genetic diversity may occur because of within-host evolution or mixed infection. To identify

158    mixed infection we used a SNP-based barcode(38) to scan all HAs for a panel of 413 robust

159    phylogenetic SNPs that can resolve *M. tuberculosis* into one of seven lineages and 55 sub-

160     lineages. We found three phylogenetic SNPs among the HAs. In all cases the heterozygous

161     phylogenetic SNP originated from the same sublineage as other SNPs present at 100%

162     frequency, and there were no cases of HAs indicating the presence of more than one lineage

163     or sublineage. We screened for mixed infection with the same sublineage by screening

164     samples by HA frequency and then using Bayseian model based clustering in samples with

165     ≥10 HAs as described previously(39). This identified mixed infection in the sputum sample

166     from patient F1096, which had 261 heterozygous alleles, greater than ten times that in any

167     other sample. This patient was therefore excluded from further analyses.

168

169     As a first step to comparing diversity between sputum and MGIT sequenced samples we

170     looked at the location of genetic diversity within the *M. tuberculosis* genome. HAs were

171     widely dispersed across the genome at similar sites in both sputum and MGIT samples. The

172     genes with the greatest density of HAs are shown in Table 2.

173

174     Notably, genetic diversity was found in the ribosomal RNA (rRNA) genes (*rrs* and *rrl*)

175     uniquely in sputum samples, compared to other genes where distribution of diversity

176     between MGIT and sputum was more balanced. As rRNA contains regions that are highly

177     conserved across bacteria(40), we considered the possibility that SureSelect baits targeting

178     rRNA genes were capturing both *M. tuberculosis* and other bacterial species. To evaluate

179     this, metagenomic taxonomic assignment was performed on all reads by sampling reads

180     that were not assigned to *M. tuberculosis* (i.e. presumed contaminants from other bacteria).

181     We then performed a BLAST search against the most diverse genes listed in Table 2 which

182     indicated that a sizeable proportion of non-*M. tuberculosis* reads from directly sequenced

183     sputum had a BLAST hit of at least 30 bases to *M. tuberculosis* *rrs* and *rrl* genes that encode

184 rRNA (330 BLAST hits from sputum sequences vs 4 BLAST hits from MGIT sequences, median

185 8.5% vs 0.0%, p<0.01, Supplementary Material: Figure 2). There were no BLAST hits against

186 any of the other genes with ≥2 sputum HAs apart from *rpoC*, for which there were 3 BLAST

187 hits from sputum sequences but none from MGIT sequences (median 0.0% for both sputum

188 and MGIT sequences), indicating that this issue appears largely specific to rRNA. To

189 determine if contaminating reads were contributing to HAs identified in intergenic regions,

190 we repeated this analysis for all intergenic regions with ≥2 sputum HAs (Supplementary

191 Material: Table 2). There were no BLAST hits to any of these regions, suggesting that this is

192 not the case. The taxonomic assignment of these contaminating reads were typical of

193 genera composing the oral flora, with a high representation of *Actinomyces, Fusobacterium,*

194 *Prevotella,* and *Streptococcus* (Supplementary Material: Figure 3).

195

196 This supported the hypothesis that the baits may enrich rRNA from other organisms so rRNA

197 genes were excluded from further analysis. The difference in diversity between sputum and

198 MGIT sequences can be explained by the selective nature of MGIT media which will enrich

199 *M. tuberculosis* sequences and the decontamination step used to kill non-mycobacteria

200 prior to culture inoculation. Importantly the frequency of HAs in other highly diverse genes

201 between sequencing strategies was more balanced (Table 2) in addition to the lack of BLAST

202 hits of contaminating reads to these genes.

203

204 After excluding the sample with mixed infection and removing rRNA gene sequences we

205 compared the frequency of HAs in sputum and MGIT. There were 265 HAs identified across

206 all sputum samples compared to 200 in MGIT samples (median 5.0 vs 4.5, p=0.04,

207 Supplementary Material: Table 1). In both sputum and MGIT samples, the majority of HAs

208     were indels, and non-synonymous mutations were more commonly frameshift than

209     missense mutations (Table 3). The distribution of HAs by patient is shown in Figure 2.

210

211     **Genetic diversity in drug resistance genes**

212

213     HAs in drug resistance associated regions, including promoters and intergenic regions, were

214     individually assessed. Four of the 32 patients with single strain infection had RAVs present

215     as HAs in at least one gene, which are shown in Table 4. Patient F1002 had three

216     compensatory mutations in *rpoC* present at HAs in both sequences. As described above, the

217     strains from patients F1066 and F1067 were highly related with only one consensus SNP

218     difference between all four sequences. Both had phenotypic high level isoniazid resistance

219     with no consensus-level *katG* or *inhA* mutation, but had frameshift *katG* mutations present

220     as HAs which have the potential to cause resistance(41). F1066 and RF021 had *Rv1979c* and

221     *pncA* mutations respectively at low frequency in sputum only which have the potential to

222     confer phenotypic resistance to clofazimine (*Rv1979c*) and pyrazinamide (*pncA*), although

223     no phenotypic testing was performed for these drugs.

