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]       c.nimmo.04@cantab.net

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

14      # Abstract

15

16      **Background**

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

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

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

20      compared *M. tuberculosis* whole genome sequences generated from 39 paired clinical

21      samples. In one sample DNA was extracted directly from sputum then enriched with

22      custom-designed SureSelect (Agilent) oligonucleotide baits and in the other it was extracted

23      from mycobacterial growth indicator tube culture.

24

25      **Results**

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

27      than that from mycobacterial growth incubator tube culture. This was demonstrated by

28      more variants present as heterozygous alleles (HAs) where both a variant and wild type

29      allele were present within a given sample ($p<0.001$) and greater within-sample Shannon

30      diversity ($p<0.001$). Seven genes with high within-sample diversity have previously been

31      identified as targets for positive selection, highlighting their potential role in adaptation to

32      survival within the host and under drug pressure. Resistance associated variants present as

33      HAs occurred in six patients, and in four cases may have provided a genotypic explanation

34      for phenotypic resistance.

35

36      **Conclusions**

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

38    diversity and may allow detection of mycobacteria that are not actively replicating.

39

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

41    sequencing; sputum; within-patient diversity; heteroresistance

## Background

42

43

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

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

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

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

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

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

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

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

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

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

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

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

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

57

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

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

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

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

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

63    genetic diversity between strains(10), but more detailed analysis of individual patient

64    samples with WGS has identified genetically separate bacterial subpopulations in sequential

65    sputum samples(11-15) and across different anatomical sites(16). This within-patient

4

66      diversity can occur as a result of mixed infection with genetically distinct strains or within-

67      host evolution of a single infecting strain(17).

68

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

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

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

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

73      bacterial population is drug susceptible while the remainder is resistant(18). Identification

74      of genetic diversity within clinical samples is important as it may improve detection of RAVs

75      over currently available genetic tests and consensus-level WGS(18). Identifying RAVs could

76      improve individualised therapy, prevent acquired resistance(12), and give insight into

77      bacterial adaptation to the host.

78

79      *M. tuberculosis* WGS is usually performed on cultured isolates to obtain sufficient purified

80      mycobacterial DNA. However, the culture process can change the population structure from

81      that of the original sample due to genetic drift (random loss of lineages) and/or the

82      selection of subpopulations more suited to growth in culture(19-21), and repeated

83      subculture leads to loss of genetic diversity and heteroresistance(22). Additionally, in the

84      normal course of *M. tuberculosis* infection, some bacteria exist as viable non-culturable

85      persister organisms that are hypothesised to cause the high relapse rate seen following

86      treatment of insufficient duration(23). These organisms are likely to be missed by any

87      sequencing method reliant on culture.

88

89     WGS directly from sputum without enrichment is challenging(24). It has recently been

90     improved by depleting human DNA during DNA extraction(25). We have previously reported

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

92     *M. tuberculosis* DNA directly from sputum(26) and demonstrated its utility in determining a

93     rapid genetic drug resistance profile(27, 28).

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(25), whereas another study of mostly drug

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

100    HAs relative to culture(27). Here we reanalyse heterozygous alleles (HAs) present in that

101    study(27) in addition to newly collected samples from patients with MDR-TB, use a more

102    sensitive analysis to measure overall within-sample genetic diversity 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 39 patients from both mycobacterial growth

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

111    predominantly of black African ethnicity (80%) and 50% were HIV positive (Table 1). First-

112    line phenotypic drug susceptibility testing (DST) results identified 22 patients with MDR-TB

6

113    and two with rifampicin monoresistance. In addition there were three isoniazid

114    monoresistant patients and ethambutol resistance was detected in 8 patients. Second-line

115    phenotypic DST was performed for patients with rifampicin-resistant or MDR-TB and

116    identified one case of kanamycin resistance (Table 2).

117

118    We observed greater median coverage depth in sputum-derived sequences than MGIT

119    sequences (164.3 vs 136.6, p=0.068). A genotypic susceptibility profile was determined by

120    evaluating MGIT WGS for consensus-level RAVs using a modified version of a publicly

121    available list(29). Genotypic RAVs predicted all rifampicin phenotypic resistance and >90% of

122    isoniazid phenotypic resistance. Ethambutol genotypic RAVs were poorly predictive of

123    phenotypic resistance in line with findings from other studies(30) (Table 2). The patient with

124    kanamycin phenotypic resistance was correctly identified by an *rrs* a1401g RAV. No full

125    phenotypic fluoroquinolone phenotypic resistance was identified, but several colonies from

126    patient F1013 did grow in the presence of ofloxacin (although not enough to be classified as

127    resistant). The consensus sequences from this patient harboured a *gyrB* E501D mutation

128    which is believed to confer resistance to moxifloxacin but not other fluoroquinolones, which

129    may explain the borderline phenotypic DST result(31).

130

131    **Genetic Diversity**

132

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

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

135    Four previously sequenced strains from KwaZulu-Natal were included(32). As expected, all

136    paired sequences were closely related, with a mean difference of 1.30 (range 0-9) single

7

137     nucleotide polymorphisms (SNPs). Samples from patients F1066 and F1067 were closely

138     related with only one consensus-level SNP separating all four consensus sequences. There

139     was no obvious epidemiological link between these patients (although this study was not

140     designed to collect comprehensive epidemiological information) and they lived 20km apart

141     in Durban. However, both patients were admitted contemporaneously to an MDR

142     treatment facility and sampled on the same day. DNA extraction and sequencing occurred

143     on different runs so the close genetic linkage may represent direct transmission within a

144     hospital setting, a community transmission chain or an unlikely cross-contamination during

145     sample collection.

