# Systematic review and meta-analysis of protocols and yield of direct from sputum sequencing of *Mycobacterium tuberculosis*

B.C. Mann[1,3], J. Loubser[1], S. Omar[1], C. Glanz[1], Y. Ektefaie[3], K.R. Jacobson[2], R.M. Warren[1*], M.R. Farhat[3*]

1. DST/NRF Centre of Excellence for Biomedical Tuberculosis Research, SAMRC Centre for Tuberculosis Research, Division of Molecular Biology and Human Genetics, Depts of Biomedical Sciences, Faculty of Medicine and Health Sciences, Stellenbosch University, Cape Town, South Africa
2. Section of Infectious Diseases, Boston University School of Medicine, Boston, MA, USA
3. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

    *Co-senior authors

## Abstract

Direct sputum whole genome sequencing (dsWGS) can revolutionize *Mycobacterium tuberculosis* (*Mtb*) diagnosis by enabling rapid detection of drug resistance and strain diversity without the biohazard of culture. We searched PubMed, Web of Science and Google scholar, and identified 8 studies that met inclusion criteria for testing protocols for dsWGS. Utilising meta-regression we identify several key factors positively associated with dsWGS success, including higher *Mtb* bacillary load, mechanical disruption, and enzymatic/chemical lysis. Specifically, smear grades of 3+ (OR = 14.7, 95% CI: 3.5, 62.1; p = 0.0005) were strongly associated with improved outcomes, whereas decontamination with sodium hydroxide (NaOH) was negatively associated (OR = 0.005, 95% CI: 0.001, 0.03; p = 7e-06), likely due to its harsh effects on *Mtb* cells. Furthermore, mechanical lysis (OR = 193.3, 95% CI: 11.7, 3197.8; p = 0.008) and enzymatic/chemical lysis (OR = 18.5, 95% CI: 1.9, 183.1; p = 0.02) were also strongly associated with improved dsWGS. Across the studies, we observed a high degree of variability in approaches to sputum pre-processing prior to dsWGS highlighting the need for standardized best practices. In particular we conclude that optimizing pre-processing steps including decontamination with the exploration of alternatives to NaOH to better preserve Mtb cells and DNA, and best practices for cell lysis during DNA extraction as priorities. Further and considering the strong association between *Mtb* load and successful dsWGS, protocol improvements for optimal sputum sample collection, handling, and storage could also further enhance the success rate of dsWGS.

## Introduction

Tuberculosis (TB) caused by the bacillus *Mycobacterium tuberculosis* (*Mtb*) remains the most common cause of death from any single infectious pathogen (1, 2). Progress in eradicating TB is hampered by the emergence of multidrug-resistant (MDR) and extensively drug resistant (XDR) *Mtb* strains. According to the latest WHO Global Tuberculosis Report approximately half a million people worldwide developed rifampicin-resistant TB (RR-TB), 78% of which had MDR-TB, defined as resistant to at least isoniazid (INH) and rifampicin (RIF) (2). Next Generation Sequencing (NGS) advances over the past decade have provided the ability to rapidly sequence the whole *Mycobacterium tuberculosis* genome; supplying an extraordinary tool to study the genetic epidemiology of this pathogen, while also detecting single nucleotide polymorphisms (SNPs) and other mutations which can be used to predict susceptibility to first-line drugs (3, 4). Several rapid genotypic drug susceptibility testing (DST) methods are currently endorsed by the WHO, but most of these tests do not provide a comprehensive summary of a patient's drug resistance (DR) profile, since most of these assays only focus on a limited number of targets involved with DR (5). Recently the introduction of targeted NGS (tNGS) has expanded the number of drugs and targets assayed substantially (6). Whole genome sequencing (WGS) however can assay the full breadth of genetic variation and is most comprehensive for predicting phenotype from genotype for *Mtb* (7, 8).

*Mtb* sequencing is still hindered by the long and cumbersome process of culturing *Mtb* for DNA extraction. This process can take weeks and sometimes even months (4). In addition to the long culture period, the culture-based approach has an additional limitation, as culture can change the population structure of the original sample due to the selection of subpopulations more suited to growth in culture and random genetic drift (9).The logical step to avoid these limitations is to sequence *Mtb* directly from clinical specimens. To date several studies have demonstrated that sequencing from direct patient specimens is possible with varying levels of success (3, 4, 9–12). Commercial tNGS assays include a targeted polymerase chain reaction (PCR) amplification step that helps improve sequencing coverage from direct samples, however even tNGS performs poorly on sputum with lower bacillary burden (including Xpert medium and low, or smear negative sputum) with 50-78% or more of test samples yielding only partial coverage of the targets (6, 8, 13). For direct whole genome sequencing, most studies have relied on the use of target capture and enrichment technology which enriches target DNA with a set of target-specific bait probes during library preparation. Others have opted for a selective lysis approach, attempting to selectively lyse contaminating host and

73    bacterial cells, followed by the depletion of contaminating DNA by enzymes such as DNase,
74    either forgoing or only turning to target enrichment in samples for which it was deemed
75    necessary by additional quality control (QC) steps (3, 4, 10, 11). This systematic review and
76    meta-analysis aims to summarize DNA target and non-target-based methods previously
77    utilized to perform whole genome sequencing *Mtb* from direct patient samples. We classify
78    the different approaches used in the literature and perform an individual sample meta-analysis
79    of the effect of specific sample processing steps on direct sequencing success. Although
80    focused on WGS, reviewing non-target-based methods can also have implications for
81    improving yield of tNGS from direct patient samples.

