

1 **Title:**

2 Comprehensive and accurate genetic variant identification from contaminated and low
3 coverage *Mycobacterium tuberculosis* whole genome sequencing data.

4 **Authors:**

5 Tim H. Heupink¹, Lennert Verboven¹, Robin M. Warren², Annelies Van Rie¹.

6 **Affiliation:**

7 ¹ Family Medicine and Population Health (FAMPOP), Faculty of Medicine and Health
8 Sciences, University of Antwerp, Antwerp, Belgium.

9 ² South African Medical Research Council Centre for Tuberculosis Research/ DST/ NRF
10 Centre of Excellence for Biomedical Tuberculosis Research, Division of Molecular Biology
11 and Human Genetics, Stellenbosch University, Stellenbosch, South Africa.

12 **Corresponding author:**

13 Tim H. Heupink, tim.heupink@uantwerpen.be

14 **Keywords:**

15 *Mycobacterium tuberculosis*, whole-genome sequencing, sputum, low coverage,
16 contamination, bacteria

17 **Repositories:**

18 Simulated sequencing data have been deposited in SRA BioProject PRJNA706121.

19 **Abstract**

20 Improved understanding of the genomic variants that allow *Mycobacterium tuberculosis*
21 (*Mtb*) to acquire drug resistance, or tolerance, and increase its virulence are important
22 factors in controlling the current tuberculosis epidemic. Current approaches to *Mtb*
23 sequencing however cannot reveal *Mtb*'s full genomic diversity due to the strict
24 requirements of low contamination levels, high *Mtb* sequence coverage, and elimination of
25 complex regions.

26 We developed the XBS (compleX Bacterial Samples) bioinformatics pipeline which
27 implements joint calling and machine-learning-based variant filtering tools to specifically
28 improve variant detection in the important *Mtb* samples that do not meet these criteria,
29 such as those from unbiased sputum samples. Using novel simulated datasets, that permit
30 exact accuracy verification, XBS was compared to the UVP and MTBseq pipelines.
31 Accuracy statistics showed that all three pipelines performed equally well for sequence
32 data that resemble those obtained from high depth coverage and low-level contamination
33 culture isolates. In the complex genomic regions however, XBS accurately identified 9.0%
34 more single nucleotide polymorphisms and 8.1% more single nucleotide insertions and
35 deletions than the WHO-endorsed unified analysis variant pipeline. XBS also had superior
36 accuracy for sequence data that resemble those obtained directly from sputum samples,
37 where depth of coverage is typically very low and contamination levels are high. XBS was
38 the only pipeline not affected by low depth of coverage (5-10 \times), type of contamination and
39 excessive contamination levels (>50%). Simulation results were confirmed using WGS
40 data from clinical samples, confirming the superior performance of XBS with a higher
41 sensitivity (98.8%) when analysing culture isolates and identification of 13.9% more
42 variable sites in WGS data from sputum samples as compared to MTBseq, without
43 evidence for false positive variants when ribosomal RNA regions were excluded.

44 The XBS pipeline facilitates sequencing of less-than-perfect *Mtb* samples. These
45 advances will benefit future clinical applications of *Mtb* sequencing, especially whole
46 genome sequencing directly from clinical specimens, thereby avoiding *in vitro* biases and
47 making many more samples available for drug resistance and other genomic analyses.
48 The additional genetic resolution and increased sample success rate will improve genome-
49 wide association studies and sequence-based transmission studies.

50 **Impact statement**

51 *Mycobacterium tuberculosis* (*Mtb*) DNA is usually extracted from culture isolates to obtain
52 high quantities of non-contaminated DNA but this process can change the make-up of the
53 bacterial population and is time-consuming. Furthermore, current analytic approaches
54 exclude complex genomic regions where DNA sequences are repeated to avoid inference
55 of false positive genetic variants, which may result in the loss of important genetic
56 information.

57 We designed the compleX Bacterial Sample (XBS) variant caller to overcome these
58 limitations. XBS employs joint variant calling and machine-learning-based variant filtering
59 to ensure that high quality variants can be inferred from low coverage and highly
60 contaminated genomic sequence data obtained directly from sputum samples.

61 Simulation and clinical data analyses showed that XBS performs better than other
62 pipelines as it can identify more genetic variants and can handle complex (low depth,
63 highly contaminated) Mtb samples. The XBS pipeline was designed to analyse Mtb
64 samples but can easily be adapted to analyse other complex bacterial samples.

65 **Data summary**

66 Simulated sequencing data have been deposited in SRA BioProject PRJNA706121. All
67 detailed findings are available in the Supplementary Material. Scripts for running the XBS
68 variant calling core are available on <https://github.com/TimHHH/XBS>
69 The authors confirm all supporting data, code and protocols have been provided within the
70 article or through supplementary data files.

71 INTRODUCTION

72 Genetic approaches are increasingly used in tuberculosis research and for the diagnosis
73 of drug resistant tuberculosis. Whole genome sequencing (WGS) of *Mycobacterium*
74 *tuberculosis* (*Mtb*) aims to investigate the entire genome of the *Mtb* strain to
75 comprehensively assess all known drug resistance conferring regions, provide maximum
76 resolution for genetic transmission studies, and investigate the role of genomic variants
77 using genome wide association studies [1]. The three key problems facing the current *Mtb*
78 WGS approaches are the need for high quantities of *Mtb* DNA, presence of contaminant
79 bacterial and human DNA, and the presence of complex regions in the *Mtb* genome.

80 *Mtb* is notoriously difficult to sequence directly from clinical samples because the DNA
81 from human cells, bacteria and viruses outnumbers that from *Mtb* bacilli. This results in
82 insufficient template *Mtb* DNA and low genomic depth of coverage when sequenced [2].
83 *Mtb* WGS therefore primarily uses cultured isolates, which requires a harsh
84 decontamination step followed by a lengthy (2 to 4 weeks) incubation to generate high
85 quantities of *Mtb*. The decontamination step not only reduces the presence of bacteria
86 other than *Mtb*, but may also reduce the *Mtb* load [3]. The culture step can increase the
87 presence of certain strains over others due to stochastic processes or when certain strains
88 are better suited at growing in culture media [4,5]. These processing steps thus result in a
89 population bias, where the inferred *in vitro* *Mtb* population may not truly represent the *in*
90 *vivo* population. To generate a rapid and unbiased result, *Mtb* would thus ideally be
91 sequenced directly from the clinical sample.

92 Despite decontamination, a small proportion of contaminants may persist in the DNA
93 extracted from the culture isolate. Current *Mtb* bioinformatic pipelines use *in silico* meta-
94 genomic classification software to identify these contaminants and exclude samples with a
95 high proportion of contaminant DNA. For example the unified analysis variant pipeline
96 (UVP) uses a cut-off of maximum 10% contamination [6], which may exclude valuable
97 samples from analysis. In addition to the low contamination threshold, the *Mtb* community
98 has adopted relatively high standards for genome coverage, with 30× to 50× and up to
99 100× depth being the most commonly used. A Poisson distribution however reveals that a
100 mean depth of coverage of 15.8× results in 95% of the genome being covered by 10× or
101 more reads (Lander and Waterman 1988), which should be sufficient for accurately calling
102 majority variants in a haploid genome.