146

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

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

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

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

151     minimum of five supporting reads there was an exponential increase in the number of HAs

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

153     we used a threshold of a minimum of five supporting reads.

154

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

156     mixed infection we used a molecular barcode(33) to scan all HAs for a panel of 413

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

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

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

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

8

161    or sublineage. This suggests that the genetic diversity identified is mostly or exclusively due

162    to within-host evolution, although there remains a small possibility that mixed infections

163    with two strains from the same sub-lineage could have occurred.

164

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

166    looked at the location of genetic diversity within the *M. tuberculosis* genome. Variants were

167    called in the MGIT and sputum sequences for each patient and classified as present in MGIT

168    only, sputum only or shared (present in both). HAs were widely dispersed across the

169    genome at similar sites in both sputum and MGIT samples but some genes had multiple HAs

170    (Table 3). The highest genetic diversity was found in the ribosomal RNA (rRNA) genes (*rrs*

171    and *rrl*) with 358 HAs, of which 98.6% were only found in sputum-derived sequences.

172

173    As rRNA contains regions that are highly conserved across bacteria, it was considered a

174    possibility that SureSelect baits targeting rRNA genes were capturing both *M. tuberculosis*

175    and other bacterial species. To evaluate this, metagenomic assignment was performed on

176    all reads. Sampling reads not assigned to *M. tuberculosis* (i.e. presumed contaminants from

177    other bacteria) and performing a BLAST search against *M. tuberculosis* 16S and 23S rRNA

178    genes indicated that a sizeable proportion of these reads from directly sequenced sputum

179    had a BLAST hit of at least 30 bases (median 11% v 0% of equivalent reads from MGIT

180    sequencing, $p < 0.001$, Supplementary Material: Figure 2). The taxonomic assignment of

181    these reads were indeed typical of genera composing the oral flora, with a high

182    representation of *Actinomyces, Fusobacterium, Prevotella,* and *Streptococcus*

183    (Supplementary Material: Figure 3).

184

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

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

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

188  *M. tuberculosis* sequences. Importantly the frequency of HAs in other highly diverse genes

189  between sequencing strategies was more balanced (Table 3). Pertinently seven of these

190  genes (*Rv1319c(34), lppB(35), Rv2082(35), ppsA(34, 36), ponA1(36), lppA(37),* and *pks12(35,*

191  *36)*) with high numbers of HAs have been previously identified as highly diverse in

192  comparative genomic studies suggesting the detected HAs are not artefactual. The

193  frequencies at which HAs in these genes were present in MGIT and sputum is shown in

194  Supplementary Material Figure 4.

195

196  After confirming the absence of mixed infections and removing rRNA gene sequences we

197  compared the frequency of HAs in sputum and MGIT. There were 2048 variants across the

198  dataset that were present as a HA in either MGIT, sputum or both sequences (Table 4).

199  Variants present in both MGIT and sputum derived sequences were more likely to be

200  present as a HA in the sputum-derived sequence (3.2% v 1.9%, p<0.0001). Of the other

201  variants present as HAs, 821 were unique to direct sputum sequencing and 153 were unique

202  to MGIT sequencing (Table 4). Variants found only in sputum were more likely to be

203  heterozygous than those found only in MGIT or in both (p<0.0001). The distribution of HAs

204  by patient is shown in Figure 2A. HAs found only by one modality were more likely to be

205  SNPs than shared HAs, where the majority were insertions or deletions. The ratio of non-

206  synonymous to synonymous HAs was similar for those that were shared or MGIT only, but

207  was lower for sputum only HAs. Frameshift mutations were most prevalent among shared

208  HAs (Table 4).

209

210    To confirm our findings of increased diversity in sputum we calculated mean within-sample

211    diversity ($H$), excluding rRNA genes and repeat regions (see methods). The mean diversity

212    was significantly greater in sputum than MGIT-derived sequences (Figure 2B: 0.116±0.078 v

213    0.054±0.026 , mean $H_{sputum}/H_{MGIT}$ = 2.66, p=3.0 x $10^{-5}$).

214

215    **Genetic diversity in drug resistance genes**

216

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

218    individually assessed. Five of the 39 patients had RAVs present as HAs in at least one gene,

219    which are shown in Table 5. F1002 had three compensatory mutations in *rpoC* present at

220    HAs in both sequences. F1007 had high-level phenotypic isoniazid resistance despite wild

221    type *katG* and *inhA* genes, but did have two *ahpC* promoter variants present as HAs. Neither

222    of these variants are reported frequently but both have been previously associated with

223    resistance in limited numbers of samples(38). As described above F1066 and F1067 were

224    highly related with only one consensus SNP difference between all four sequences. Both had

225    phenotypic high level isoniazid resistance with no consensus-level *katG* or *inhA* mutation,

226    but had frameshift *katG* mutations present as HAs which have the potential to cause

227    resistance(39). F1066 and RF021 had *Rv1979c* and *pncA* mutations respectively at low

228    frequency in sputum only which have the potential to confer phenotypic resistance to

229    clofazimine (*Rv1979c*) and pyrazinamide (*pncA*), although no phenotypic testing was

230    performed for these drugs.

231

232    Discussion

233

234    In this study we whole genome sequenced DNA from sputum and MGIT culture in paired

235    samples from 39 patients and compared within-patient genetic diversity of the bacterial

236    genome identified from each source. All paired sequences were closely related at the

237    consensus level, and WGS predicted phenotypic drug susceptibility with over 90% sensitivity

238    and specificity for rifampicin and isoniazid in line with published data(40).