224

225     Discussion

226

227     In this study we performed whole genome sequencing using DNA from sputum and MGIT

228     culture in paired samples from 33 patients and compared within-patient genetic diversity

229     between methods. All paired sequences were closely related at the consensus level, and

230     WGS predicted phenotypic drug susceptibility with over 95% sensitivity and specificity for

231     rifampicin and isoniazid in line with published data(42).

10

232

233    We find that the rRNA genes have high levels of diversity in sputum samples, but believe

234    this is due to non-mycobacterial DNA hybridising to the capture baits. This conclusion is

235    borne out by the taxonomic assignment of reads aligning to these genes in common oral

236    bacteria. We therefore excluded these from further analysis, and recommend others using

237    enrichment from sputum do similarly. We find more diversity when sequencing directly

238    from sputum with significantly more unique heterozygous alleles (HAs) than sequencing

239    from MGIT culture (p=0.04).

240

241    The understanding of within-patient *M. tuberculosis* genetic diversity is becoming

242    increasingly important as the detection of rare variants has been shown to improve the

243    correlation between phenotypic and genotypic drug resistance profiles(19) and can identify

244    emerging drug resistance(11, 12). Not including a culture step avoids the introduction of

245    bias towards culture-adapted subpopulations and the impact of random chance and is also

246    likely to incorporate DNA from viable non-culturable mycobacteria. A reduction in genetic

247    diversity has previously been shown with sequential *M. tuberculosis* subculture(25, 28), but

248    was not confirmed by a study performing WGS directly from sputum(31). However, the 16

249    paired sputum and MGIT samples compared by Votintseva(31) had a minimum of 5x

250    coverage compared to a minimum 60x coverage in this study, and were likely to contain less

251    genetic material as they were surplus clinical rather than dedicated research samples.

252

253    Two-thirds of the patients with MDR-TB had already been treated for drug susceptible-TB

254    (DS-TB), and additional diversity in sputum samples may represent early adaptation to drug

255    pressure. As direct sputum sequencing does not rely on live mycobacteria, DNA from

256     recently killed *M. tuberculosis* is likely to also be sequenced, meaning that recent genomic

257     mutations are likely to be represented as HAs.

258

259     In two patients, RAVs present as HAs provided a likely genotypic basis for otherwise

260     unexplained phenotypic resistance. Given the small total number of resistance mutations in

261     this study, it is not possible to draw conclusions about the frequency of heterozygous RAVs

262     in directly sequenced sputum. However the presence of heterozygous RAVs in both MGIT

263     and sputum sequences reinforces the biological importance of these mutations.

264

265     To reduce the risk of sample cross contamination, paired samples were extracted on

266     different days, prepared in different sequencing libraries and sequenced on different runs.

267     However it is not possible to completely exclude the possibility of contamination during

268     sample collection and between different samples processed in batches. A further limitation

269     of this study is that it can be difficult to distinguish low frequency variants from sequencing

270     error. The SureSelect library preparation protocol for sputum sequencing incorporates more

271     PCR cycles than that used for MGIT sequencing, which may increase the risk of error. Where

272     possible this could be evaluated further by performing technical sequencing replicates on

273     extracted DNA samples, although this was not possible due to insufficient surplus material

274     and financial constraints. To reduce the risk of sequencing errors we used high read and

275     mapping quality thresholds, and required a stringent 98% identity between sequenced

276     reads and the reference genome. Low frequency variants of particular clinical importance

277     could be confirmed by resequencing the same DNA samples.

278

279　## Conclusions

280

281　Directly sequencing *M. tuberculosis* from sputum is able to identify more genetic diversity

282　than sequencing from culture. Understanding within-patient genetic diversity is important

283　to understand bacterial adaptation to drug treatment and the acquisition of drug resistance.

284　It also has potential to identify low frequency RAVs that may further enhance the prediction

285　of drug resistance phenotype from genotype.

286

287　## Methods

288

289　**Patient enrolment**

290　Adult patients presenting with a new diagnosis of sputum culture positive TB were included

291　in the study. Patients were recruited in London, UK (n=12) and Durban, South Africa (n=21).

292　All patients recruited in Durban were Xpert MTB/RIF (Cepheid, CA, USA) positive for

293　rifampicin resistance. Two sputum samples were collected prior to starting the current

294　treatment regimen, with one inoculated into mycobacterial growth indicator tube (MGIT)

295　culture (BD, NJ, USA) and the other used for direct DNA extraction. Therefore for patients

296　with drug susceptible-TB (DS-TB), sputum was collected prior to taking any TB therapy,

297　while patients starting MDR-TB treatment may have already taken treatment for DS-TB if

298　this was intiated prior to resistance results being available.