103 A third problem that complicates *Mtb* WGS is the abundance of complex regions including
104 repeats, transposons, duplicates and phage genes, and the numerous PE/PPE genes.
105 These complex regions are generally excluded from analyses by *Mtb* pipelines. The core
106 genome multi-locus sequence typing (cgMLST) method goes even further as only the most
107 trustable regions are analysed by cgMLST [7]. While these strategies ensure the accuracy
108 of the genome assembly and variant calling, they can result in the loss of a significant
109 proportion of genomic information (~9% when using the UVP pipeline [6]) that may be
110 important for defining transmission events and identification of variants that affect
111 pathogenicity.

112 There is thus a need for novel bioinformatics tools that overcome the current requirements
113 of low contamination, high *Mtb* DNA sequence coverage and exclusion of complex

114 genomic regions [8]. The Genome Analysis ToolKit (GATK), originally designed for human
115 genome studies (i.e. diploid), now allow for processing of haploid genomes such that of
116 *Mtb*. The Base Quality Score Recalibration and Indel Realigner tools and single sample
117 variant calling using the now superseded Unified Genotyper tool have already been
118 applied in *Mtb* genome studies [6,9]. The GATK's 'Germline short variant discovery Best
119 Practices workflow' however includes joint genotyping and machine-learning-based variant
120 filtering and has seen little to no implementation in bacterial and *Mtb* genome assembly
121 pipelines. The major advantage of joint variant calling, as opposed to single sample variant
122 calling, is a greater sensitivity for variants at low frequency in the population and detection
123 of variants in low coverage samples for which there would be insufficient confidence if the
124 sample had been analysed on its own [10]. Joint variant calling also enables the
125 calculation of various statistical annotations (including depth of coverage, strand bias and
126 read mapping quality) for alleles in the population rather than for those in a single strain.
127 These population variants are more numerous and their annotations suffer from fewer
128 stochastic deviations, thereby improving subsequent variant filtering. The machine-
129 learning-based variant filtering (VQSR) in the GATK [11] eliminates the need for hard-
130 filtering of variants, which is commonly applied through the use of rather arbitrary cut-offs
131 for strand bias and coverage depth.

132 We hypothesize that the GATK's tools are suitable for distinguishing contaminant variants
133 from *Mtb* and to score and identify variants in complex regions of the *Mtb* genome. We
134 developed a novel *Mtb* pipeline integrating the GATK's tools to improve the identification of
135 genetic variants in less-than-perfect *Mtb* samples and thereby greatly increase our power
136 to capture the diversity of within-patient and within-bacterial genetic information. We also
137 tested the variant calling core of this pipeline for its accuracy to identify genetic variants in
138 comparison to existing pipelines.

139 **METHODS**

140 **Development of the XBS pipeline**

141 The complex Bacterial Sample (XBS) pipeline was designed to perform analyses of
142 Illumina FASTQ sequence data. The pipeline was primarily designed to analyse *Mtb*
143 samples but can easily be adapted to analyse other complex bacterial samples. XBS was
144 realised through coupling published software packages with custom Bash and Python
145 scripts.

146 Pipelines typically start with identifying the level of contaminants and/or removing
147 contaminants before mapping the sequence reads. In XBS, all FASTQ sequence data,
148 whether single read or paired-end, are directly mapped to the reference genome (H37Rv:
149 NC_000962.3) however using BWA mem [12] (Figure 1). XBS does not employ an adapter
150 trimming step because BWA mem locally aligns sequence reads, which masks the portions
151 of the read that do not align well with the reference genome. Skipping the step of removing
152 contaminants saves considerable computing time but does require sophisticated
153 downstream variant filtering to distinguish genuine *Mtb* variants from those introduced by
154 contaminants.

155 Next, the mapped sequence library is merged with other independently mapped sequence
156 libraries from the same sample using Samtools (<https://www.htslib.org/>). The GATK
157 MarkDuplicates (Picard) is then used to mark duplicated reads in the merged bam file.
158 Unlike other *Mtb* pipelines XBS does not employ Base Quality Score Recalibration (BQSR)
159 to avoid that variants in contaminant DNA are interpreted as systematic error by BQSR
160 which would result in reduced base qualities, including for genuine *Mtb* variants.

161 The mapped sequences are then locally reassembled to correctly identify possible
162 haplotypes and their variants using the GATK HaplotypeCaller. At this point, the statistics
163 of depth of coverage, breadth of coverage, multiple infection and level of nontuberculous
164 Mycobacteria (NTM) contamination are assessed to judge if a sample is suitable for
165 subsequent joint variant calling. The coverage statistics are inferred using the GATK
166 CollectWgsMetrics (Picard). Quality approved samples' Genomic Variant Call Format
167 (GVCF) files are then merged with the GATK CombineGVCFs and the genotypes are joint
168 called using the GATK GenotypeGVCFs. This results in a VCF file with the unfiltered
169 variants for all quality approved samples. GATK is run with a ploidy of 1 for the variant
170 calling processes so that the allele with the highest confidence is identified as the allele
171 representing the haploid genotype for each variant site.

172 Next, the machine-learning-based variant filtering (VQSR) in the GATK is employed to
173 identify the likely true variants [11]. This step requires a truth set of variants known to occur
174 in *Mtb*, which can for example consist of DR conferring mutations. Single nucleotide
175 polymorphisms (SNPs) and insertions and deletions (INDELS) are processed separately
176 for variant filtering. The annotated statistics calculated during the genotyping are used to
177 build a positive statistical model for the variants in the dataset that also occur in the truth
178 set. Similarly, a negative variant model is built for the variants with the most inferior
179 annotated statistics. The remaining variants not consulted for either the positive or the
180 negative model construction are then confidence scored according to the placement of
181 their annotated statistics in relation to these models. To identify as many variants as
182 possible, variants are then filtered by applying a target sensitivity of 99.9%, calculated as
183 the percentage of identified variants from the present truth-set variants. The filtered SNP
184 and INDEL VCFs are then further processed as appropriate for constructing annotated
185 phylogenies and inferences of multiple infection, drug (hetero-)resistance, lineage and
186 transmission clusters.

187 **In silico development of a simulated dataset**

188 A dataset of 1,200 simulated samples representing 50,000 SNPs, 2,500 insertions and
189 2,500 deletions was developed (Figure 2). Of these, 600 were designed to resemble WGS
190 data from mono-culture isolates and 600 to resemble WGS data obtained directly from
191 sputum samples, the latter including high levels of various contaminants.

192 First, a set of 50 simulated strains with the exact mutations known in respect to the
193 reference genome was created *in silico* using SNP Mutator v.1.2.0 [13]. Each simulated
194 genome was created by randomly introducing 1000 SNPs, 50 single nucleotide insertions
195 and 50 single nucleotide deletions in H37Rv (NC_000962.3). Multi-nucleotide INDELS
196 could not be introduced using the SNP Mutator software. SNPs and INDELS were
197 introduced randomly throughout each simulated strain genome to ensure that random

198 bases were affected and/or introduced and to create genetically varying strains. The GATK
199 LeftAlignAndTrimVariants was used so that the truth set its INDELS were in the standard
200 notation.