82

# Methods

84

## Search strategy

86

87    In brief we searched PubMed and Web of Science, looking for English articles published on
88    *Mtb* direct sputum whole genome sequencing (dsWGS) to compare various approaches that
89    have previously been applied for successful dsWGS of *Mtb*. This was done to highlight gaps
90    in the current applied methodology that can be improved upon to better facilitate dsWGS. The
91    literature search was carried out using the following keywords "direct + sputum + sequencing
92    + tuberculosis + mycobacteria" and was conducted by three independent reviewers. The three
93    reviewers reviewed titles, abstract, key words and subsequently the full texts to identify and
94    include articles meeting the study inclusion and not meeting the exclusion criteria below.

95    An ethical review was deemed unnecessary as this was a secondary analysis of published
96    articles. We conducted an additional search of Google scholar to identify any relevant articles
97    not identified in the original search. No grey literature, conference papers or unpublished
98    works were included in this analysis because of uncertainty over the relevance and validity of
99    the presented data.

100

101

102

103

## Inclusion and exclusion criteria

The inclusion criteria for eligible publications were defined as all articles that:

- Studied *Mtb*.
- Attempted direct whole genome sequencing or target capture/target enrichment followed by whole genome sequencing (i.e. no targeted PCR) directly on sample without intervening culture.
- Reported on *Mtb* input sputum smear grade or Xpert CT or other *MTB* input DNA quantification.
- Reported on one or more of the following outputs:
  - *Mtb* genomic coverage of drug resistance genes (any subset) or the whole genome.
  - Resistance mutation recall relative to phenotype or to sequencing after culture.
  - Lineage mutation recall relative to sequencing after culture.
- Provided methodological detail on the sputum processing protocol.

The exclusion criteria included articles that:

- Focused on the application to pathogens other than *Mtb*.
- Focused on samples other than sputum.

## Data extraction

Three reviewers independently extracted data from the included studies.The following background information was extracted from eligible papers: Author details, year of publication and population. Sixteen technical variables were extracted under three categories including Sample data and characteristics, Methodology and Results (Table 1). Data extracted were compiled into several predesigned spreadsheets. In instances where data was represented graphically, the authors were contacted to provide the original numerical data used to generate the graphs. If the authors could not be reached or the information was no longer available, we extracted the numerical values from the graphical representations using available software (PlotDigitizer).

**Table 1 Summary of data items extracted from articles fulfilling inclusion and exclusion criteria**

| Data group | Number | Data extracted |
|---|---|---|
| **Sample data and characteristics** | 1 | Population group |
| | 2 | Sample number |
| | 3 | Sample types |
| | 4 | Sputum culture pairs |
| | 5 | Xpert or qPCR data (quantitative or semiquantitative) |
| | 6 | Smear data |
| **Methodology** | 7 | Sample pre-treatment and enrichment methodology |
| | 8 | DNA extraction methodology |
| | 9 | DNA concentrations |
| | 10 | Bioinformatics methodology |
| **Results** | 11 | Coverage for both culture and direct sputum samples |
| | 12 | Percentage on target reads for direct sputum samples |

## Bioinformatics analysis

Sequencing data was downloaded for each included study, either from NCBI GenBank or the European Nucleotide Archive (ENA). Initial quality assessment was done using fastQC, version 0.11.9, followed by adapter removal, quality filtering and per-read low quality base trimming using fastP, version 0.20.1 (14). Following quality control, reads were taxonomically classified using the metagenomic classification tool Kraken2 using the standard database (version 2.0.8) (Wood *et al.*, 2019). Mycobacterial reads were extracted and aligned using bwa-mem2 (version 2.2.1) to the H37Rv reference genome AL123456 (15). Duplicates were removed using Picard and excluded in downstream analyses. Alignment statistics, including number of reads, depth and breadth of coverage and GC-content were determined and visualised using Qualimap (version 2.2.2c) (16).

## Statistical analysis

To assess the effects of various factors on genome and DR region coverage (DR regions comprised of 73 genomic regions strongly associated with a DR phenotyope, identified and curated by the well-known TB-Profiler tool) (17), we employed regression comparing generalized linear mixed models (GLMMs) to support batch control across studies, and a

158  generalized linear model (GLM) without batch control. We found no significant batch effects

159  across the studies and report the GLM results in the main text and the GLMM results in the

160  supplement (Supplementary Tables Y and Z). We assessed the fixed effects of several factors

161  including: smear grade, mechanical disruption, enzymatic/chemical lysis, decontamination,

162  and heat treatment on whole genome coverage (>5x at >90%) and coverage of DR conferring

163  regions (>5x at >95%). Given the limited number of samples sequenced directly we limited

164  the analysis to samples that underwent target capture and enrichment and excluded directly

165  sequenced samples. Associations were assessed using the Wald test with significance

166  assessed at a P-value <0.05. Two processing steps homogenisation and contamination

167  depletion were coded but excluded from the final model due to their utilization in all but one

168  study, or only two target capture studies respectively making it difficult to evaluate their effect.
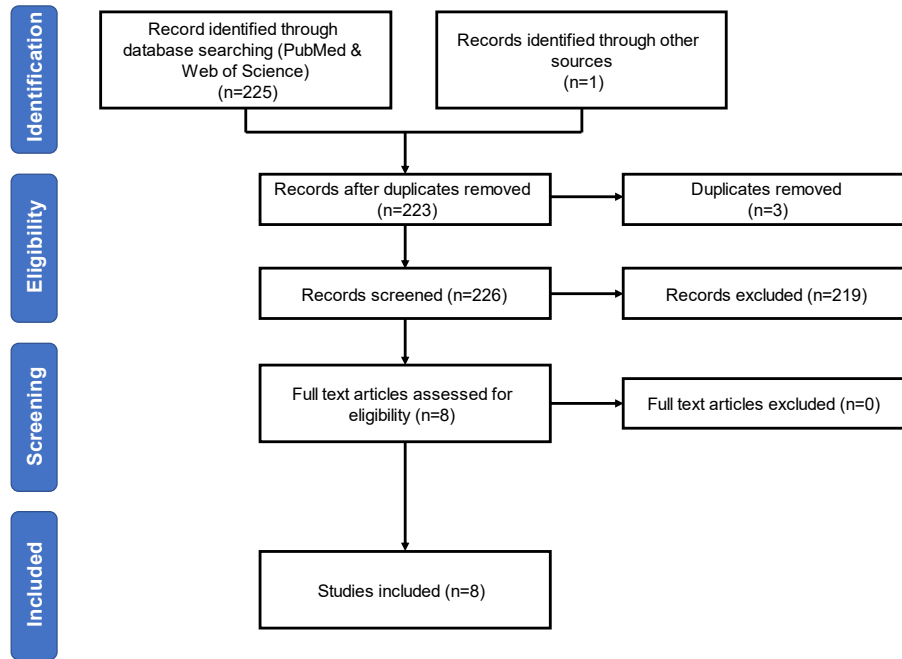
169