239

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

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

242    correlation between phenotypic and genotypic drug resistance profiles(18) and can identify

243    emerging drug resistance(11, 12). Here we have demonstrated that significantly more

244    genetic diversity is identified by WGS performed directly from enriched sputum than MGIT

245    culture. Not including a culture step avoids the introduction of bias towards culture-adapted

246    subpopulations and the impact of random chance and is also likely to incorporate DNA from

247    viable non-culturable mycobacteria. A reduction in genetic diversity has previously been

248    shown with sequential *M. tuberculosis* subculture(19, 22), but was not confirmed by a study

249    performing WGS directly from sputum(25). However, the 16 paired sputum and MGIT

250    samples compared by Votintseva(25) had a minimum of 5x coverage compared to a

251    minimum 40x coverage in this study, and were likely to contain less genetic material as they

252    were surplus clinical rather than dedicated research samples.

253

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

255    this is due to non-mycobacterial DNA hybridising to the capture baits — a conclusion borne

12

256     out by the taxonomic assignment of reads aligning to these genes in common oral bacteria.

257     We therefore exclude these from further analysis, and recommend others using enrichment

258     from sputum do similarly. We use two methods to evaluate within-sample *M. tuberculosis*

259     genetic diversity. First, we demonstrate increased diversity when sequencing directly from

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

261     MGIT culture. We also observed significantly higher genetic diversity in sputum-derived

262     sequences by comparing the Shannon diversity of variable sites across pairs of samples.

263

264     Many of the genes with high levels of within-sample diversity are also reported to be targets

265     for convergent evolution, independently accumulating similar mutations on a global scale.

266     This implies that diversity seen on a macroevolutionary scale has a basis in microevolution,

267     and reinforces the importance of accurately characterising the biological function of these

268     genes and their products to aid the identification of new therapeutic targets. Two-thirds of

269     the patients with MDR-TB had already been treated for drug-sensitive TB, and the diversity

270     identified in sputum samples may therefore represent early adaptation to drug pressure.

271     Importantly, as direct sputum sequencing does not rely on live mycobacteria, DNA from

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

273     mutations are likely to be represented as HAs.

274

275     In four patients, RAVs present as HAs provided a likely genotypic basis for otherwise

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

277     this study, the excess of heterozygous known RAVs in directly sequenced sputum is not

278     statistically significant. However the presence of heterozygous RAVs in both MGIT and

279     sputum sequences reinforces the biological importance of these mutations.

280

281    A limitation of this study is that it can be difficult to distinguish low frequency variants from

282    sequencing error. Ideally low frequency variants could be confirmed by resequencing the

283    same DNA samples. To reduce the risk of sequencing errors yet still identify genetic diversity

284    we used the lowest minimum read threshold at which the number of HAs remained stable.

285    Also, it is reassuring that of all fixed and heterozygous variants called, more than 93% were

286    identified from both DNA sources.

287

288    Conclusions

289

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

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

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

293    It also has potential to identify low frequency RAVs that may further enhance genotypic-

294    phenotypic drug resistance correlation.

295

296    Methods

297

298    **Patient enrolment**

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

300    in the study. Patients were recruited in London, UK (n=15) and Durban, South Africa (n=24).

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

302    rifampicin resistance. Two sputum samples were collected prior to initiating treatment, with

303    one inoculated into mycobacterial growth indicator tube (MGIT) culture (BD, NJ, USA) and

304    the other used for direct DNA extraction.

305

306    **Ethics, Consent and Permissions**

307

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

309    the London study was granted by NHS National Research Ethics Service East Midlands—

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

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

312    BE022/13).

313

314    **Microbiology**

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

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

317    Health England. Phenotypic DST for Durban samples was performed using the solid agar

318    proportion method (Supplementary Material: Methods).

319

320    **DNA extraction and sequencing**

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

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

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

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

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

326     with DiaSorin Liaison Ixt (DiaSorin, Italy) or CTAB(41). NEBNext Ultra II DNA (New England

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

328

329     Sputum samples for direct sequencing were similarly heat killed processed as for MGIT

330     samples. DNA extraction was performed with mechanical ribolysis followed by purification

331     using DiaSorin Liaison Ixt (DiaSorin, Italy) or DNeasy blood & tissue kit (Qiagen,

332     Germany)(41). Target enrichment was performed using SureSelect with a custom-designed

333     bait set providing coverage of the entire *M. tuberculosis* genome as described

334     previously(27). Batches of 48 multiplexed samples were sequenced on a NextSeq (Illumina,

335     CA, USA).

336

337     **Bioinformatic analysis**

338     Bioinformatic analysis was performed with CLC Genomics Workbench v11.0 (Qiagen,

339     Germany). DNA sequence reads were aligned to an H37Rv reference genome as detailed in

340     the Supplementary Material Methods section (GenBank accession NC_000962.3). All

341     samples had minimum 98% 1x reference genome coverage and mean coverage depth 40x

342     across the genome. Variants falling within or near hypervariable elements were excluded

343     (Supplementary Material: Table 1). A consensus sequence was extracted and used to

344     determine the genotypic drug susceptibility profile. To construct the maximum likelihood

345     phylogenetic tree, variants were called against the reference genome using VarScan v2.3.9

346     (Supplementary Material: Methods).

347

348     For the initial analysis of genetic diversity, variants were included if supported by $\geq 2$ reads,

349     with $\geq 1$ forward and reverse read. The minimum supporting read threshold was increased in

16

350    a stepwise fashion from 2 to 20. Further analyses were performed on variant tracks where

351    variants were supported ≥5 supporting reads including ≥1 forward and reverse read.