201 For each dataset, 100 strains were randomly drawn from the 50 simulated strains to
202 ensure that some strains and their variants occurred more than once, as would be the
203 case for clinical datasets were specific strains and drug resistance variants often occur
204 more than once. To simulate WGS data obtained from culture isolates, 6 datasets
205 representing 5×, 10×, 20×, 30×, 50× or 100× depth were generated using a 100 randomly
206 drawn simulated strains each (Figure 2). In order to investigate the impact of low-level
207 contamination, no or low-level contamination (0, 1, 2, 3, 4, or 5%) and contamination type
208 was randomly assigned to each of the simulated strains. The eight contamination types
209 used were *Mycobacterium intracellulare* (NC_016946.1), *Mycobacterium abscessus*
210 (NC_010397.1, including plasmid), *Mycobacterium avium* (NC_002944.2), the three most
211 common NTM, *Pseudomonas aeruginosa* (NC_002516.2) and *Staphylococcus*
212 *epidermidis* (NC_004461.1), *Homo sapiens* (GRCh38), a mixture of the three NTM (*M.*
213 *intracellulare*, *M. abscessus* and *M. avium*) and a mixture of all six contaminants. Because it
214 was not possible to represent the full diverse spectrum of contaminating bacteria in the
215 simulations, *Pseudomonas aeruginosa* and *Staphylococcus epidermidis* were selected as
216 these are the most common bacterial contaminants [2], NTMs were included because
217 these pose a serious challenge for Mtb variant calling. The simulated samples included no
218 or low-level contamination in order to resemble WGS data from culture isolates and to be
219 able to investigate the effect of the various low levels of such contamination. Simulated
220 contaminant sequence reads were added to the simulated *Mtb* reads so that the final
221 contamination level matched the assigned contamination percentage.

222 ART v.2.5.8 software [14] was then used to emulate Illumina 150 bp paired-end sequence
223 reads with an HiSeq error profile and a Poisson distributed 300bp average library insert
224 size for each simulated strain and its contaminant(s). In total, 600 cultured WGS samples
225 were generated in six datasets of 100 simulated samples with each dataset representing a
226 different level of coverage (5×, 10×, 20×, 30×, 50× or 100× depth) to allow assessment of
227 the relation between coverage and accuracy of variant identification.

228 To simulate *Mtb* WGS data obtained directly from sputum samples, another six datasets
229 were generated, each with a set number of paired-end sequence reads per sample,
230 ranging from 500,000 to 3,000,000 PE reads (Figure 2). *Mtb* and contamination levels
231 were randomly sampled from a beta distribution around 0.01 to 78.63% *Mtb* DNA to match
232 levels observed for direct-from-sputum WGS data [8]. The contamination type was either
233 *P. aeruginosa*, *S. epidermidis*, *H. sapiens*, a mixture of *P. aeruginosa*, *S. epidermidis* and
234 *H. sapiens* or a NTM mixture (*M. intracellulare*, *M. abscessus* or *M. avium* up to 20% of the
235 *Mtb* fraction with the remaining contamination consisting of *P. aeruginosa*, *S. epidermidis*
236 or *H. sapiens*).

237 ART v.2.5.8 software [14] was used as described previously to emulate sequence reads
238 for each simulated strain and its contaminant(s). In total, 600 simulated WGS data directly
239 from sputum samples were generated in six datasets of 100 simulated samples, each
240 dataset representing a number of paired-end sequence reads and had various levels and

241 types of contamination, allowing the study of the relation between contaminant nature,
242 read number and variant identification.

243

244 **Assessment of XBS pipeline performance**

245 The performance of XBS was compared to UVP [6] and MTBseq [15], two well-established
246 and commonly used *Mtb* pipelines . Each pipeline was evaluated using their standard
247 settings. For UVP, variants in the GATK filtered VCF file were used for accuracy
248 calculations. For MTBseq, two approaches were assessed. In 'MTBseq-basic', the 'GATK
249 position variants' file was used for accuracy calculations. In 'MTBseq-exrep', the variants
250 marked 'repetitive' in MTBseq's 'MTB_Gene_Categories.txt' were excluded to assess the
251 effect of this commonly applied filtering step. From here on XBS, UVP and the two
252 MTBseq approaches will be referred to as pipelines.

253 For XBS, VQSR was run with a truth set consisting of 5,000 SNPs, 250 insertions and 250
254 deletions randomly selected from the mutations known to occur in the 50 simulated strains.
255 The truth set therefore represented 10% (5,500/ 55,000) of the total variants introduced *in*
256 *silico*. The GATK VariantRecalibrator was run to score each SNP according to the inferred
257 positive and negative models, built on the depth, mapping quality, mapping quality rank-
258 sum and quality by depth statistical annotations. These annotations were processed in an
259 allele specific fashion to distinguish between genuine and contaminant variants occurring
260 on the same genomic location. Annotations that showed insufficient variance, as
261 determined by VQSR, were excluded. A logit transform and jitter were applied to improve
262 mapping quality-based filtering. The FS, ReadPosRanksum and SOR annotations were
263 excluded because they are more applicable for real sequence than for simulated data.
264 Where possible, four Gaussians were used for the positive model. The GATK ApplyVQSR
265 was applied with a truth sensitivity level of 99.9%. The same process was followed for
266 INDELS except that allele specificity was disabled and, where possible, two Gaussians
267 were used for the positive model. The MQ annotation was not taken into consideration
268 following the GATK Best Practices Workflow for Germline short variant discovery. To avoid
269 over-representation of contaminant alleles in the filtered dataset, INDELS in the sputum
270 simulations were filtered using a VQSLOD score of 0 rather than a truth sensitivity level.
271 This ensures that only those INDELS that are most likely to fall in the positive and not the
272 negative model are kept. The resulting SNP and INDEL variants were used for the
273 accuracy calculations. Version 4.1.9.0 of the GATK was used.

274 The six datasets simulating WGS from culture were analysed using all four pipelines (XBS,
275 MTBseq-basic, MTBseq-exrep, UVP). The six datasets simulating WGS from sputum
276 samples only by three pipelines as UVP can not analyse samples with contamination
277 exceeding 10%. The inferred variants identified by each pipeline were compared with the
278 truth in terms of genome position and allelic nature (bases involved and length). The
279 pipeline's accuracy was calculated in terms of precision, recall and their harmonic mean
280 (F_1 score).

281 For simulation of WGS from culture isolates, accuracy scores were averaged over the 100
282 samples in each dataset and calculated for each combination of variant type (SNP or
283 INDEL) and genomic region (complete, complex or non-complex). F_1 scores were plotted

284 for the four pipelines at six levels of depth for the complete genome, and for complex and
285 non-complex regions separately. UVP's list of excluded loci was used to define regions as
286 complex, these were filtered using the GATK SelectVariants.

287 For simulation of WGS from sputum, the range of F_1 scores was calculated, separately for
288 SNPs and INDELS, and the proportion of samples with an F_1 score >0.9 was estimated.
289 The number of *Mtb* reads was converted to the theoretical depth of coverage (number of
290 simulated *Mtb* reads multiplied by average read length) and plotted against F_1 scores for
291 each pipeline, contamination level and type after excluding samples with a theoretical
292 coverage of $<20\times$ so that the lesser performance of such samples did not distort these
293 figures.