299

300　**Ethics, Consent and Permissions**

301　All patients gave written informed consent to participate in the study. Ethical approval for

302　the London study was granted by NHS National Research Ethics Service East Midlands–

13

303　　Nottingham 1 (reference 15/EM/0091). Ethical approval for the Durban study was granted

304　　by University of KwaZulu-Natal Biomedical Research Ethics Committee (reference

305　　BE022/13).

306

307　　**Microbiology**

308　　MGIT samples were incubated in a BACTEC MGIT 960 (BD, NJ, USA) until flagging positive.

309　　Phenotypic DST data for London samples were those provided to treating hospitals by Public

310　　Health England. Phenotypic DST were performed using equivalent standardised methods.

311　　For Durban samples this was the solid agar proportion method (Supplementary Material:

312　　Methods) and for London samples the resistance ratio method(43).

313

314　　**DNA extraction and sequencing**

315　　Positive MGIT tubes were centrifuged at 16,000g for 15 minutes and the supernatant

316　　removed. Cells were resuspended in phosphate-buffered saline before undergoing heat

317　　killing at 95°C for 1 hour followed by centrifugation at 16,000g for 15 minutes. The

318　　supernatant was removed and the sample resuspended in 1mL sterile saline (0.9% w/v). The

319　　wash step was repeated. DNA was extracted with mechanical ribolysis before purification

320　　with DiaSorin Liaison Ixt (DiaSorin, Italy) or CTAB(44). NEBNext Ultra II DNA (New England

321　　Biolabs, MA, USA) was used for DNA library preparation.

322

323　　Sputum samples for direct sequencing were heat killed, centrifuged at 16,000g for 15

324　　minutes and the supernatant was removed. DNA extraction was performed with mechanical

325　　ribolysis followed by purification using DiaSorin Liaison Ixt (DiaSorin, Italy) or DNeasy blood

326　　& tissue kit (Qiagen, Germany)(44). Target enrichment was performed using SureSelect with

14

327    a custom-designed bait set covering the entire positive strand of the *M. tuberculosis*

328    genome as described previously(33). Batches of 48 multiplexed samples were sequenced on

329    NextSeq 500 (Illumina, CA, USA) 300-cycle paired end runs with a mid-output kit.

330    Sequencing was performed by the Pathogen Genomics Unit at University College London in

331    a dedicated laboratory where one sequencing run was processed at a time. All paired

332    samples were extracted, prepared and sequenced on different days. The National Center for

333    Biotechnology Information Sequence Read Archive (NCBI SRA) accession number for each

334    sample is shown in Supplementary Material: Table 3.

335

336    **Read mapping**

337    DNA sequence reads were adapter and quality trimmed then aligned to the H37Rv

338    reference genome (GenBank accession NC_000962.3) with Trim Galore v0.4.4(45) and

339    BBMap v38.32(46), with mapped reads stored in an output bam file. Duplicate reads were

340    removed with Picard tools v1.130(47) MarkDuplicates and coverage statistics generated

341    with Qualimap v2.2.1(48). For each sample pair, the bam file with greater mean genome

342    coverage was randomly downsampled to that of the paired sample with Picard tools

343    v1.130(47) DownsampleSam. All further analyses were performed using these

344    downsampled bam files. Command line parameters used are specified in the Supplementary

345    Material: Methods.

346

347    **Variant calling**

348    Variant calling for comparison for HA counts was performed with FreeBayes v1.2(49).

349    Variants falling in or within 50 bases of PE/PPE family genes and repeat elements were

350    excluded using vcfinteresect in vcflib(50). For the initial analysis of genetic diversity, variants

351     were included if supported by ≥2 reads, with ≥1 forward and reverse read, no read position

352     bias, a minimum mapping quality of 30 and base quality of 30. The minimum supporting

353     read threshold was then increased in a stepwise fashion from 2 to 15. Variant calling files

354     where variants were supported ≥4 supporting reads including ≥1 forward and reverse read

355     were used to compare HA frequency and location and to screen for mixed infection.

356

357     The phylogenetic tree was constructed by calling variants with VarScan v2.4.0(51)

358     mpileup2cns as this is able to generate consensus-level calls at each reference sequence

359     base. SNPs were then used to generate a sequence of equal length to the reference using a

360     custom perl script and these sequences were combined in a multi-alignment fasta file. SNP

361     sites were extracted from this alignment using snp-sites v2.4.1(52), and pairwise SNP

362     differences calculated using snp-dists v0.6.3(53). Extracted SNP sites were used to generate

363     a maximum likelihood phylogenetic tree using RaxML v8.2.12(54) which was visualised using

364     FigTree v1.4.3.

365

366     **Identification of Mixed Infection**

367     All samples were screened for evidence of mixed infection using described methods(39). In

368     brief, any sample with 10 or fewer heterozygous SNPs, or between 11 and 20 heterozygous

369     SNPs where heterozygous SNPs were ≤1.5% of all SNPs was classified as not mixed. For

370     other samples, the Baysian mixture model analysis(39) was used where samples with a

371     Bayesian information criterion value >20 for presence of more than one strain were

372     assumed to be mixed.