352

353    To compare diversity between paired samples, we first mapped reads to the reference

354    genome using bwa mem v0.7.12(42). After verifying all samples had adequate coverage with

355    qualimap(43) (mean ± standard deviation coverage at 10x: 98.0 ± 1.8%) and realigning

356    indels, variants were called with HaplotypeCaller in GATK v3.3.0(44) (Supplementary

357    Material: Methods). The gvcf files were combined for each pair of samples with

358    CombineGVCFs in GATK then screened to remove sites in variable regions and rRNA genes

359    with vcfintersect in vcflib, resulting in 39 paired gvcf files containing allele depths at variable

360    positions for diversity analysis.

361

362    **Metagenomic assignment**

363    Sequencing reads were classified using Kraken v0.10.6(45) against a custom Kraken

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

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

366    three human genome reference sequences(46). The size of the final database after shrinking

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

368

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

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

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

372    v2.2.28(47) against rRNA genes from the H37Rv reference genome. We only analysed hits of

373    at least 30 bases.

374

375 **Statistics**

376 Statistical analyses were performed with Prism v7.0 (Graphpad, CA, USA). The number of

377 HAs in paired samples were compared using a two-tailed Wilcoxon matched-pairs signed

378 rank test. Numbers of HAs found between groups were compared with chi-squared.

379

380 Within-sample diversity ($H$) was calculated using Shannon diversity from the allele

381 frequencies ($p$). The Shannon index ($H_n$) expresses the positional entropy at each position

382 ($n$), with the mean positional entropy ($H$) indicating greater within-sample diversity:

$$H = \sum_n H_n = \sum_n \sum_{i \in \{A,C,G,T\}} p_{n,i} \log (p_{n,i})$$

383 To make a fair comparison for each pair of samples, after removing indels, sites were

384 included if they contained a variant in at least one and had a depth coverage $\geq$30 in both.

385 We found that the depth coverage cutoff had no qualitative effect on the conclusions. The

386 difference in mean within-sample diversity depending on DNA source was compared with a

387 two-tailed Wilcoxon matched-pairs signed rank test.

388

389 Abbreviations

390

| DST | drug susceptibility testing |
|---|---|
| HA | heterozygous allele |
| MDR-TB | multidrug resistant-tuberculosis |
| MGIT | mycobacterial growth indicator tube |

18

| RAV | resistance-associated variant |
|-----|-------------------------------|
| rRNA | ribosomal RNA |
| SNP | single nucleotide polymorphism |
| TB | tuberculosis |
| WGS | whole genome sequencing |

391

392

393

394     # Declarations

395

396     **Ethics approval and consent to participate**

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

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

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

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

401     BE022/13).

402

403     **Consent for publication**

404     Not applicable

405

406     **Availability of data and materials**

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

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

409

410     **Competing interests**

411     The authors declare that they have no competing interests.

412

413     **Funding**

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

415     203583/Z/16/Z. This work was additionally funded by National Institute for Health Research

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

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

420

421    **Authors' contributions**

422    Study conception: JB, ASP

423    Data collection: CB, KB

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

425    Drafting of manuscript: CN, LPS

426    Revision of manuscript: FB, JB, ASP

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

428

429    **Acknowledgements**

21

432     Tables

433

|  | Count / mean | Range/percentage |
|---|---|---|
| **Age** | 36.7 | 22 – 64 |
| **Male sex** | 23/37 | 59.0% |
| **Ethnicity** <br><br> Asian <br><br> Black African <br><br> Caucasian | <br><br> 3/35 <br><br> 28/35 <br><br> 4/35 | <br><br> 8.6% <br><br> 80.0% <br><br> 11.4% |
| **HIV positive** <br><br> CD4 count (median) <br><br> On antiretroviral therapy at time of diagnosis | 19/38 <br><br> 296.5* <br><br> 8/19** | 50.0% <br><br> 17 – 707 <br><br> 42.1% |

434

435     Table 1. Baseline patient characteristics for 39 patients (or as otherwise specified where

436     data were missing). *Data missing for 1 patient.

437

22

438

| Drug | Resistance by phenotypic DST | Resistance by genotypic DST | Genotypic DST sensitivity | Genotypic DST specificity |
|---|---|---|---|---|
| *First-line drugs* | | | | |
| Rifampicin | 24/37 (64.9%) | 24/39 | 24/24 (100%)* | 24/24 (100%) |
| Isoniazid | 25/37 (67.6%) | 24/39 | 23/25 (92.0%) | 23/24 (95.8%) |
| Ethambutol | 8/37 (21.6%) | 17/39 | 8/8 (100%) | 8/17 (47.1%) |
| *Second-line drugs* | | | | |
| Ofloxacin | 0/24 (0.0%) | 1/24 | N/A | 0/1 (0%)** |
| Kanamycin | 1/24 (4.2%) | 1/24 | 1/1 (100%) | 1/1 (100%) |

439

440    Table 2. Phenotypic and genotypic drug susceptibility testing (DST) results and sensitivity

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

442    first-line drugs for 37 of the 39 patients, and for second-line drugs for 24 patients who

443    demonstrated rifampicin drug resistance. *In two directly-sequenced sputum samples

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

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

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

447    fluoroquinolones (see main text).