294 Plots were created in R using the ggplot2 and gridExtra packages.

295 **Analysis of WGS from clinical samples**

296 The performance of the XBS variant caller was examined using two published WGS
297 datasets obtained from clinical samples and compared to UVP and MTBseq. Data
298 published by Roetzer et al. [16] was used to test the pipelines' sensitivity by identifying a
299 set of known variants (Sanger confirmed) in DNA extracted from cultured *Mtb* samples.
300 Data from Goig et al. [8] was used to evaluate the ability to call variants in WGS data from
301 DNA extracted directly from sputum. Only samples with $\geq 5\times$ genomic coverage depth
302 (S02, S26, S17, S21, S31, S20, S27, S67, S09 and S69) were analysed to ensure
303 sufficient width of coverage and prevent problems with phylogenetic inference. A reference
304 set of 125 diverse cultured *Mtb* strains with high coverage WGS data was included in the
305 analyses to provide a reference in the phylogenetic tree and increase the variation in
306 statistical annotations, thereby improving VQSR for XBS. XBS VQSR was run in SNP
307 mode with a truth-set containing lineage and DR variants [17–19]. These variants reflect
308 the diversity of the bacterium (*Mtb* lineages) and the entirety of the genome, enabling
309 VQSR to build a model to identify variants in exactly such regions and as such avoid bias.
310 The variants from the Goig et al. and reference dataset were converted to FASTA format,
311 where positions represented by fewer than 95% of the samples were excluded. MTBseq
312 and UVP were run in default mode.

313 IQ-TREE v2.1.2 was used to construct the Maximum Likelihood trees [20] which was
314 plotted with Figtree v1.4.3 [21] and the resulting branch lengths were used to evaluate the
315 potential presence of false positive variants.

316 **RESULTS**

317 **Pipeline performance for analysis of WGS data from simulated culture** 318 **isolates**

319 At the highest coverage ($100\times$) and with the low levels of contamination ($\leq 5\%$), all pipeline
320 approaches detected very few false positives and missed few variants, resulting in a 100%
321 precision for SNPs and INDELS across the genome, except for UVP which obtained a
322 slightly lower precision of 98% for SNP calling (Tables 1 and 2). Recall scores were
323 highest for XBS and MTBseq-basic (97-99% for the SNPs and INDELS) compared to
324 MTBseq-exrep and UVP which missed some variants (92% for SNPs and 91% for INDELS

325 by MTBseq-exrep; 90% for SNPs and 89% for INDELS by UVP). The overall variant calling
326 accuracy was highest for XBS and MTBseq-basic (F_1 score 0.99 for SNPs, ≥ 0.98 for
327 INDELS), somewhat lower for MTBseq-exrep (F_1 score 0.96 for SNPs, 0.95 for INDELS),
328 and lowest for UVP (F_1 score 0.94 for SNPs and INDELS) (Figure 3). At 100 \times coverage
329 MTBseq-basic identified an average 9.2% more true positive SNPs and 9.8% more
330 INDELS per genome when compared to UVP (Table 1 and 2). XBS identified an average
331 9.0% more true positive SNPs and 8.1% more INDELS per genome when compared to
332 UVP at 100 \times coverage.

333 Lowering the depth of *Mtb* genome coverage from 100 \times to 20 \times had minimal effect on
334 accuracy scores. XBS's precision and recal did not change and the F_1 score deviated by \leq
335 0.01 for SNPs and INDELS. The performance of MTBseq-basic and MTBseq-exrep also
336 did not differ for these lower coverages with similar precision for SNPs and INDELS, a drop
337 in recall by 2% for SNPs and 3% for INDELS, and the F_1 score lowered by 0.01 for both
338 SNPs and INDELS. UVP's accuracy statistics did not change for INDELS, but precision,
339 recall and F_1 score for SNPs dropped slightly by 4%, 1% and 0.01, respectively.

340 At depths ranging from 20 \times to 100 \times , the performance for SNP and INDEL calling in the
341 non-complex regions of the *Mtb* genome was high for all four pipelines (Figure 4 and
342 Supplementary Table 1). Performance for variant calling in the complex regions was
343 similar for MTBseq-basic and XBS, with an average SNP and INDEL precision of 100%,
344 recall around 91% and the F_1 around 0.95. Accuracy statistics for variant calling in the
345 complex regions could not be calculate for UVP as complex regions are excluded from
346 from its standard output. MTBseq-exrep's exclusion of repetitive regions was less strict
347 than UVP's excluded loci and hence the former was able to identify a small number of
348 variants in the complex regions.

349 At 10 \times depth of coverage, all four pipelines retained their precision but recall was affected.
350 UVP and both MTBseq approaches missed many variants resulting in a recall of around
351 50% for SNPs and INDELS. Consequently, F_1 scores dropped drastically for MTBseq-basic
352 (0.69 for SNPs and 0.66 for INDELS), MTBseq-exrep (0.66 for SNPs and 0.63 for INDELS)
353 and UVP (0.64 for both SNPs and INDELS). In contrast, XBS's accuracy remained high,
354 with recall at 99% for SNPs and 98% for INDELS and F_1 scores of 0.99 for SNPs and
355 INDELS (Figure 3). Only at an *Mtb* genome coverage of 5 \times was the performance of XBS
356 noticeably affected, although performance remained largely accurate with F_1 scores of
357 0.96 for SNPs and 0.95 for INDELS, a precision of 100% for both, and recall of 93% for
358 SNPs and 91% for INDEL calling. XBS's ability to call variants in both low coverage and
359 complex regions was retained (Figure 4).

360 The type of low-level contaminant (*M. intracellulare*, *M. abscessus*, *M. avium*, *P.*
361 *aeruginosa*, *S. epidermidis*, *H. sapiens*, NTM mixture or mixture of all 6 contaminants) only
362 affected the F_1 estimates of UVP due to false positive SNP calls when NTMs or *S.*
363 *epidermidis* were present (Supplementary Figure 1 and Supplementary Table 1). The level
364 of contamination (varying from 0-5%) also only affected UVP's performance, with a
365 decrease in precision and SNP F_1 scores at higher levels of contamination (Supplementary
366 Figure 2 and Supplementary Table 1).

367 Pipeline performance for analysis of WGS data from simulated sputum 368 samples

369 XBS outperformed both MTBseq approaches when analysing the data simulated to
370 represent WGS directly from sputum. For SNPs, F_1 scores ranged from 0.63-0.84 for XBS
371 compared to 0.33-0.58 for MTBseq-basic and 0.31-0.59 for MTBseq-exrep. Using XBS, 49
372 to 77% of samples achieved F_1 scores above 0.90, compared to 20 to 53% and 15 to 45%
373 for MTBseq-basic and MTBseq-exrep, respectively. (Table 1 and 2). For INDELS F_1 scores
374 ranged from 0.61-0.81 for XBS, 0.32-0.63 for MTBseq-basic and 0.31-0.60 for MTBseq-
375 exrep. Using XBS, 47 to 73 % of samples achieved F_1 scores above 0.9, compared to 21
376 up to 58% and 16 up to 53% for MTBseq-basic and MTBseq-exrep, respectively. Plotting
377 the theoretical depth of coverage against F_1 score showed that XBS calls SNPs and
378 INDELS with higher accuracy at low genomic depth of coverage compared to both MTBseq
379 approaches (Figure 5).