373

374     **Metagenomic assignment**

16

375     Sequencing reads were classified using Kraken v0.10.6(55) against a custom Kraken

376     database previously constructed from all available RefSeq genomes for bacteria, archaea,

377     viruses, protozoa, and fungi, as well as all RefSeq plasmids (as of September 19[th] 2017) and

378     three human genome reference sequences(56). The size of the final database after shrinking

379     was 193 Gb, covering 38,190 distinct NCBI taxonomic IDs.

380

381     To assess the proportion of contaminating reads that could generate spurious diversity

382     when mapped to *M. tuberculosis* ribosomal genes, we randomly subsampled 100 reads

383     taxonomically assigned as non-*M. tuberculosis* and performed a BLAST search with blastn

384     v2.2.28(57) against each gene as described from the H37Rv reference genome. We only

385     analysed hits of at least 30 bases.

386

387     **Statistics**

388     Statistical analyses were performed with Prism v8.0 (Graphpad, CA, USA). Mean coverage

389     depth statistics, number of HAs and BLAST hits of contaminating reads in paired samples

390     were compared using a two-tailed Wilcoxon matched-pairs signed rank test.

391

392     Abbreviations

393

| DST | drug susceptibility testing |
|---|---|
| DS-TB | drug susceptible-tuberculosis |
| HA | heterozygous allele |
| MDR-TB | multidrug resistant-tuberculosis |

17

| MGIT | mycobacterial growth indicator tube |
|------|-------------------------------------|
| RAV | Resistance associated variant |
| rRNA | ribosomal RNA |
| SNP | single nucleotide polymorphism |
| TB | tuberculosis |
| WGS | whole genome sequencing |

394

395

396

397    # Declarations

398

399    **Ethics approval and consent to participate**

400    All patients gave written informed consent to participate in the study. Ethical approval for

401    the London study was granted by NHS National Research Ethics Service East Midlands–

402    Nottingham 1 (reference 15/EM/0091). Ethical approval for the Durban study was granted

403    by University of KwaZulu-Natal Biomedical Research Ethics Committee (reference

404    BE022/13).

405

406    **Consent for publication**

407    Not applicable

408

409    **Availability of data and materials**

410    Original fastq files are available at NCBI Sequence Read Archive with BioProject reference

411    PRJNA486713: https://www.ncbi.nlm.nih.gov/bioproject/PRJNA486713/

412

413    **Competing interests**

414    The authors declare that they have no competing interests.

415

416    **Funding**

417    Camus Nimmo is funded by a Wellcome Trust Research Training Fellowship reference

418    203583/Z/16/Z. This work was additionally funded by National Institute for Health Research

419    via the UCLH/UCL Biomedical Research Centre (grant number BRC/176/III/JB/101350) and

420    the PATHSEEK European Union's Seventh Programme for research and technological

423

424     **Authors' contributions**

425     Study conception: JB, ASP

426     Data collection: CB, KB

427     Analysis and interpretation: CN, LPS, RD, RW

428     Drafting of manuscript: CN, LPS

429     Revision of manuscript: FB, JB, ASP

430     Final approval of manuscript: CN, LPS, RD, RW, KB, CB, JB, FB, ASP

431

432     **Acknowledgements**

435    Tables

436

| Drug | Resistance by phenotypic DST | Resistance by genotypic DST | Genotypic DST sensitivity | Genotypic DST specificity |
|---|---|---|---|---|
| *First line drugs* | | | | |
| Rifampicin | 21/32 (65.6%) | 21/33 (63.6%) | 21/21 (100%)* | 21/21 (100%) |
| Isoniazid | 22/32 (68.8%) | 24/36 (66.7%) | 21/22 (95.5%) | 23/24 (95.8%) |
| Ethambutol | 7/31 (22.6%) | 15/34 (44.1%) | 7/7 (100%) | 7/15 (46.7%) |
| *Second line drugs* | | | | |
| Ofloxacin | 0/22 (0.0%) | 1/22 (4.5%) | N/A | 0/1 (0%)** |
| Kanamycin | 1/22 (4.5%) | 1/22 (4.5%) | 1/1 (100%) | 1/1 (100%) |

437

438    Table 1. Phenotypic and genotypic drug susceptibility testing (DST) results and sensitivity

439    and specificity of genotypic DST relative to phenotypic DST. Phenotypic DST available for

440    first line drugs for 32 of the 33 patients, and for second line drugs for 22 patients who

441    demonstrated rifampicin drug resistance. *In one directly-sequenced sputum samples

442    rifampicin RAVs were missed due to low coverage, although they were identified in the

443    corresponding MGIT sample. **This sample had <1% of colonies grow in the presence of

444    ofloxacin, so is categorised as sensitive but may have low-level or heteroresistance to

445    fluoroquinolones (see main text).