448

| Gene | Heterozygous Allele Count | | | | Gene length (base pairs) | Hypothesised gene function |
| | Shared | MGIT only | Sputum only | Total | | |
|---|---|---|---|---|---|---|
| rrs | 2 | 3 | 180 | 185 | 3138 | 23S rRNA |
| rrl | 0 | 0 | 173 | 173 | 1537 | 16S rRNA |
| Rv1319c | 70 | 1 | 24 | 95 | 1608 | Metabolism and respiration |
| lppB | 7 | 6 | 10 | 23 | 663 | Surface lipoprotein |
| Rv2561 | 21 | 0 | 0 | 21 | 294 | Unknown function |
| Rv3424c | 1 | 1 | 19 | 21 | 363 | Unknown function |
| Rv2082 | 16 | 1 | 2 | 19 | 2166 | Unknown function |
| ppsA | 7 | 0 | 11 | 18 | 1059 | GGPP synthetase (lipid synthesis) |
| Rv1435c | 3 | 9 | 6 | 18 | 609 | Secreted protein |
| ponA1 | 5 | 3 | 9 | 17 | 2037 | Cell wall biosynthesis |
| Rv2277c | 2 | 0 | 15 | 17 | 906 | Metabolism and respiration |
| vapC31 | 5 | 0 | 12 | 17 | 429 | Possible toxin |
| Rv2823c | 5 | 2 | 9 | 16 | 2430 | Unknown function |
| lppA | 1 | 3 | 11 | 15 | 660 | Surface lipoprotein |
| pks12 | 5 | 5 | 4 | 14 | 12456 | MPM synthesis (lipid metabolism) |

449

450     Table 3. Genes with the most heterozygous alleles (HAs) identified across the entire dataset.

451

452

|  | Shared variants | MGIT only variants | Sputum only variants | Total |
|---|---|---|---|---|
| *All variants vs H37Rv (fixed or heterozygous)* | | | | |
| Total variants | 33 153 | 1162 | 1217 | 35532 |
| *Variants vs H37Rv present as heterozygous alleles (HAs) only* | | | | |
| Total variants present as HAs (% of total variants) | MGIT 645 (1.9%) Sputum 1074 (3.2%) | 153 (13.2%) | 821 (67.5%) | 2048 (5.8%) |
| Median HAs per sample | 21 | 3 | 15 | 40 |
| Variant type (% all HAs) SNP MNP Insertion Deletion Replacement | 500 (46.6%) 12 (1.1%) 303 (28.2%) 259 (24.1%) 0 (0.0%) | 127 (83.0%) 1 (0.7%) 8 (5.2%) 16 (10.5%) 1 (0.7%) | 708 (86.2%) 24 (2.9%) 31 (3.8%) 57 (6.9%) 1 (0.1%) | 1335 (65.2%) 37 (1.8%) 342 (16.7%) 332 (16.2%) 2 (0.1%) |
| Coding change (% all HAs) Non-synonymous Synonymous Intergenic | 395 (36.8%) 159 (14.8%) 520 (48.4%) | 79 (51.6%) 32 (20.9%) 42 (27.5%) | 318 (38.7%) 171 (20.8%) 332 (40.4%) | 792 (38.7%) 362 (17.7%) 894 (43.7%) |
| Non-synon/synon ratio | 2.48 | 2.47 | 1.86 | 2.19 |
| Stop codon (% of all non-synonymous HAs) | 4 (1.0%) | 1 (1.3%) | 9 (2.8%) | 14 (1.8%) |
| Frameshift (% of all non-synonymous HAs) | 185 (46.8%) | 19 (24.1%) | 47 (14.8%) | 251 (31.7%) |

453

454 Table 4. Variants identified in MGIT derived, sputum derived, or both sequences from paired

455 samples. Values given represent totals for the 39 paired samples. SNP = single nucleotide

456 polymorphism; MNP = multi-nucleotide polymorphism.

457

458

| Patient ID | Phenotypic resistance | Mutation | Frequency (MGIT/sputum) | Description |
|---|---|---|---|---|
| F1002 | Rifampicin | *rpoB* S450L | 100%/100% | High confidence resistance mutation |
| F1002 | Rifampicin | *rpoC* G332R(48) | 82.6%/21.7% | Putative compensatory mutations |
| F1002 | Rifampicin | *rpoC* L516P(48) | 12.7%/7.7% | |
| F1002 | Rifampicin | *rpoC* P1040S(49) | 21.7%/12.3% | |
| F1007 | Isoniazid (high) | *ahpC* c-52t(38) | 60.0%/50.7% | Rare, have been associated with resistance |
| F1007 | Isoniazid (high) | *ahpC* g-48a(38) | 28.6%/30.3% | |
| F1061 | Rifampicin | *rpoB* H445D | 16.1%/0.0%* | High confidence resistance mutation |
| F1061 | Rifampicin | *rpoB* S450W | 84.4%/0.0%* | High confidence resistance mutation |
| 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% | |

459

27

460      Table 5. Resistance-associated variants present as heterozygous alleles (HAs). *These

461      mutations were also present in sputum but due to low coverage of the area (3 and 4 reads

462      respectively) variant calling criteria were not met.

463

464

465   ## Figure legends

466

467   Figure 1. Variation in total number of heterozygous alleles (HAs) identified across all 39

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

469   read count threshold. Direct sputum samples indicated by red squares, MGIT samples blue

470   circles.

471

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

473   (sputum), MGIT (MGIT) only or in both samples (shared) by patient. (B) Mean Shannon

474   diversity at variable positions across pairs of samples ($H$) as calculated for MGIT and

475   sputum-derived sequences. Size of point indicates number of variable positions considered

476   (see Methods).