380 SNP accuracy of XBS was unaffected by contamination level or type (Figure 6). In
381 contrast, accuracy for MTBseq-basic and MTBseq-exrep depended on type and level of
382 contamination. *H. sapiens* contamination did not affect the F_1 score, *S. epidermidis*
383 lowered the F_1 score to 0.90 when the contamination level was $\geq 50\%$, and NTM and
384 bacterial/human contamination mixtures reduced the F_1 score when the contamination
385 level was $\geq 75\%$. For INDELS, the MTBseq-basic pipeline performed slightly better than
386 XBS when *Mtb* depth of coverage was $\geq 20\times$, with average F_1 scores of 0.99 for MTBseq-
387 basic, 0.95 for MTBseq-exrep and 0.96 for XBS respectively. (Supplemental Figure 3,
388 Figure 5B).

389 Pipeline performance for analysis of WGS data from clinical culture 390 isolates

391 MTBseq and XBS could analyse all samples from the Roetzer et al. dataset, whereas UVP
392 excluded 33 of the 86 (38%) of samples. Of the 85 Sanger confirmed mutations, MTBseq-
393 basic recovered 81, MTBseq-exrep 79, UVP 61 and XBS 84, corresponding to sensitivities
394 of 95.3, 92.9, 71.8 and 98.8% respectively (Supplementary Table 3). The single variant
395 missed by XBS was located right on the border of a repetitive region, resulting in reads
396 with sub-optimal mapping qualities.

397 Pipeline performance for analysis of WGS data from clinical sputum 398 samples

399 UVP failed to analyse any sample included in the Goig et al. dataset as the contamination
400 levels was above the 10% threshold for all samples. For the 10 Goig et al. samples and
401 the 125 reference samples, XBS reported 11,977 variant positions (after exclusion of the
402 ribosomal RNA regions), 13.9% more than the 10,514 variants reported by MTBseq. The
403 number of variants called by MTBseq further reduced to 10,114 when variants within 12bp
404 of each other were excluded.

405 There was no evidence of false positive variants when using XBS (no obvious branch
406 extension for any sputum samples) with the highly conserved ribosomal RNA regions
407 removed (Figure 7). When including the genes coding for ribosomal RNA obvious
408 extended branch lengths were present for three samples (S02, S26 and S20,

409 Supplemental Figure 4) due to VQSLOD scores for such variants that had fallen just within
410 the within the positive VQSR model. When using MTBseq, there were also no obvious
411 branch extensions but one sample (S26) showed a shorter branch length compared to its
412 nearest-neighbours in the phylogenetic tree. This was the case for FASTA files in- and
413 excluding SNPs within 12bp distance from each other (Supplementary Figures 5 and 6).

414

415 **DISCUSSION**

416 We developed XBS and applied the joint variant calling and machine-learning-based
417 variant filtering approaches, initially designed for human genome analyses, to a pipeline
418 for *Mtb* WGS analyses. Using 1,200 simulated samples representing characteristics of
419 WGS data from *Mtb* culture or directly from sputum samples, we demonstrated that XBS
420 increases the performance in variant calling compared to existing pipelines (UVP and
421 MTBseq), especially for WGS data from less-than-perfect, contaminated low *Mtb* burden
422 samples. The strain simulation, variant calling and filtering approach presented here may
423 also benefit the study of other bacteria where sequence coverage, complex genomes or
424 contamination hinder accurate genetic variant identification.

425 We showed that current pipeline approaches perform well for SNP and INDEL calling
426 when sequencing DNA extracts from decontaminated cultures with high ($\geq 20\times$) depth, but
427 accuracy decreases when depth of coverage is low (5-10 \times) or contamination levels are
428 high (>50%). The novel XBS pipeline substantially outperformed other MTB pipelines for
429 SNP and INDEL calling in *Mtb* at low coverage depth culture samples and highly
430 contaminated, low coverage depth sputum samples.

431 When analysing WGS data from culture isolates at the current standard 30 to 100 \times depth
432 coverage, all pipelines accurately called SNP and INDEL (F_1 scores >0.90). Of the three
433 pipelines assessed, UVP's performance was slightly inferior given its lower precision (false
434 positives variants) at higher (5%) contamination, particularly when the contaminant was an
435 NTM. XBS and MTBseq-basic were not affected by low level (0-5%) contamination levels
436 and identified on average 9% more SNPs and INDELS compared to UVP by investigating
437 *Mtb*'s complex genomic regions. Identifying 9% more variants could greatly benefit
438 transmission and genome-wide association studies.

439 At lower coverage (<20 \times), XBS was the only pipeline that could accurately call SNPs and
440 INDELS, likely due to the joint calling and filtering processes that permit lower allele
441 coverages. XBS's accuracy remained high at 5 \times depth, where the modest drop in F_1 score
442 was due to a slightly lower recall rate and difficulty in accurately calling genuine INDELS
443 due to coverage gaps. This is expected as, according to the Poisson distribution, only
444 99.3% of the genome is covered by at least one read at 5 \times coverage. The low-level
445 contamination simulated for the culture samples did not affect the accuracy of the XBS or
446 MTBseq pipelines. The UVP pipeline was however affected by both the level and the type
447 of low-level contamination, such effects have been observed previously [22]. Considering
448 these findings it is understandable that UVP uses a strict contamination cut-off, but the
449 other pipelines show that variants can be identified more accurately despite the absence of
450 such cut-offs. Sequencing at 5 \times depth using XBS resulted in average SNP F_1 score of 0.96
451 (minimum 0.95) and average INDEL F_1 scores of 0.95 (minimum 0.91), whereas the F_1

452 scores of the other pipelines were ≤ 0.10 for both SNPs and INDELS. Such low-coverage
453 sequencing could lower the costs by a factor of 10 compared to standard 50 \times coverage
454 sequencing. In combination with low-cost library preparations, which is the main driver of
455 sequencing cost, this could open the door to large-scale population sequencing projects in
456 high TB-burden settings.

457 WGS data obtained directly from sputum is characterized by a low number of *Mtb* reads
458 (theoretical coverage), a high the level of contamination, and presence of a mix of
459 contaminants. The novel XBS pipeline showed superior performance for analysing such
460 impure sequencing data. Due to the joint calling approach, XBS could analyse samples
461 with much lower genomic coverage than the two MTBseq approaches. XBS successfully
462 identified SNPs and INDELS in an average 73% of samples with 2,5 to 5 million paired-end
463 reads, where MTBseq-basic only successfully analysed 50% and MTBseq-exrep 45% of
464 such samples. By employing VQSR filtering, which identifies contaminant reads based on
465 a multitude of statistical annotations, XBS's performance was not affected by level or type
466 of contaminants. Hard filtering, as is implemented in MTBseq-basic and -exrep, was not
467 sufficient at high levels (>50%) of contamination because contaminant alleles may be
468 interpreted as the most likely and therefore genuine allele, leading to false positives, once
469 they reach coverage levels greater than the *Mtb* allele. For MTBseq, the type of
470 contaminant affected the accuracy. High levels of human DNA did not affect accuracy as
471 these are unlikely to map to the reference genome, but *S. epidermidis* contamination
472 started to have an effect from 50% upwards as contaminant alleles then outnumber that of
473 *Mtb*. The NTM-mix only affected accuracy at high contamination despite NTMs great
474 genome similarity. This counterintuitive finding is likely because high levels of
475 contamination are required before one of the NTM contaminants present in the mix
476 approach 50% allele frequency.