446

447

| Gene | Heterozygous alleles per base | | Total number of heterozygous alleles | | Functional category |
|---|---|---|---|---|---|
| | Sputum | MGIT | Sputum | MGIT | |
| rv1319c | 0.021 | 0.021 | 33 | 33 | Metabolism and respiration |
| rrs | 0.016 | 0.000 | 25 | 0 | 16S ribosomal RNA |
| rrl | 0.006 | 0.000 | 19 | 0 | 23S ribosomal RNA |
| ppsA | 0.003 | 0.001 | 15 | 4 | Lipid metabolism |
| rv2082 | 0.006 | 0.006 | 13 | 14 | Unknown function |
| accE5 | 0.006 | 0.000 | 3 | 0 | Lipid metabolism |
| lppB | 0.005 | 0.005 | 3 | 3 | Probable surface lipoprotein |
| pks12 | 0.000 | 0.001 | 3 | 10 | Lipid metabolism |
| rv2319c | 0.003 | 0.005 | 3 | 4 | Stress protein |
| lppA | 0.003 | 0.002 | 2 | 1 | Probable surface lipoprotein |
| rpoC | 0.001 | 0.001 | 2 | 3 | RNA polymerase beta' subunit |
| rv3888c | 0.002 | 0.001 | 2 | 1 | Probable membrane protein |
| vapC25 | 0.005 | 0.000 | 2 | 0 | Possible toxin |
| vapC31 | 0.005 | 0.002 | 2 | 1 | Possible toxin |

448

449     Table 2. Genes with ≥2 heterozygous alleles (HAs) across all sputum samples, ordered by

450     greatest number of HAs per base.

451

| | Sputum variants | MGIT variants |
|---|---|---|
| Total variants | 24480 | 25465 |
| Total variants present as HAs (% of total variants) | 265 (1.1%) | 200 (0.8%) |
| Median HAs per sample | 5.0 | 4.5 |
| Variant type (% all HAs) | | |
| SNP | 217 (81.9%) | 174 (87.0%) |
| MNP | 2 (0.8%) | 0 (0.0%) |
| Insertion | 4 (1.5%) | 1 (0.5%) |
| Deletion | 24 (9.1%) | 15 (7.5%) |
| Complex | 18 (6.8%) | 10 (5.0%) |
| Coding change (% all HAs) | | |
| Non-synonymous (missense) | 93 (35.1%) | 77 (38.5%) |
| Non-synonymous (frameshift) | 6 (2.3%) | 7 (3.5%) |
| Synonymous | 57 (21.5%) | 57 (28.5%) |
| Intergenic | 109 (41.1%) | 59 (29.5%) |

452

453    Table 3. Variants identified in MGIT and sputum derived sequences from paired samples.

454    Values given represent totals for 32 paired samples. SNP = single nucleotide polymorphism;

455    MNP = multi-nucleotide polymorphism.

456

24

457

| Patient ID | Phenotypic resistance | Mutation | Frequency (MGIT/sputum) | Description |
|---|---|---|---|---|
| F1002 | Rifampicin | *rpoB* S450L | 100%/100% | High confidence resistance mutation |
| F1002 | Rifampicin | *rpoC* G332R(58) | 82.6%/21.7% | Putative compensatory mutations |
| F1002 | Rifampicin | *rpoC* L516P(58) | 12.7%/7.7% | |
| F1002 | Rifampicin | *rpoC* P1040S(59) | 21.7%/12.3% | |
| F1066 | Isoniazid (high) | *katG* N218fs | 0.0%/6.9% | Possible resistance mutations, not previously described |
| F1066 | Clofazimine – not tested | *Rv1979c* G376D | 0.0%/0.5% | |
| F1067 | Isoniazid (high) | *katG* N218fs | 10.7%/7.6% | |
| RF021 | Pyrazinamide – testing failed | *pncA* Q122H | 0%/2.5% | |