# Results

170

171

## Search strategy and study selection

172

173

174  We identified 134 and 91 studies respectively from the PubMed and Web of Science searches.

175  One additional record was identified using a Google scholar search. Three independent

176  reviewers read titles, abstracts, and keywords to assess for duplicates and adherence to

177  inclusion and exclusion criteria. Records that were not excluded at this stage were reviewed

178  in full text to assess the same; 219 records were excluded leaving a total of 8 studies for

179  inclusion in this review.

180

Figure 1: Flow diagram outlining study selection. Abbreviation: PRISMA, Preferred Reporting Items for Systematic Reviews and Meta-Analyses.

## Characteristics of included studies

The number of participants per study ranged from 34 to 100; 5 studies sequenced directly from sputum and compared with sequencing after culture. The largest number of sputum/ culture pairs available from a single study was 43. All studies except 1 included smear data as a semi quantitative measure of bacilli load, and 3 studies included Xpert cycle threshold (Ct) values. Only 2 studies quantified DNA in the input sample by means of qPCR, an additional 2 studies report that qPCR was done to assess bacillary load, but results were not reported, or accessible after contacting the authors (Supplementary Table 1).

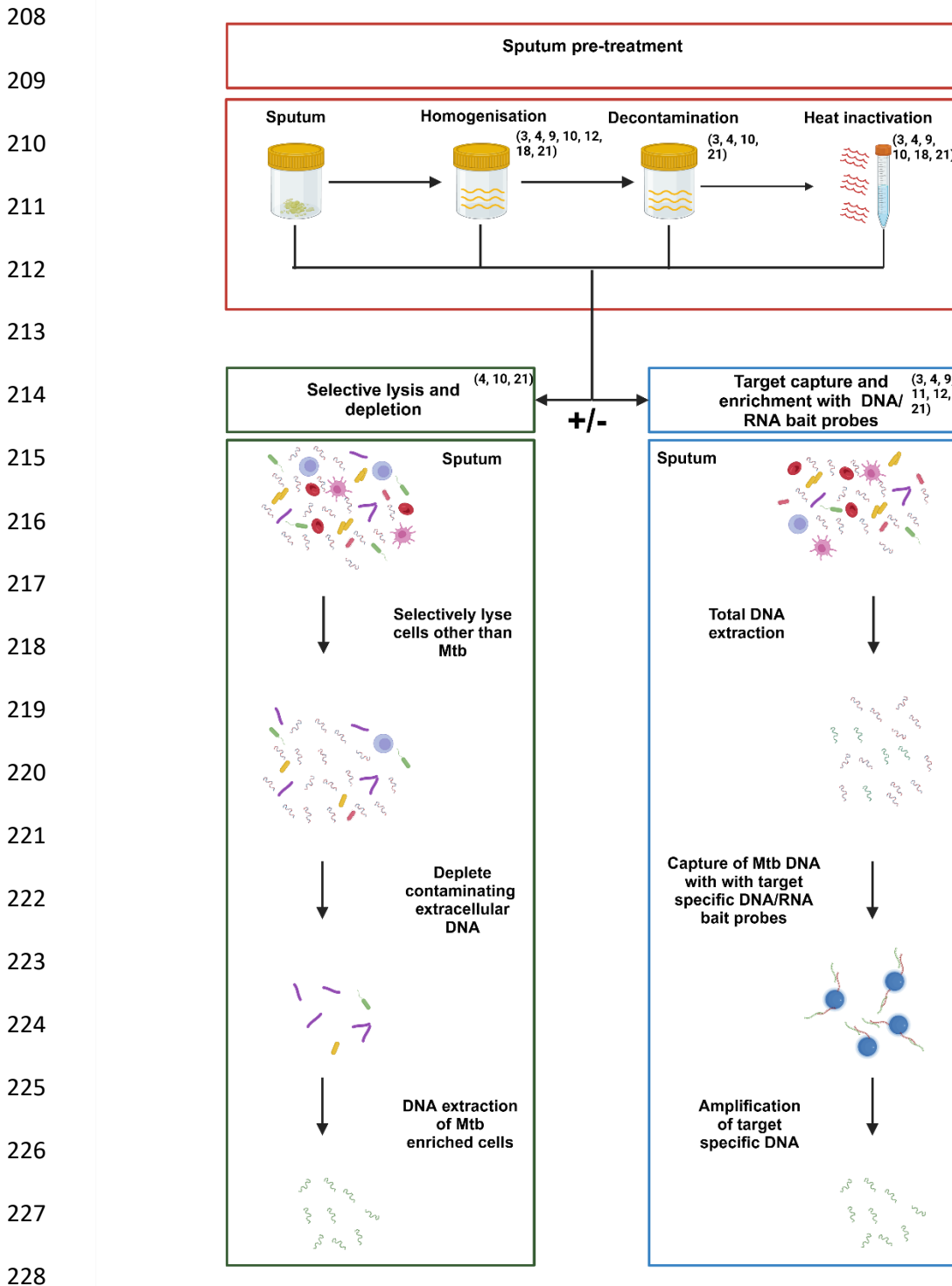## Current approaches facilitating direct sputum sequencing from clinical specimens