477 Since it cannot be said exactly which samples underperform in terms of variant
478 identification accuracy in a clinical dataset, it is best to ensure an acceptable minimum
479 accuracy instead. When employing XBS on WGS data from real-life sputum, our data
480 suggests that it may be prudent to restrict the analyses to those samples that present a
481 coverage of $\geq 10\times$. A 10 \times cut-off would result in average SNP F_1 scores of 0.99, minimum
482 0.98, and average INDEL F_1 scores of 0.97, minimum 0.91, for 60% (60/100) of the
483 samples in the 3,000,000 PE read dataset. MTBseq-basic would result in average SNP F_1
484 scores of 0.91, minimum 0.53, and average INDEL F_1 scores of 0.95, minimum 0.69, for
485 the same samples.

486 Comparisons of the performance of *Mtb* pipelines is important but hampered by the
487 absence of large datasets for which the true variants are known. To date, studies
488 assessing the performance of *Mtb* pipelines have compared pipelines' ability to identify
489 transmission clusters as established through contact tracing or older molecular methods,
490 or by comparing the detection of genomic drug resistance in relation to phenotypic tests or
491 Sanger-confirmed variants [6,15,23,24]. These approaches suffer from important
492 limitations. Contact tracing is complex and may not necessarily identify all clusters
493 correctly [25]. Older molecular methods have significantly lower resolution than WGS so
494 that all pipelines call clusters identified by these older methods with relative ease [26].
495 Using genomic drug resistance to compare pipelines is affected by the reference drug

496 resistance mutation list used by each pipeline, whereas focussing on a limited number of
497 Sanger-confirmed variants is not representative for the entire genome. The only other
498 study that used simulated read datasets to compare combinations of mapper, caller and
499 filtering methods found that the GATK variant caller in combination with VQSR consistently
500 had the highest precision scores [27], supporting the findings of our study. This study
501 however had multiple limitations. First, the GATK calling was performed for one sample at
502 a time, which is not optimal for VQSR or low coverage samples. Second, for the VQSR
503 truth sets, half of the samples' variants with the best quality score were taken, an approach
504 is problematic when high frequency contaminant alleles are present. Third, the use of the
505 clinical CDC1551 strain prohibited accuracy assessment for the complex regions as the
506 exact location of CDC1551 variants in relation to H37Rv cannot be established with
507 certainty for complex regions.

508 We successfully overcame the limitations of prior *Mtb* pipeline accuracy studies by using
509 an *in silico* approach to construct a fully representative variant truth set. The simulated
510 dataset resembled clinical datasets by ensuring that some strains occurred once while
511 others occurred several times. VQSR benefited from the presence of clonal strains as it
512 improves the identification of variants in low coverage samples observed in other samples.
513 Our datasets simulating culture isolates represented the range of depth (5 to 100×) and
514 low levels of contamination (0-5%). Our simulated sputum datasets contained more than
515 21.4% contamination (mean 82.12%, maximum 99.99%) and thus very low levels of *Mtb*
516 DNA, which correspond to the findings of a recent study of the clinical samples where
517 most (51%) WGS data showed less than 5% *Mtb* DNA sequence reads [8]. The use of the
518 simulated datasets allowed us not only to accurately quantify the performance of different
519 pipelines for variant calling throughout the entire genome, including the complex regions,
520 but also assess the effect of important characteristics that determine accuracy such as
521 mycobacterial burden, level and type of contamination.

522 The excellent performance of XBS for the analysis of complex samples was confirmed
523 when analysing WGS data obtained from clinical samples. The analyses of the WGS data
524 from clinical culture isolates showed that XBS outperformed other pipelines (including
525 pipelines not investigated in this study) in terms of sensitivity (Supplementary Table 3). The
526 high specificity of XBS matches the findings of the culture simulations where the four
527 pipeline approaches show similar recall scores. The analysis of WGS data obtained from
528 DNA extracted directly from sputum samples confirmed that only XBS and MTBseq, but
529 not UVP, could successfully analyse such data. While both pipeline showed high specificity
530 (no evidence of false positive variants resulting in branch extension on phylogeny), the
531 performance of XBS was superior to MTBseq as it allowed the identification of 13.9% more
532 variants.

533 Several limitations remain to the novel XBS variant caller. First, XBS (and other pipelines)
534 cannot analyse samples when NTM contaminant sequences exceed 20%. Such samples
535 would require analyses by programs such as QuantTB that can potentially filter out NTM
536 contaminants before *Mtb* variant identification as they resemble multiple infections [28].
537 Second, XBS requires multiple samples for each run. Previously inferred Genomic VCF's
538 can however be included from the Combine VCF step just before the VQSR, eliminating
539 the need to batch new samples. Third, the highly conserved ribosomal RNA regions had to

540 be excluded for optimal specificity as sequences in these regions from contaminating
541 bacteria can map to the *Mtb* reference genome with very high confidence, making the
542 variants in such regions indistinguishable in terms of the statistical annotations used by
543 VQSR. Eliminating these regions may result in some loss of the genomic information,
544 albeit small as this region represent only 0.1% of the genome Fourth, we used H37Rv as
545 the reference genome for all analysis. If the *Mtb* ancestral genome would be used as the
546 reference genome instead a VQSR truth-set could be constructed by aligning the ancestral
547 genome to H37Rv and translating the latter its lineage and DR truth variants. When
548 applying the XBS approach to other bacteria, one should employ well-established variants
549 that occur throughout the entire genome for constructing the VQSR truth set. Finally, we
550 were not able to compare the run time of XBS to the other pipelines as there were
551 important differences in the analyses run other than the variant calling core and the
552 number of samples analysed differed due to pipeline restrictions. However, elimination of
553 the adapter removal and base recalibration steps reduces the overall processing time and
554 exclusion of meta-genomic classification further reduces computing time. It also prevents
555 the need for computing infrastructure with large memory requirements.

556 Direct-from-sputum WGS data contains a wealth of diverse bacterial contaminants besides
557 that of human origin and all these contaminants occur at widely varying levels. It was not
558 possible to represent this endless variety of bacterial contaminants in our simulation
559 experiments hence the most commonly observed bacterial contaminants where used
560 instead [2], NTMs were included because these pose a serious challenge for *Mtb* variant
561 calling. Contaminant levels were simulated to match levels observed for genuine sputum
562 samples [8]. As such it was possible to study the effect of the various contaminants and
563 their levels on variant calling, this would not have been possible had all the endless
564 contaminants observed for direct-from-sputum samples been used. To show that XBS was
565 able to handle such genuinely diverse contamination two clinical WGS datasets were
566 studied, one from cultured samples and one direct-from-sputum. Further studies are
567 required to study the effect of the full diversity of contaminants that are observed for direct-
568 from-sputum samples.