477

# References

1.      Global Tuberculosis Report 2018. Geneva: World Health Organization; 2018 2018.

2.      Murray CJ, Ortblad KF, Guinovart C, Lim SS, Wolock TM, Roberts DA, et al. Global,
regional, and national incidence and mortality for HIV, tuberculosis, and malaria during
1990-2013: a systematic analysis for the Global Burden of Disease Study 2013. Lancet.
2014;384(9947):1005-70.

3.      Dheda K, Gumbo T, Maartens G, Dooley KE, McNerney R, Murray M, et al. The
epidemiology, pathogenesis, transmission, diagnosis, and management of multidrug-
resistant, extensively drug-resistant, and incurable tuberculosis. Lancet Respir Med. 2017.

4.      WHO treatment guidelines for drug-resistant tuberculosis. World Health
Organization; 2016.

5.      Trauner A, Liu Q, Via LE, Liu X, Ruan X, Liang L, et al. The within-host population
dynamics of Mycobacterium tuberculosis vary with treatment efficacy. Genome Biol.
2017;18(1):71.

6.      Olaru ID, Lange C, Heyckendorf J. Personalized medicine for patients with MDR-TB. J
Antimicrob Chemother. 2016;71(4):852-5.

7.      Pasipanodya JG, McIlleron H, Burger A, Wash PA, Smith P, Gumbo T. Serum drug
concentrations predictive of pulmonary tuberculosis outcomes. J Infect Dis.
2013;208(9):1464-73.

8.      Cegielski JP, Kurbatova E, van der Walt M, Brand J, Ershova J, Tupasi T, et al.
Multidrug-Resistant Tuberculosis Treatment Outcomes in Relation to Treatment and Initial
Versus Acquired Second-Line Drug Resistance. Clin Infect Dis. 2016;62(4):418-30.

501  9.      Satta G, Lipman M, Smith GP, Arnold C, Kon OM, McHugh TD. Mycobacterium

502     tuberculosis and whole-genome sequencing: how close are we to unleashing its full

503     potential? Clin Microbiol Infect. 2017.

504     10.     Sreevatsan S, Pan X, Stockbauer KE, Connell ND, Kreiswirth BN, Whittam TS, et al.

505     Restricted structural gene polymorphism in the Mycobacterium tuberculosis complex

506     indicates evolutionarily recent global dissemination. Proceedings of the National Academy

507     of Sciences of the United States of America. 1997;94(18):9869-74.

508     11.     Sun G, Luo T, Yang C, Dong X, Li J, Zhu Y, et al. Dynamic population changes in

509     Mycobacterium tuberculosis during acquisition and fixation of drug resistance in patients. J

510     Infect Dis. 2012;206(11):1724-33.

511     12.     Merker M, Kohl TA, Roetzer A, Truebe L, Richter E, Rüsch-Gerdes S, et al. Whole

512     genome sequencing reveals complex evolution patterns of multidrug-resistant

513     Mycobacterium tuberculosis Beijing strains in patients. PLoS One. 2013;8(12):e82551.

514     13.     Operario DJ, Koeppel AF, Turner SD, Bao Y, Pholwat S, Banu S, et al. Prevalence and

515     extent of heteroresistance by next generation sequencing of multidrug-resistant

516     tuberculosis. PLoS One. 2017;12(5):e0176522.

517     14.     Black PA, de Vos M, Louw GE, van der Merwe RG, Dippenaar A, Streicher EM, et al.

518     Whole genome sequencing reveals genomic heterogeneity and antibiotic purification in

519     Mycobacterium tuberculosis isolates. BMC Genomics. 2015;16(1):857.

520     15.     Eldholm V, Norheim G, von der Lippe B, Kinander W, Dahle UR, Caugant DA, et al.

521     Evolution of extensively drug-resistant Mycobacterium tuberculosis from a susceptible

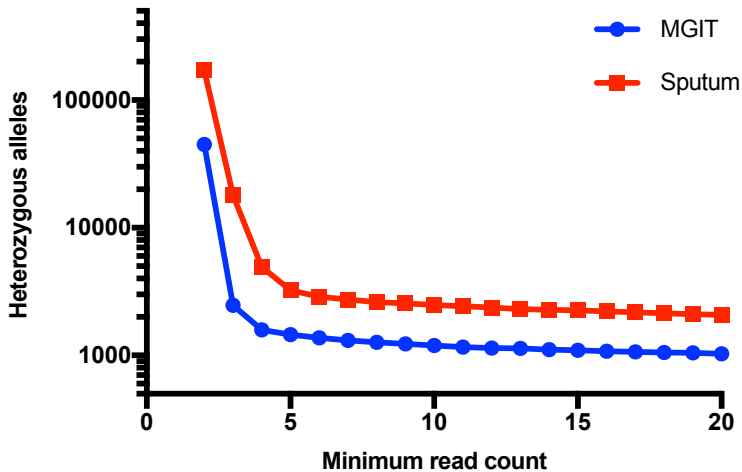522     ancestor in a single patient. Genome Biol. 2014;15(11):490.

31

523     16.     Lieberman TD, Wilson D, Misra R, Xiong LL, Moodley P, Cohen T, et al. Genomic

524     diversity in autopsy samples reveals within-host dissemination of HIV-associated

525     Mycobacterium tuberculosis. Nat Med. 2016;22(12):1470-4.

526     17.     Ford C, Yusim K, Ioerger T, Feng S, Chase M, Greene M, et al. Mycobacterium

527     tuberculosis--heterogeneity revealed through whole genome sequencing. Tuberculosis

528     (Edinb). 2012;92(3):194-201.