The selected studies spanned three primary avenues to facilitate dsWGS with various ways of pre-treating samples before proceeding to one or the other (Figure 2). We define these

201  methods broadly as: 1) non- (DNA) target methods: including selective lysis or other physical
202  or chemical enrichment, which involves the enrichment of *Mtb* cells by breaking down
203  contaminating host and/or commensal bacterial cells, followed by depletion of contaminating
204  DNA by washing or enzymatic degradation and then sequencing; 2) DNA targeting methods
205  specifically the bait capture approach that capture and enrich for *Mtb* DNA with specific
206  DNA/RNA bait probes and; 3) a combination approach involving both selective lysis and
207  contaminating DNA depletion, while also employing bait capture probes (3, 4, 9–12, 18).



the pre-treatment steps. This includes 1) the selective lysis approach for further *Mtb* enrichment where non *Mtb* cells and their DNA are depleted by chemical and/or enzymatic means, or 2) the target capture and enrichment approach where target specific DNA is enriched by means of DNA/RNA bait probes. Studies numbered in the figure as they appear in the bibliography. Created in https://BioRender.com

229

230

231

## Key pre-processing steps that may contribute to the success of dsWGS

233

234 We aimed to identify key steps that may contribute to the enrichment of the *Mtb* target and
235 thus success of dsWGS (Table 2). A step was marked with X if it formed part of the workflow
236 described in the paper. We identified seven steps within either a targeted or non-targeted
237 approach that were employed across the eight studies: specifically, sputum homogenisation,
238 sputum decontamination, heat inactivation, host DNA depletion, commensal microbe DNA
239 depletion, DNA extraction and target enrichment. We briefly summarize the approaches taken
240 under each step below. A more detailed summary of the methodology applied at each step
241 can be found in Supplementary Table 2.

242

243 **Table 2. Direct sputum processing workflow as outlined for each individual**
244 **study**

| Study (reference) | Homogenisation | Decontamination | Heat inactivation | Host DNA depletion | Depletion of commensal microflora | DNA extraction | | | DNA/RNA bait capture |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Mechanical | Enzymatic | Chemical | |
| Brown *et al.* 2015 (3) | X | X | X | | | X | | | X |
| Votintseva *et al.* 2017 (10) | X | X | X | X | | X | | | |
| Doyle *et al.* 2018 (11) | | | | | | X | X | X | X |
| Nimmo *et al.* 2019 (9) | X | | X | | | | X | X | X |
| Soundararajan *et al.* 2020 (12) | X | | | X | | X | X | X | X |
| Goig *et al.* 2020 (4) | X | X | X | X | X | X | | | X |
| George *et al.* 2020 (18) | X | | X | | | X | | | |
| Macedo *et al.* 2023 (21) | X | X | X | | | | X | X | X |

245

## Sputum homogenisation and decontamination

247

248 Although none of the studies utilised homogenisation or decontamination with the aim of
249 enriching for the *Mtb* target, these steps were still incorporated as part of sample pre-
250 processing prior to DNA extraction. All studies reviewed, except one, included a
251 homogenisation step, but approaches varied from the use of the mucolytic agent N-acetyl-L-
252 cysteine (NALC) to the reducing agent dithiothreitol (DTT) or other DTT containing products

253 such as Sputasol. After homogenisation sputum samples are decontaminated in four studies

254 using sodium hydroxide (NaOH) with the goal of preferentially killing other bacteria, fungi, and

255 viruses, thereby reducing the risk of these contaminants influencing diagnostic testing

256 (Supplementary Table 3).

### Heat inactivation

257

258

259 All the studies reviewed except two applied a heat inactivation step commonly used to reduce

260 the biohazard involved with downstream processing (Table 2 and 3). Heat inactivation of *Mtb*

261 specimens involved exposing samples to high temperatures ranging between 80 and 95°C for

262 times ranging from 15 min to 1 hour defined period (3, 4, 9, 11, 12, 18, 21). One of the studies

263 (George *et al.* 2020) demonstrated that heat inactivation can achieve enrichment of tough to

264 lyse cells such *Mtb* but required the addition of a specialised thermal protection buffer to

265 maintain the integrity of *Mtb* DNA during extensive heating (30 min at 99°C) which also

266 subsequently lead to the degradation of any extracellular host DNA (18).

267

### Host and commensal microbe DNA depletion

268

269

270 Votinseva *et al.* (2017) and Soundararajan *et al.* (2020), aimed to enrich for *Mtb* by applying

271 commercially available kits namely the MolYsis Basic5 kit and the Ultra Deep Microprep DNA

272 isolation kit (Molzym, Germany) for the depletion of host DNA prior to *Mtb* WGS. Goig *et al.*

273 (2020), used GTC solution (4M guanidinium thiocyanate 4M, 0·5% w/v sodium N-lauryl

274 sarcosine, 25mM trisodium citrate, 0·1M 2mercaptoethanol, 0·5% w/v Tween 80) instead to

275 lyse eukaryotic cells in conjunction with DNase. In addition to this, the study by Galo *et al.*

276 (2020), was also the only study to directly preform additional depletion with the aim to not only

277 lyse eukaryotic cells but also gram-negative bacterial cells utilising a GTC buffer, while leaving