458

459    Table 4. Resistance associated variants present as heterozygous alleles (HAs).

460

461    References

462

463    1.    Global Tuberculosis Report 2018. Geneva: World Health Organization; 2018 2018.
464    2.    Murray CJ, Ortblad KF, Guinovart C, Lim SS, Wolock TM, Roberts DA, et al. Global,
465    regional, and national incidence and mortality for HIV, tuberculosis, and malaria during
466    1990-2013: a systematic analysis for the Global Burden of Disease Study 2013. Lancet.
467    2014;384(9947):1005-70.
468    3.    Dheda K, Gumbo T, Maartens G, Dooley KE, McNerney R, Murray M, et al. The
469    epidemiology, pathogenesis, transmission, diagnosis, and management of multidrug-
470    resistant, extensively drug-resistant, and incurable tuberculosis. Lancet Respir Med. 2017.
471    4.    WHO treatment guidelines for drug-resistant tuberculosis. World Health
472    Organization; 2016.
473    5.    Trauner A, Liu Q, Via LE, Liu X, Ruan X, Liang L, et al. The within-host population
474    dynamics of Mycobacterium tuberculosis vary with treatment efficacy. Genome Biol.
475    2017;18(1):71.
476    6.    Olaru ID, Lange C, Heyckendorf J. Personalized medicine for patients with MDR-TB. J
477    Antimicrob Chemother. 2016;71(4):852-5.
478    7.    Pasipanodya JG, McIlleron H, Burger A, Wash PA, Smith P, Gumbo T. Serum drug
479    concentrations predictive of pulmonary tuberculosis outcomes. J Infect Dis.
480    2013;208(9):1464-73.
481    8.    Cegielski JP, Kurbatova E, van der Walt M, Brand J, Ershova J, Tupasi T, et al.
482    Multidrug-Resistant Tuberculosis Treatment Outcomes in Relation to Treatment and Initial
483    Versus Acquired Second-Line Drug Resistance. Clin Infect Dis. 2016;62(4):418-30.
484    9.    Satta G, Lipman M, Smith GP, Arnold C, Kon OM, McHugh TD. Mycobacterium
485    tuberculosis and whole-genome sequencing: how close are we to unleashing its full
486    potential? Clin Microbiol Infect. 2017.
487    10.    Sreevatsan S, Pan X, Stockbauer KE, Connell ND, Kreiswirth BN, Whittam TS, et al.
488    Restricted structural gene polymorphism in the Mycobacterium tuberculosis complex
489    indicates evolutionarily recent global dissemination. Proceedings of the National Academy
490    of Sciences of the United States of America. 1997;94(18):9869-74.
491    11.    Sun G, Luo T, Yang C, Dong X, Li J, Zhu Y, et al. Dynamic population changes in
492    Mycobacterium tuberculosis during acquisition and fixation of drug resistance in patients. J
493    Infect Dis. 2012;206(11):1724-33.
494    12.    Merker M, Kohl TA, Roetzer A, Truebe L, Richter E, Rüsch-Gerdes S, et al. Whole
495    genome sequencing reveals complex evolution patterns of multidrug-resistant
496    Mycobacterium tuberculosis Beijing strains in patients. PLoS One. 2013;8(12):e82551.
497    13.    Operario DJ, Koeppel AF, Turner SD, Bao Y, Pholwat S, Banu S, et al. Prevalence and
498    extent of heteroresistance by next generation sequencing of multidrug-resistant
499    tuberculosis. PLoS One. 2017;12(5):e0176522.
500    14.    Black PA, de Vos M, Louw GE, van der Merwe RG, Dippenaar A, Streicher EM, et al.
501    Whole genome sequencing reveals genomic heterogeneity and antibiotic purification in
502    Mycobacterium tuberculosis isolates. BMC Genomics. 2015;16(1):857.
503    15.    Eldholm V, Norheim G, von der Lippe B, Kinander W, Dahle UR, Caugant DA, et al.
504    Evolution of extensively drug-resistant Mycobacterium tuberculosis from a susceptible
505    ancestor in a single patient. Genome Biol. 2014;15(11):490.