529     18.     Metcalfe JZ, Streicher E, Theron G, Colman RE, Allender C, Lemmer D, et al. Cryptic

530     Micro-heteroresistance Explains M. tuberculosis Phenotypic Resistance. Am J Respir Crit

531     Care Med. 2017.

532     19.     Depledge DP, Palser AL, Watson SJ, Lai IY, Gray ER, Grant P, et al. Specific capture

533     and whole-genome sequencing of viruses from clinical samples. PLoS One.

534     2011;6(11):e27805.

535     20.     Hanekom M, Streicher EM, Van de Berg D, Cox H, McDermid C, Bosman M, et al.

536     Population structure of mixed Mycobacterium tuberculosis infection is strain genotype and

537     culture medium dependent. PLoS One. 2013;8(7):e70178.

538     21.     Martin A, Herranz M, Ruiz Serrano MJ, Bouza E, Garcia de Viedma D. The clonal

539     composition of Mycobacterium tuberculosis in clinical specimens could be modified by

540     culture. Tuberculosis (Edinb). 2010;90(3):201-7.

541     22.     Metcalfe JZ, Streicher E, Theron G, Colman RE, Penaloza R, Allender C, et al.

542     Mycobacterium tuberculosis subculture results in loss of potentially clinically relevant

543     heteroresistance. Antimicrob Agents Chemother. 2017.

544     23.     Mukamolova GV, Turapov O, Malkin J, Woltmann G, Barer MR. Resuscitation-

545     promoting factors reveal an occult population of tubercle Bacilli in Sputum. Am J Respir Crit
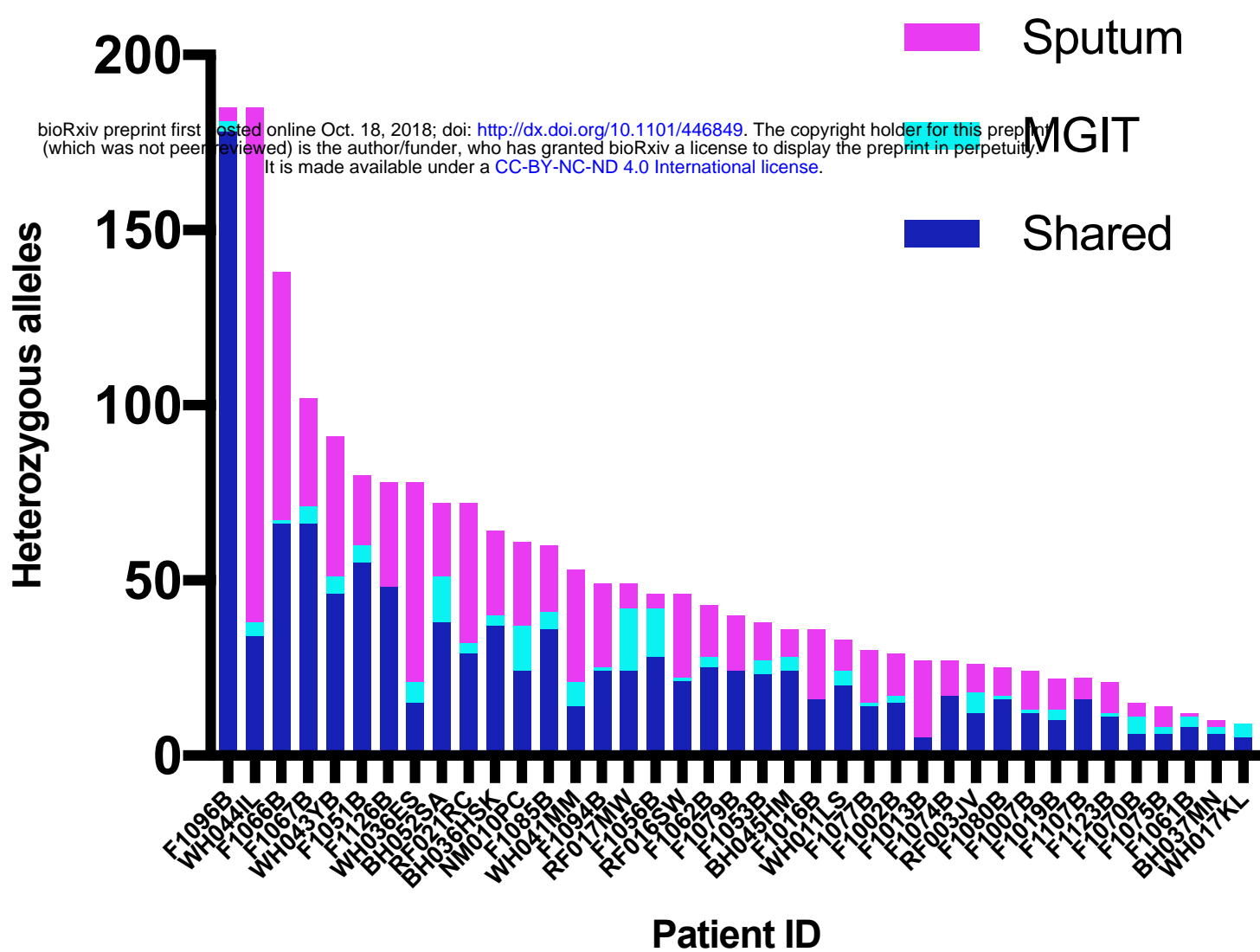
546     Care Med. 2010;181(2):174-80.