278 tough-walled *Mtb* cells intact (4).
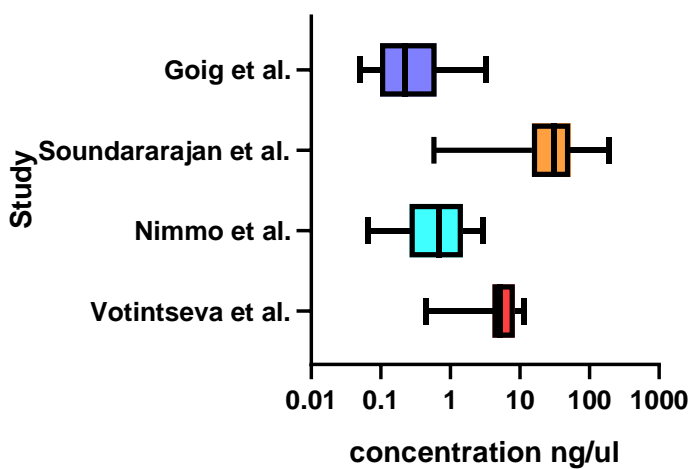
279

### Lysis and DNA extraction

280

281

282 Mycobacteria are known to be difficult to lyse (22). Failure to achieve adequate lysis in the

283 context of dsWGS can result in a biased representation of target *Mtb* DNA relative to

284 contaminating DNA (10, 22). DNA concentrations post extraction were available for four of the

285 studies included in this review (Figure 3). DNA extraction methods varied across the reviewed

286 studies (Supplementary Table 3). All studies, except Soundararajan *et al.* (2020) and Macedo

287 *et al.* (2023), employed a combination of chemical and mechanical cell lysis. Soundararjan *et*

288   *al.*, extracted DNA utilising the Ultra Deep Microprep DNA isolation kit (Molzym, Bremen,
289   Germany), while Macedo *et al.* utilised the QIAmp DNA Mini Kit, both of which according to
290   the manufacturer's instructions include both chemical and enzymatic means to facilitate lysis
291   but omit any mechanical steps. None of the studies assessed the effect of specific pre-
292   processing/DNA extraction steps on total DNA concentration by molecular quantification of
293   DNA before and after any given step.

294

295

296



| Study (Reference) | Average DNA yield ng/ul | Processing steps that may impact DNA yield |
|---|---|---|
| Goig *et al* 2020 (4) | 0.57±0.86 | Heat inactivation, decontamination, selective lyses with GTC, DNase treatment |
| Sundaresan *et al.* 2020 (12) | 39.42±34.93 | Ultra-Deep Microprep DNA isolation kit |
| Nimmo *et al.* 2019 (9) | 0.97±0.97 | Heat inactivation, decontamination |
| Votintseva *et al.* 2017 (10) | 6.18±3.1 | Heat inactivation, decontamination, treatment with Molzym Basic5 |

297

298

299   Figure 3: DNA concentration measured after DNA extraction and prior to target capture (if any) for the 4 studies
300   with available data.

301

302   Target enrichment

303

304   Six studies used DNA or RNA bait capture to enrich for *Mtb* DNA (3, 4, 9–12, 18, 21).
305   However, only two studies (Brown *et al.* (2015) and *Goig et al.* (2020)), compared sequencing
306   yield with or without bait capture to measure the increase of sequencing reads attributed to
307   the *Mtb* target . The former compared the percentages of on-target reads (%OTR), and the
308   mean sequencing depths for two sputum samples. Brown *et al.* reported a percent of on-target
309   reads (%OTR) of 0.3%, with a sequencing depth of 4.6x without bait capture, compared to
310   82% and 200x respectively with bait capture. Goig *et al.*, quantified the target *Mtb* DNA in the
311   input and used this information to target bait enrichment to the lowest input samples. We used
312   the raw data provided with this publication to assess percentage of *Mtb* target pre and post

313 bait capture (Figure 4). The average % of *Mtb* target pre-bait capture was 1,67% compared to
314 48,5% post-capture (4).
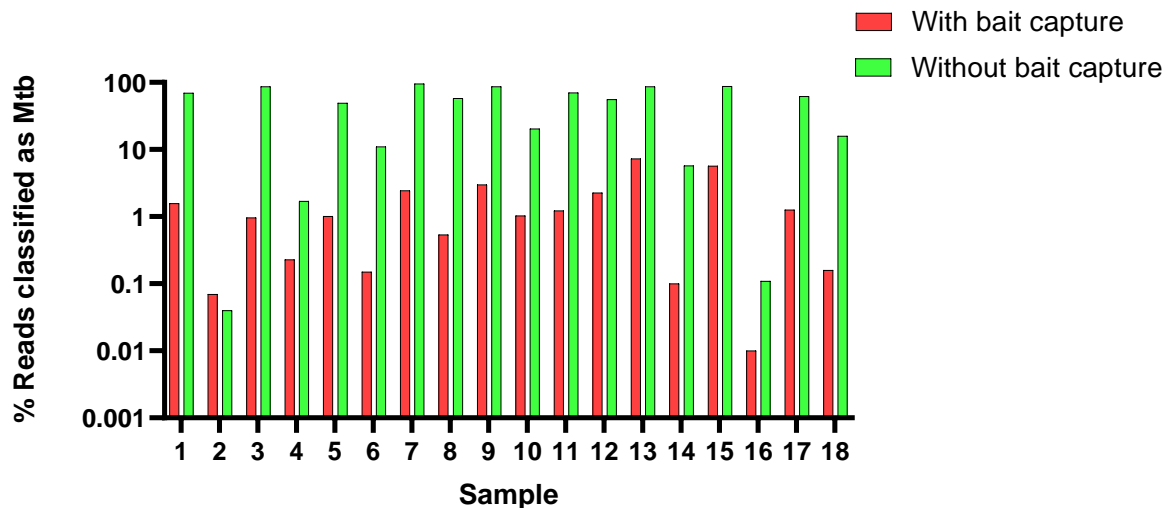
315

316



**Goig et al. With - without bait capture**

Figure 4 Increase in *Mtb* target read % in samples sequenced with and without bait capture from Goig *et al.,* 2020 (4). Note the y-axis is on the logarithmic scale.