569 In conclusion, all pipelines studied (MTBseq, UVP, XBS) accurately analysed WGS data
570 from *Mtb* culture isolates. Only XBS and MTBseq could accurately identify variants in low
571 *Mtb* coverage and highly contaminated samples and XBS achieved higher performance
572 parameters of all pipelines studies. High performance at low depth could decrease
573 sequencing cost and improve WGS analysis directly from sputum samples. The accurate
574 identification of variants in the complex genomic *Mtb* regions allow for improved resolution
575 in transmission studies through increased genetic resolution and creates the ability to
576 explore the functional role of variants in these complex regions. Taken together, the novel
577 XBS pipeline sets the stage for the next generation of *Mtb* WGS studies.

578 **Authors and contributors**

579 THH, RMW and AVR conceptualised the project and methodology. THH and LV curated
580 the data and designed and scripted the software. RMW and AVR acquired the funding.
581 THH, LV, RMW and AVR wrote, reviewed and edited the manuscript.

582 **Conflicts of interest**

583 The authors declare that there are no conflicts of interest.

584 **Funding information**

585 This work was supported by the Research Foundation Flanders (FWO) [grant number
586 FWO Odysseus G0F8316N].

587 **Acknowledgements**

588 We would like to thank the members of the Tuberculosis Omics ResearCH (TORCH)
589 consortium for helpful discussions and particularly Elise De Vos for discussions
590 surrounding sputum simulation. We thank the reviewers for helpful comments and Galo A.
591 Goig et al. for providing additional direct-from-sputum sequencing data.

592 **References**

593

- 594 1. Meehan CJ, Goig GA, Kohl TA, et al. Whole genome sequencing of *Mycobacterium*
595 tuberculosis: current standards and open issues. *Nat. Rev. Microbiol.* 2019; 17:533–545
- 596 2. McClean M, Stanley T, Stanley S, et al. Identification and characterization of
597 breakthrough contaminants associated with the conventional isolation of *Mycobacterium*
598 tuberculosis. *J. Med. Microbiol.* 2011; 60:1292–1298
- 599 3. Rachow A, Saathoff E, Mtafya B, et al. The impact of repeated NALC/NaOH-
600 decontamination on the performance of Xpert MTB/RIF assay. *Tuberculosis* 2018; 110:56–
601 58
- 602 4. Farmanfarmaei G, Kamakoli MK, Sadegh HR, et al. Bias in detection of *Mycobacterium*
603 tuberculosis polyclonal infection: Use clinical samples or cultures? *Mol. Cell. Probes* 2017;
604 33:1–3
- 605 5. Nimmo C, Shaw LP, Doyle R, et al. Whole genome sequencing *Mycobacterium*
606 tuberculosis directly from sputum identifies more genetic diversity than sequencing from
607 culture. *BMC Genomics* 2019; 20:389
- 608 6. Ezewudo M, Borens A, Chiner-Oms Á, et al. Integrating standardized whole genome
609 sequence analysis with a global *Mycobacterium tuberculosis* antibiotic resistance
610 knowledgebase. *Sci. Rep.* 2018; 8:1–10
- 611 7. Kohl TA, Diel R, Harmsen D, et al. Whole-genome-based *Mycobacterium tuberculosis*
612 surveillance: a standardized, portable, and expandable approach. *J. Clin. Microbiol.* 2014;
613 52:2479–2486
- 614 8. Goig GA, Cancino-Muñoz I, Torres-Puente M, et al. Whole-genome sequencing of
615 *Mycobacterium tuberculosis* directly from clinical samples for high-resolution genomic
616 epidemiology and drug resistance surveillance: an observational study. *The Lancet*
617 *Microbe* 2020; 1:e175–e183
- 618 9. Kato-Maeda M, Ho C, Passarelli B, et al. Use of whole genome sequencing to
619 determine the microevolution of *Mycobacterium tuberculosis* during an outbreak. *PLoS*
620 *One* 2013; 8:e58235
- 621 10. Poplin R, Ruano-Rubio V, DePristo MA, et al. Scaling accurate genetic variant
622 discovery to tens of thousands of samples. *BioRxiv* 2017; 201178
- 623 11. DePristo MA, Banks E, Poplin R, et al. A framework for variation discovery and
624 genotyping using next-generation DNA sequencing data. *Nat. Genet.* 2011; 43:491
- 625 12. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
626 MEM. *arXiv Prepr. arXiv1303.3997* 2013;

- 627 13. Davis S, Pettengill JB, Luo Y, et al. CFSAN SNP Pipeline: an automated method for
628 constructing SNP matrices from next-generation sequence data. *PeerJ Comput. Sci.* 2015;
629 1:e20
- 630 14. Huang W, Li L, Myers JR, et al. ART: a next-generation sequencing read simulator.
631 *Bioinformatics* 2012; 28:593–594
- 632 15. Kohl TA, Utpatel C, Schleusener V, et al. MTBseq: a comprehensive pipeline for whole
633 genome sequence analysis of *Mycobacterium tuberculosis* complex isolates. *PeerJ* 2018;
634 6:e5895
- 635 16. Roetzer A, Diel R, Kohl TA, et al. Whole genome sequencing versus traditional
636 genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a longitudinal
637 molecular epidemiological study. *PLoS Med* 2013; 10:e1001387
- 638 17. Coll F, McNerney R, Guerra-Assunção JA, et al. A robust SNP barcode for typing
639 *Mycobacterium tuberculosis* complex strains. *Nat. Commun.* 2014; 5:1–5
- 640 18. Coll F, Phelan J, Hill-Cawthorne GA, et al. Genome-wide analysis of multi-and
641 extensively drug-resistant *Mycobacterium tuberculosis*. *Nat. Genet.* 2018; 50:307–316
- 642 19. Napier G, Campino S, Merid Y, et al. Robust barcoding and identification of
643 *Mycobacterium tuberculosis* lineages for epidemiological and clinical studies. *Genome*
644 *Med.* 2020; 12:1–10
- 645 20. Minh BQ, Schmidt HA, Chernomor O, et al. IQ-TREE 2: new models and efficient
646 methods for phylogenetic inference in the genomic era. *Mol. Biol. Evol.* 2020; 37:1530–
647 1534
- 648 21. Rambaut A. FigTree v.1.4.3. <https://github.com/rambaut/figtree/>
- 649 22. Goig GA, Blanco S, Garcia-Basteiro AL, et al. Contaminant DNA in bacterial
650 sequencing experiments is a major source of false genetic variability. *BMC Biol.* 2020;
651 18:1–15
- 652 23. Schleusener V, Köser CU, Beckert P, et al. *Mycobacterium tuberculosis* resistance
653 prediction and lineage classification from genome sequencing: Comparison of automated
654 analysis tools. *Sci. Rep.* 2017; 7:1–9
- 655 24. Jajou R, Kohl TA, Walker T, et al. Towards standardisation: Comparison of five whole
656 genome sequencing (WGS) analysis pipelines for detection of epidemiologically linked
657 tuberculosis cases. *Eurosurveillance* 2019; 24:
- 658 25. Nikolayevskyy V, Kranzer K, Niemann S, et al. Whole genome sequencing of
659 *Mycobacterium tuberculosis* for detection of recent transmission and tracing outbreaks: a
660 systematic review. *Tuberculosis* 2016; 98:77–85
- 661 26. Meehan CJ, Moris P, Kohl TA, et al. The relationship between transmission time and
662 clustering methods in *Mycobacterium tuberculosis* epidemiology. *EBioMedicine* 2018;
663 37:410–416