547    24.     Doughty EL, Sergeant MJ, Adetifa I, Antonio M, Pallen MJ. Culture-independent

548    detection and characterisation of Mycobacterium tuberculosis and M. africanum in sputum

549    samples using shotgun metagenomics on a benchtop sequencer. PeerJ. 2014;2:e585.

550    25.     Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K, et al.

551    Same-day diagnostic and surveillance data for tuberculosis via whole genome sequencing of

552    direct respiratory samples. J Clin Microbiol. 2017.

553    26.     Brown AC, Bryant JM, Einer-Jensen K, Holdstock J, Houniet DT, Chan JZ, et al. Rapid

554    Whole-Genome Sequencing of Mycobacterium tuberculosis Isolates Directly from Clinical

555    Samples. J Clin Microbiol. 2015;53(7):2230-7.

556    27.     Doyle RM, Burgess C, Williams R, Gorton R, Booth H, Brown J, et al. Direct whole

557    genome sequencing of sputum accurately identifies drug resistant Mycobacterium

558    tuberculosis faster than MGIT culture sequencing. J Clin Microbiol. 2018.

559    28.     Nimmo C, Doyle R, Burgess C, Williams R, Gorton R, McHugh TD, et al. Rapid

560    identification of a Mycobacterium tuberculosis full genetic drug resistance profile through

561    whole genome sequencing directly from sputum. Int J Infect Dis. 2017;62:44-6.

562    29.     Coll F, McNerney R, Preston MD, Guerra-Assunção JA, Warry A, Hill-Cawthorne G, et

563    al. Rapid determination of anti-tuberculosis drug resistance from whole-genome sequences.

564    Genome Med. 2015;7(1):51.

565    30.     Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, et al. Whole-

566    genome sequencing for prediction of Mycobacterium tuberculosis drug susceptibility and

567    resistance: a retrospective cohort study. Lancet Infect Dis. 2015;15(10):1193-202.

568    31.     Malik S, Willby M, Sikes D, Tsodikov OV, Posey JE. New insights into fluoroquinolone

569    resistance in Mycobacterium tuberculosis: functional genetic analysis of gyrA and gyrB

570    mutations. PLoS One. 2012;7(6):e39754.

571    32.    Cohen KA, Abeel T, Manson McGuire A, Desjardins CA, Munsamy V, Shea TP, et al.

572    Evolution of Extensively Drug-Resistant Tuberculosis over Four Decades: Whole Genome

573    Sequencing and Dating Analysis of Mycobacterium tuberculosis Isolates from KwaZulu-

574    Natal. PLoS Med. 2015;12(9):e1001880.

575    33.    Coll F, McNerney R, Guerra-Assuncao JA, Glynn JR, Perdigao J, Viveiros M, et al. A

576    robust SNP barcode for typing Mycobacterium tuberculosis complex strains. Nat Commun.

577    2014;5:4812.

578    34.    Gonzalo X, Drobniewski F. Is there a place for beta-lactams in the treatment of

579    multidrug-resistant/extensively drug-resistant tuberculosis? Synergy between meropenem

580    and amoxicillin/clavulanate. J Antimicrob Chemother. 2013;68(2):366-9.

581    35.    Grandjean L, Gilman RH, Iwamoto T, Koser CU, Coronel J, Zimic M, et al. Convergent

582    evolution and topologically disruptive polymorphisms among multidrug-resistant

583    tuberculosis in Peru. PLoS One. 2017;12(12):e0189838.

584    36.    Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, et al. Genomic

585    analysis identifies targets of convergent positive selection in drug-resistant Mycobacterium

586    tuberculosis. Nat Genet. 2013;45(10):1183-9.

587    37.    Phelan JE, Coll F, Bergval I, Anthony RM, Warren R, Sampson SL, et al. Recombination

588    in pe/ppe genes contributes to genetic variation in Mycobacterium tuberculosis lineages.

589    BMC Genomics. 2016;17:151.

590    38.    Ruesen C, Riza AL, Florescu A, Chaidir L, Editoiu C, Aalders N, et al. Linking minimum

591    inhibitory concentrations to whole genome sequence-predicted drug resistance in

592    Mycobacterium tuberculosis strains from Romania. Sci Rep. 2018;8(1):9676.

593    39.      Heym B, Alzari PM, Honore N, Cole ST. Missense mutations in the catalase-

594    peroxidase gene, katG, are associated with isoniazid resistance in Mycobacterium

595    tuberculosis. Mol Microbiol. 1995;15(2):235-45.

596    40.      Coll F, Phelan J, Hill-Cawthorne GA, Nair MB, Mallard K, Ali S, et al. Genome-wide

597    analysis of multi- and extensively drug-resistant Mycobacterium tuberculosis. Nat Genet.

598    2018;50(2):307-16.

599    41.      Larsen MH, Biermann K, Tandberg S, Hsu T, Jacobs WR. Genetic Manipulation of

600    Mycobacterium tuberculosis. Curr Protoc Microbiol. 2007;Chapter 10:Unit 10A.2.

601    42.      Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler

602    transform. Bioinformatics. 2009;25(14):1754-60.

603    43.      Garcia-Alcalde F, Okonechnikov K, Carbonell J, Cruz LM, Gotz S, Tarazona S, et al.

604    Qualimap: evaluating next-generation sequencing alignment data. Bioinformatics.

605    2012;28(20):2678-9.

606    44.      McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The

607    Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA

608    sequencing data. Genome Res. 2010;20(9):1297-303.

609    45.      Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using

610    exact alignments. Genome Biol. 2014;15(3):R46.

611    46.      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    47.      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    48.    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    49.    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