## Processing steps predictive of successful dsWGS

329 The effects of sample characteristic (smear grade), and processing steps prior to target
330 enrichment and sequencing (mechanical disruption for lysis, enzymatic/chemical lysis,
331 decontamination with NaOH and heat treatment) were studied as predictors dsWGS success
332 using logistic regression (Methods) with data extracted for 289 samples. The effect of
333 processing steps was consistent across the dsWGS success metric used (>95% of drug
334 resistance regions covered at >5x read depth or >90% of the whole genome covered at >5x
335 depth, Table 3, Supplementary Table Y) and we observed no evidence of batch effects by
336 study (Supplementary table Z). The results identify samples with higher smear grade as more
337 likely to be successfully sequenced directly, as expected. In addition, mechanical disruption,
338 and enzymatic/chemical lysis were also associated with higher dsWGS success prior to target
339 capture. Sputum decontamination with NaOH on the other hand resulted in lower dsWGS
340 success when used prior to target capture (Table 3).

341

**Table 3 Sputum processing step that are significantly associated with dsWGS success based on the whole genome coverage generalized linear model.**

| Characteristic | OR | 95% CI | p-value |
|---|---|---|---|
| Smear neg | ref | | |
| 2+,1+,scanty | 3.4 | (0.98, 12.1) | 0.07 |
| 3+ | 14.7 | (3.5, 62.1) | 0.0005 |
| | | | |
| Mechanical lysis | 193.3 | (11.7, 3197.8) | 0.008 |
| Enzymatic/Chemical lysis | 18.5 | (1.9, 183.1) | 0.02 |
| Decontamination | 0.005 | (0.001, 0.03) | 7e-06 |
| Heat Inactivation | 2.25 | (1.1, 4.5) | 0.02 |

Associations were performed using a sample-level logistic regression model pooled across studies with dsWGS success defined as whole genome coverage >90% at read depth of >5x per site. OR = Odds Ratio, CI = Confidence Interval.

342

# Discussion

344

345 This systematic review and meta-analysis summarize key experimental steps that may
346 influence the success of dsWGS for *Mtb*. We identified and reviewed 226 study records,
347 ultimately including 8 studies that met our inclusion criteria. Included studies varied in sample
348 size, ranging from 34 to 100 participants/*Mtb* isolates, and utilized various approaches for the
349 pre-process of sputum samples prior to sequencing. The current review demonstrates that
350 although there is overlap in the applied methodology across studies there is simultaneously a
351 lot of variability in the specific processing steps employed for sequencing directly from sputum.
352 This review classifies key pre-processing steps that may or may not contribute to the success
353 of dsWGS, these include sputum homogenization, decontamination, heat inactivation,
354 depletion of host and commensal microbial DNA, lysis and DNA extraction, and target
355 enrichment through DNA/RNA bait capture (3, 9, 11, 12, 21).

356

357 Due to the amount of contaminating cells and DNA present in sputum the main overlapping
358 trend across studies is the utilisation of bait capture and enrichment probes which was utilised
359 by all the reviewed studies except that of Votintseva *et al.* (2017). Bait capture relies on the
360 hybridisation of *Mtb* specific biotinylated DNA/RNA probes that bind to complimentary *Mtb*
361 DNA, bait/target DNA hybrids are then captured with streptavidin magnetic beads and pulled
362 down magnetically allowing for the selective enrichment of the *Mtb* target (23). Studies by
363 Brown *et al.* (2015) and Goig *et al.* (2020), have sequenced samples both enriched and

364  unenriched for comparison clearly demonstrating the effectiveness of the addition of target
365  capture and enrichment systems for dsWGS, however the impact of additional pretreatment
366  steps has not yet been elucidated (3, 4). Homogenization, heat inactivation, and the depletion
367  of contaminating host and bacterial cells/DNA are key steps identified during the literature
368  review that could potentially enhance *Mtb* enrichment. Currently reviewed studies have not
369  critically assessed the impact of various sample treatment/preparation steps on the final
370  outcome/success of dsWGS, but have consistently highlighted correlation between smear
371  grade/*Mtb* load and improved sequencing performance (11, 12, 21).

372

373  Other than bait capture, there is limited data on the impact of specific sputum processing steps
374  on dsWGS success (3, 4, 9–12, 21). To address this, we used logistic regression to perform
375  a meta-analysis of the effect of a subset of these steps on dsWGS success controlling for the
376  *Mtb* load in the sputum as measured by smear microscopy. Study data allowed the evaluation
377  of the effects of four processing steps and only prior to target capture: specifically, sputum
378  decontamination with NaOH, mechanical disruption, enzymatic/chemical lysis, and heat
379  treatment. The results support the latter three steps as significantly increasing the success of
380  sputum dsWGS. The *Mtb* cell wall is difficult to lyse and this may explain why a combination
381  approach which involves chemical, enzymatic and mechanical lysis contributes to improved
382  *Mtb* DNA recovery and thus also potentially improved sequencing results (4, 24). Although
383  heat inactivation was employed with the goal of sterilizing the sample and thereby reduce the
384  biohazard risk of downstream processing, George *et al.* (2020) (18) demonstrated that heat
385  inactivation can enrich for tough-to-lyse cells like *Mtb* in the setting of thermal protection buffer
386  supporting the meta-analysis association between heat duration and temperature and dsWGS
387  success that we observe across studies. The available data limited our ability to study the
388  effect of other processing steps, such as host DNA depletion,  using the MolYsis Basic5 kit,
389  the Ultra Deep Microprep DNA isolation kit, or a GTC solution combined with DNase treatment
390  (4, 10, 12). Confirming the effectiveness of these methods in depleting host DNA and enriching
391  *Mtb* will require additional future study.