506    16.    Bloemberg GV, Keller PM, Stucki D, Stuckia D, Trauner A, Borrell S, et al. Acquired
507    Resistance to Bedaquiline and Delamanid in Therapy for Tuberculosis. N Engl J Med.
508    2015;373(20):1986-8.
509    17.    Lieberman TD, Wilson D, Misra R, Xiong LL, Moodley P, Cohen T, et al. Genomic
510    diversity in autopsy samples reveals within-host dissemination of HIV-associated
511    Mycobacterium tuberculosis. Nat Med. 2016;22(12):1470-4.
512    18.    Ford C, Yusim K, Ioerger T, Feng S, Chase M, Greene M, et al. Mycobacterium
513    tuberculosis--heterogeneity revealed through whole genome sequencing. Tuberculosis
514    (Edinb). 2012;92(3):194-201.
515    19.    Metcalfe JZ, Streicher E, Theron G, Colman RE, Allender C, Lemmer D, et al. Cryptic
516    Micro-heteroresistance Explains M. tuberculosis Phenotypic Resistance. Am J Respir Crit
517    Care Med. 2017.
518    20.    Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, et al. Rapid antibiotic-
519    resistance predictions from genome sequence data for Staphylococcus aureus and
520    Mycobacterium tuberculosis. Nat Commun. 2015;6:10063.
521    21.    Feuerriegel S, Schleusener V, Beckert P, Kohl TA, Miotto P, Cirillo DM, et al.
522    PhyResSE: a Web Tool Delineating Mycobacterium tuberculosis Antibiotic Resistance and
523    Lineage from Whole-Genome Sequencing Data. J Clin Microbiol. 2015;53(6):1908-14.
524    22.    Coll F, McNerney R, Preston MD, Guerra-Assunção JA, Warry A, Hill-Cawthorne G, et
525    al. Rapid determination of anti-tuberculosis drug resistance from whole-genome sequences.
526    Genome Med. 2015;7(1):51.
527    23.    Doughty EL, Sergeant MJ, Adetifa I, Antonio M, Pallen MJ. Culture-independent
528    detection and characterisation of Mycobacterium tuberculosis and M. africanum in sputum
529    samples using shotgun metagenomics on a benchtop sequencer. PeerJ. 2014;2:e585.
530    24.    Bjorn-Mortensen K, Zallet J, Lillebaek T, Andersen AB, Niemann S, Rasmussen EM, et
531    al. Direct DNA Extraction from Mycobacterium tuberculosis Frozen Stocks as a Reculture-
532    Independent Approach to Whole-Genome Sequencing. J Clin Microbiol. 2015;53(8):2716-9.
533    25.    Depledge DP, Palser AL, Watson SJ, Lai IY, Gray ER, Grant P, et al. Specific capture
534    and whole-genome sequencing of viruses from clinical samples. PLoS One.
535    2011;6(11):e27805.
536    26.    Hanekom M, Streicher EM, Van de Berg D, Cox H, McDermid C, Bosman M, et al.
537    Population structure of mixed Mycobacterium tuberculosis infection is strain genotype and
538    culture medium dependent. PLoS One. 2013;8(7):e70178.
539    27.    Martin A, Herranz M, Ruiz Serrano MJ, Bouza E, Garcia de Viedma D. The clonal
540    composition of Mycobacterium tuberculosis in clinical specimens could be modified by
541    culture. Tuberculosis (Edinb). 2010;90(3):201-7.
542    28.    Metcalfe JZ, Streicher E, Theron G, Colman RE, Penaloza R, Allender C, et al.
543    Mycobacterium tuberculosis subculture results in loss of potentially clinically relevant
544    heteroresistance. Antimicrob Agents Chemother. 2017.
545    29.    Jain P, Weinrick BC, Kalivoda EJ, Yang H, Munsamy V, Vilcheze C, et al. Dual-Reporter
546    Mycobacteriophages (Phi2DRMs) Reveal Preexisting Mycobacterium tuberculosis Persistent
547    Cells in Human Sputum. MBio. 2016;7(5).
548    30.    Mukamolova GV, Turapov O, Malkin J, Woltmann G, Barer MR. Resuscitation-
549    promoting factors reveal an occult population of tubercle Bacilli in Sputum. Am J Respir Crit
550    Care Med. 2010;181(2):174-80.

551    31.    Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K, et al.
552    Same-day diagnostic and surveillance data for tuberculosis via whole genome sequencing of
553    direct respiratory samples. J Clin Microbiol. 2017.
554    32.    Brown AC, Bryant JM, Einer-Jensen K, Holdstock J, Houniet DT, Chan JZ, et al. Rapid
555    Whole-Genome Sequencing of Mycobacterium tuberculosis Isolates Directly from Clinical
556    Samples. J Clin Microbiol. 2015;53(7):2230-7.
557    33.    Doyle RM, Burgess C, Williams R, Gorton R, Booth H, Brown J, et al. Direct whole
558    genome sequencing of sputum accurately identifies drug resistant Mycobacterium
559    tuberculosis faster than MGIT culture sequencing. J Clin Microbiol. 2018.
560    34.    Nimmo C, Doyle R, Burgess C, Williams R, Gorton R, McHugh TD, et al. Rapid
561    identification of a Mycobacterium tuberculosis full genetic drug resistance profile through
562    whole genome sequencing directly from sputum. Int J Infect Dis. 2017;62:44-6.
563    35.    Consortium CR, the GP, Allix-Beguec C, Arandjelovic I, Bi L, Beckert P, et al. Prediction
564    of Susceptibility to First-Line Tuberculosis Drugs by DNA Sequencing. N Engl J Med.
565    2018;379(15):1403-15.
566    36.    Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, et al. Whole-
567    genome sequencing for prediction of Mycobacterium tuberculosis drug susceptibility and
568    resistance: a retrospective cohort study. Lancet Infect Dis. 2015;15(10):1193-202.
569    37.    Malik S, Willby M, Sikes D, Tsodikov OV, Posey JE. New insights into fluoroquinolone
570    resistance in Mycobacterium tuberculosis: functional genetic analysis of gyrA and gyrB
571    mutations. PLoS One. 2012;7(6):e39754.
572    38.    Coll F, McNerney R, Guerra-Assuncao JA, Glynn JR, Perdigao J, Viveiros M, et al. A
573    robust SNP barcode for typing Mycobacterium tuberculosis complex strains. Nat Commun.
574    2014;5:4812.
575    39.    Sobkowiak B, Glynn JR, Houben R, Mallard K, Phelan JE, Guerra-Assuncao JA, et al.
576    Identifying mixed Mycobacterium tuberculosis infections from whole genome sequence
577    data. BMC Genomics. 2018;19(1):613.
578    40.    Konstantinidis KT, Tiedje JM. Prokaryotic taxonomy and phylogeny in the genomic
579    era: advancements and challenges ahead. Curr Opin Microbiol. 2007;10(5):504-9.
580    41.    Heym B, Alzari PM, Honore N, Cole ST. Missense mutations in the catalase-
581    peroxidase gene, katG, are associated with isoniazid resistance in Mycobacterium
582    tuberculosis. Mol Microbiol. 1995;15(2):235-45.
583    42.    Coll F, Phelan J, Hill-Cawthorne GA, Nair MB, Mallard K, Ali S, et al. Genome-wide
584    analysis of multi- and extensively drug-resistant Mycobacterium tuberculosis. Nat Genet.
585    2018;50(2):307-16.
586    43.    Sam IC, Drobniewski F, More P, Kemp M, Brown T. Mycobacterium tuberculosis and
587    rifampin resistance, United Kingdom. Emerg Infect Dis. 2006;12(5):752-9.
588    44.    Larsen MH, Biermann K, Tandberg S, Hsu T, Jacobs WR. Genetic Manipulation of
589    Mycobacterium tuberculosis. Curr Protoc Microbiol. 2007;Chapter 10:Unit 10A.2.
590    45.    Krueger F. TrimGalore  [Available from: https://github.com/FelixKrueger/TrimGalore.
591    46.    Bushnell B. BBMap  [cited 2019 18 Mar]. Available from: https://jgi.doe.gov/data-
592    and-tools/bbtools/.
593    47.    Picard Tools: Broad Institute;  [Available from:
594    http://broadinstitute.github.io/picard/.
595    48.    Okonechnikov K, Conesa A, Garcia-Alcalde F. Qualimap 2: advanced multi-sample
596    quality control for high-throughput sequencing data. Bioinformatics. 2016;32(2):292-4.