- 664 27. Walter KS, Colijn C, Cohen T, et al. Genomic variant-identification methods may alter
665 Mycobacterium tuberculosis transmission inferences. *Microb. Genomics* 2020;
666 6:mgen000418
- 667 28. Anyansi C, Keo A, Walker BJ, et al. QuantTB--A method to classify mixed
668 Mycobacterium tuberculosis infections within whole genome sequencing data. *BMC*
669 *Genomics* 2020; 21:80
670

671 **Figures and tables**

672 **Table 1: SNP calling accuracies across the entire genome for four Mtb pipelines.**

673 **Table 2: INDEL calling accuracies across the entire genome for four Mtb pipelines.**

674 **Figure 1: Flow chart for XBS's variant calling core.**

675 **Figure 2: Flow chart for dataset construction.**

676 **Figure 3: Performance (F_1 scores) of four bioinformatic pipelines for SNP calling in simulated *Mtb* culture isolates at six levels of depth of *Mtb* genomic coverage.**

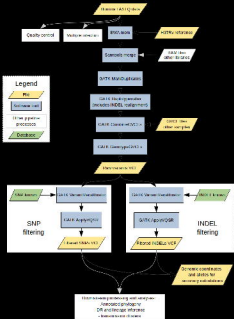
677
678 **Figure 4: Performance (F_1 scores) of four bioinformatic pipelines for SNP calling in simulated *Mtb* culture isolates at six levels of depth of *Mtb* genomic coverage, stratified by complex and non-complex regions of the genome.**

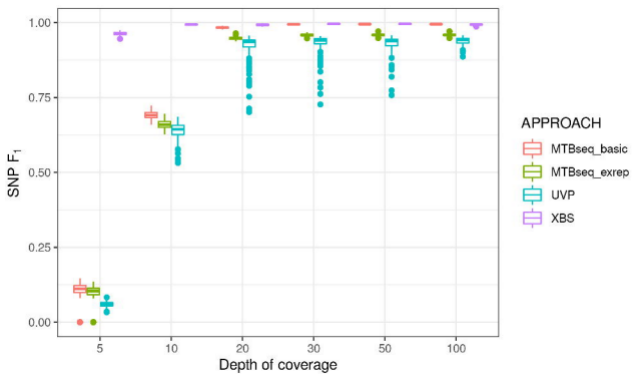
680
681 **Figure 5: Performance (F_1 scores) of different bioinformatic pipelines for SNP and INDEL calling from contaminated sputum samples at various levels of theoretical depth of *Mtb* genomic coverage.**

682
683
684 **Figure 6: Performance (F_1 scores) of three bioinformatic pipelines for SNP calling in sputum samples with minimum 20× *Mtb* coverage and various types and levels of contamination.**

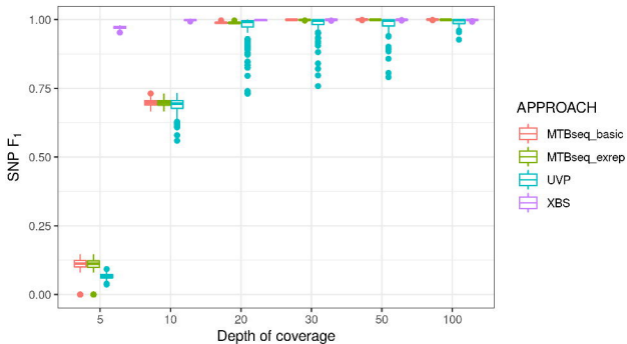
685
686
687 **Figure 7: XBS Maximum Likelihood tree showing the location of Goig et al.'s sputum samples (marked in red) in relation to the reference dataset.**

688

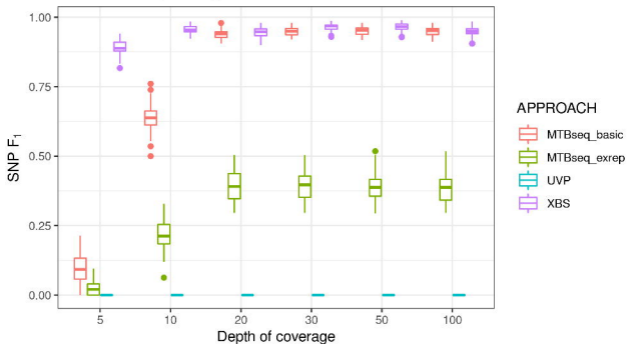




SNP F1 non-complex regions

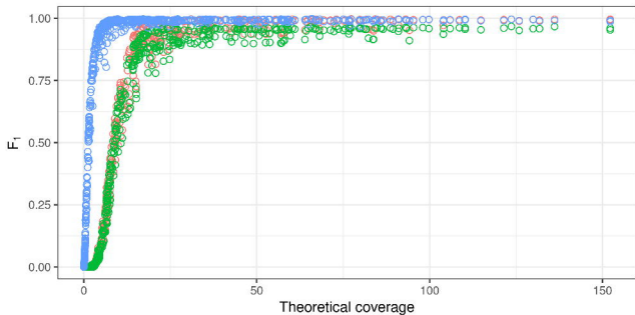


SNP F1 complex regions

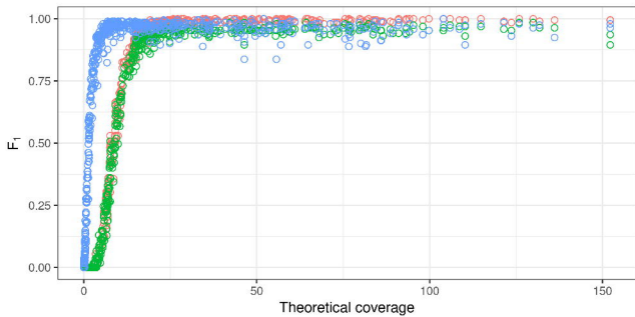


APPROACH ○ MTBseq-basic ○ MTBseq-exrep ○ XBS

SNP F1

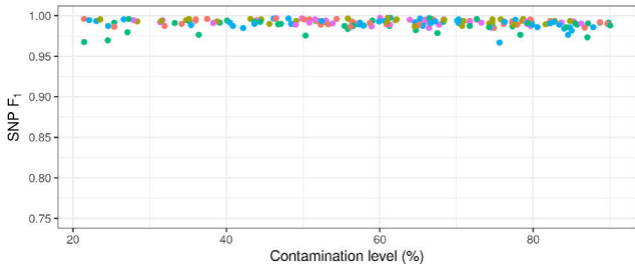


INDEL F1

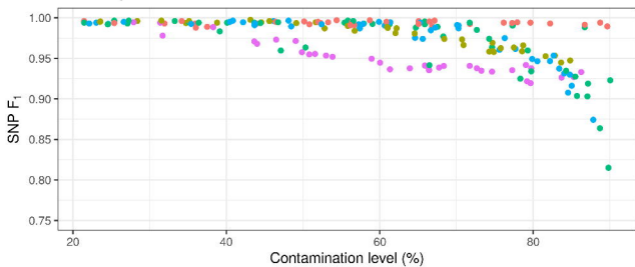


Contaminant: ● H.sap ● H.sap/P.aer./S.Epi. ● NTM mixture ● P.aer. ● S.epi

XBS



MTBseq-basic



MTBseq-exrep