392

393  Our meta-analysis supports sputum decontamination using NaOH treatment as decreasing
394  dsWGS success in the setting of target capture. The use of NaOH is standard practice to
395  reduce live contaminants in sputum prior to *Mtb* culture (19). While it is known to create a
396  highly alkaline environment that is inhospitable to most microorganisms except *Mtb*, its impact
397  on dsWGS, specifically in terms of *Mtb* enrichment and potential loss of target cells and DNA,
398  has not been thoroughly evaluated (25, 26). NaOH can selectively lyse contaminating cells

399  that are not of interest for downstream analysis, but studies support a risk of *Mtb* cell loss (4,

400  19, 25, 26). Our finding raises a need to reevaluate the effect of NaOH treatment on sample

401  composition and the exploration of alternative methods as potentially more suitable for WGS

402  which does not require viable *Mtb* bacilli.

403

## Conclusion

404
405

406  Despite the observed heterogeneity of approaches for dsWGS, a common trend observed

407  during the course of this systematic review and meta-analysis is the utilization of target capture

408  and enrichment probes, which is observed to be highly effective in enhancing *Mtb* sequencing

409  from direct patient samples (3, 4, 9, 12). The efficacy of these probes nevertheless depends

410  on the overall *Mtb/Mtb* DNA load in sputum. Target capture probes are expensive and

411  alternative or additional processing steps that can deplete contaminants or directly enrich for

412  *Mtb* DNA will be beneficial to facilitate dsWGS and reduce cost (4, 10). Future research should

413  thus focus on refining the identified pre-processing steps to enhance the robustness and also

414  reliability of dsWGS with the ultimate aim of developing standardized pre-processing protocols

415  to advancing DR profiling directly from clinical sputum samples (4, 9, 10, 12). Considering the

416  importance of *Mtb* load highlighted in the current study a suggested additional research

417  direction is to optimise sputum collection, standardize sputum volume, storage, handling and

418  transport with the aim of further improving *Mtb* bacilli yield prior to the application of target

419  capture and enrichment (27, 28).

420

## Acknowledgements

421
422

427

428

429

# References

1. Chakaya J, Khan M, Ntoumi F, Aklillu E, Fatima R, Mwaba P, Kapata N, Mfinanga S, Hasnain SE, Katoto PDMC, Bulabula ANH, Sam-Agudu NA, Nachega JB, Tiberi S, McHugh TD, Abubakar I, Zumla A. 2021. Global Tuberculosis Report 2020 – Reflections on the Global TB burden, treatment and prevention efforts. Int J Infect Dis 113:S7–S12.

2. Global Tuberculosis Report 2024. https://www.who.int/teams/global-tuberculosis-programme/tb-reports/global-tuberculosis-report-2024. Retrieved 13 November 2024.

3. Brown AC, Bryant JM, Einer-Jensen K, Holdstock J, Houniet DT, Chan JZM, Depledge DP, Nikolayevskyy V, Broda A, Stone MJ, Christiansen MT, Williams R, McAndrew MB, Tutill H, Brown J, Melzer M, Rosmarin C, McHugh TD, Shorten RJ, Drobniewski F, Speight G, Breuer J. 2015. Rapid Whole-Genome Sequencing of Mycobacterium tuberculosis Isolates Directly from Clinical Samples. J Clin Microbiol 53:2230–2237.

4. Goig GA, Cancino-Muñoz I, Torres-Puente M, Villamayor LM, Navarro D, Borrás R, Comas I. 2020. Whole-genome sequencing of Mycobacterium tuberculosis directly from clinical samples for high-resolution genomic epidemiology and drug resistance surveillance: an observational study. The Lancet Microbe 1:e175–e183.

5. WHO consolidated guidelines on tuberculosis: Module 3: diagnosis – rapid diagnostics for tuberculosis detection . 2021. WHO consolidated guidelines on tuberculosis: Module 3: diagnosis – rapid diagnostics for tuberculosis detection. World Health Organization, Geneva. http://www.ncbi.nlm.nih.gov/books/NBK572344/. Retrieved 19 July 2023.

6. Mansoor H, Hirani N, Chavan V, Das M, Sharma J, Bharati M, Oswal V, Iyer A, Morales M, Joshi A, Ferlazzo G, Isaakidis P, Ndlovu Z, England K. 2023. Clinical utility of target-

454      based next-generation sequencing for drug-resistant TB. Int J Tuberc Lung Dis 27:41–

455      48.

456   7.   Ness TE, DiNardo A, Farhat MR. 2022. High Throughput Sequencing for Clinical

457      Tuberculosis: An Overview. Pathogens 11:1343.

458   8.   Dookie N, Khan A, Padayatchi N, Naidoo K. 2022. Application of Next Generation

459      Sequencing for Diagnosis and Clinical Management of Drug-Resistant Tuberculosis:

460      Updates on Recent Developments in the Field. Frontiers in Microbiology 13.

461   9.   Nimmo C, Shaw LP, Doyle R, Williams R, Brien K, Burgess C, Breuer J, Balloux F,

462      Pym AS. 2019. Whole genome sequencing Mycobacterium tuberculosis directly from

463      sputum identifies more genetic diversity than sequencing from culture. BMC Genomics

464      20:389.