597    49.    Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing.
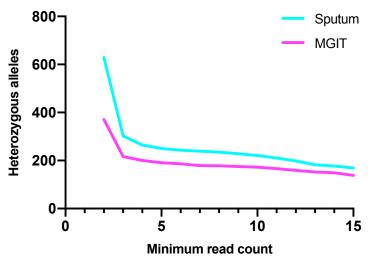598    arXiv preprint arXiv:12073907. 2012.
599    50.    Garrison E. Vcflib  [Available from: https://github.com/vcflib/vcflib.
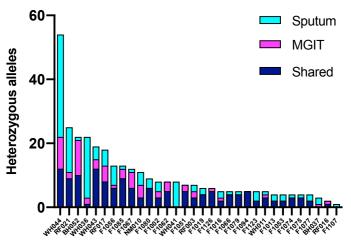600    51.    Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2:
601    somatic mutation and copy number alteration discovery in cancer by exome sequencing.
602    Genome Res. 2012;22(3):568-76.
603    52.    Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites: rapid
604    efficient extraction of SNPs from multi-FASTA alignments. Microb Genom.
605    2016;2(4):e000056.
606    53.    Seemann T. snp-dists  [Available from: https://github.com/tseemann/snp-dists.
607    54.    Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
608    large phylogenies. Bioinformatics. 2014;30(9):1312-3.
609    55.    Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using
610    exact alignments. Genome Biol. 2014;15(3):R46.
611    56.    Lassalle F, Spagnoletti M, Fumagalli M, Shaw L, Dyble M, Walker C, et al. Oral
612    microbiomes from hunter-gatherers and traditional farmers reveal shifts in commensal
613    balance and pathogen load linked to diet. Mol Ecol. 2018;27(1):182-95.
614    57.    Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
615    architecture and applications. BMC Bioinformatics. 2009;10:421.
616    58.    Yang C, Luo T, Shen X, Wu J, Gan M, Xu P, et al. Transmission of multidrug-resistant
617    Mycobacterium tuberculosis in Shanghai, China: a retrospective observational study using
618    whole-genome sequencing and epidemiological investigation. Lancet Infect Dis.
619    2017;17(3):275-84.
620    59.    Eldholm V, Monteserin J, Rieux A, Lopez B, Sobkowiak B, Ritacco V, et al. Four
621    decades of transmission of a multidrug-resistant Mycobacterium tuberculosis outbreak
622    strain. Nat Commun. 2015;6:7119.
623

624

## Figure legends

626

627 Figure 1. Variation in total number of heterozygous alleles (HAs) identified across all 36

628 patients in sequences generated from sputum and MGIT depending on minimum supporting

629 read count threshold.

630

631 Figure 2. Number of heterozygous alleles (HAs) found in directly sequenced sputum only

632 (sputum), MGIT (MGIT) only or in both samples (shared) by patient.

633