465   10.   Votintseva AA, Bradley P, Pankhurst L. 2017. Same-Day Diagnostic and Surveillance

466      Data for Tuberculosis via Whole-Genome Sequencing of Direct Respiratory Samples.

467      Journal of Clinical Microbiology 55.

468   11.   Doyle RM, Burgess C, Williams R, Gorton R, Booth H, Brown J, Bryant JM, Chan J,

469      Creer D, Holdstock J, Kunst H, Lozewicz S, Platt G, Romero EY, Speight G, Tiberi S,

470      Abubakar I, Lipman M, McHugh TD, Breuer J. 2018. Direct Whole-Genome

471      Sequencing of Sputum Accurately Identifies Drug-Resistant Mycobacterium

472      tuberculosis Faster than MGIT Culture Sequencing. J Clin Microbiol 56:e00666-18.

473   12.   Soundararajan L, Kambli P, Priyadarshini S, Let B, Murugan S, Iravatham C, Tornheim

474      JA, Rodrigues C, Gupta R, Ramprasad VL. 2020. Whole genome enrichment approach

475      for rapid detection of Mycobacterium tuberculosis and drug resistance-associated

476      mutations from direct sputum sequencing. Tuberculosis 121:101915.

477  13.  Colman RE, Seifert M, Rossa AD la, Georghiou SB, Hoogland C, Uplekar S, Laurent S,

478      Rodrigues C, Kambli P, Tukvadze N, Maghradze N, Omar SV, Joseph L, Suresh A,

479      Rodwell TC. 2024. Evaluating culture-free targeted next-generation sequencing for

480      diagnosing drug-resistant tuberculosis: a multicentre clinical study of two end-to-end

481      commercial workflows. The Lancet Infectious Diseases 0.

482  14.  Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ

483      preprocessor. Bioinformatics 34:i884–i890.

484  15.  Vasimuddin Md, Misra S, Li H, Aluru S. 2019. Efficient Architecture-Aware Acceleration

485      of BWA-MEM for Multicore Systems, p. 314–324. *In* 2019 IEEE International Parallel

486      and Distributed Processing Symposium (IPDPS). IEEE, Rio de Janeiro, Brazil.

487  16.  Okonechnikov K, Conesa A, García-Alcalde F. 2016. Qualimap 2: advanced multi-

488      sample quality control for high-throughput sequencing data. Bioinformatics 32:292–294.

489  17.  Phelan JE, O'Sullivan DM, Machado D, Ramos J, Oppong YEA, Campino S, O'Grady

490      J, McNerney R, Hibberd ML, Viveiros M, Huggett JF, Clark TG. 2019. Integrating

491      informatics tools and portable sequencing technology for rapid detection of resistance

492      to anti-tuberculous drugs. Genome Medicine 11:41.

493  18.  George S, Xu Y, Rodger G, Morgan M, Sanderson ND, Hoosdally SJ, Thulborn S,

494      Robinson E, Rathod P, Walker AS, Peto TEA, Crook DW, Dingle KE. 2020. DNA

495      Thermo-Protection Facilitates Whole-Genome Sequencing of Mycobacteria Direct from

496      Clinical Samples. J Clin Microbiol 58:e00670-20.

497  19.  Dippenaar A, Ismail N, Grobbelaar M, Oostvogels S, de Vos M, Streicher EM, Heupink

498      TH, van Rie A, Warren RM. 2022. Optimizing liquefaction and decontamination of

499      sputum for DNA extraction from Mycobacterium tuberculosis. Tuberculosis

500      132:102159.

501    20.   Burdz TVN, Wolfe J, Kabani A. 2003. Evaluation of sputum decontamination methods

502          for Mycobacterium tuberculosis using viable colony counts and flow cytometry.

503          Diagnostic Microbiology and Infectious Disease 47:503–509.

504    21.   Macedo R, Isidro J, Ferreira R, Pinto M, Borges V, Duarte S, Vieira L, Gomes JP.

505          2023. Molecular Capture of Mycobacterium tuberculosis Genomes Directly from

506          Clinical Samples: A Potential Backup Approach for Epidemiological and Drug

507          Susceptibility Inferences. IJMS 24:2912.

508    22.   Kok NA, Peker N, Schuele L, de Beer JL, Rossen JWA, Sinha B, Couto N. 2022. Host

509          DNA depletion can increase the sensitivity of Mycobacterium spp. detection through

510          shotgun metagenomics in sputum. Front Microbiol 13:949328.

511    23.   Mann BC, Jacobson KR, Ghebrekristos Y, Warren RM, Farhat MR. 2023. Assessment

512          and validation of enrichment and target capture approaches to improve Mycobacterium

513          tuberculosis WGS from direct patient samples. bioRxiv

514          https://doi.org/10.1101/2023.03.12.530724.

515    24.   Epperson LE, Strong M. 2020. A scalable, efficient, and safe method to prepare high

516          quality DNA from mycobacteria and other challenging cells. Journal of Clinical

517          Tuberculosis and Other Mycobacterial Diseases 19:100150.

518    25.   Shehadul Islam M, Aryasomayajula A, Selvaganapathy PR. 2017. A Review on

519          Macroscale and Microscale Cell Lysis Methods. Micromachines (Basel) 8:83.

520    26.   Zhou J, Zhang M, Li X, Wang Z, Pan D, Shi Y. 2021. Performance comparison of four

521          types of target enrichment baits for exome DNA sequencing. Hereditas 158:10.

522    27.   Datta S, Shah L, Gilman RH, Evans CA. 2017. Comparison of sputum collection

523          methods for tuberculosis diagnosis: a systematic review and pairwise and network

524          meta-analysis. The Lancet Global Health 5:e760–e771.

525    28.  Meyer AJ, Atuheire C, Worodria W, Kizito S, Katamba A, Sanyu I, Andama A, Ayakaka

526         I, Cattamanchi A, Bwanga F, Huang L, Davis JL. 2017. Sputum quality and diagnostic

527         performance of GeneXpert MTB/RIF among smear-negative adults with presumed

528         tuberculosis in Uganda. PLOS ONE 12:e0180572.

529

530

531

532

533

534

535

536