

# This item is the archived peer-reviewed author-version of:

Whole genome sequencing of Mycobacterium tuberculosis : current standards and open issues

## **Reference:**

Meehan Conor J., Goig Galo A., Kohl Thomas A., Verboven Lennert, Dippenaar Anzaan, Ezew udo Matthew , Farhat Maha R., Guthrie Jennifer L., Laukens Kris, Miotto Paolo, ....- Whole genome sequencing of Mycobacterium tuberculosis : current standards and open issues Nature review s: microbiology - ISSN 1740-1526 - London, Nature publishing group, 17:9(2019), p. 533-545 Full text (Publisher's DOI): https://doi.org/10.1038/S41579-019-0214-5 To cite this reference: https://hdl.handle.net/10067/1601130151162165141

uantwerpen.be

Institutional repository IRUA

### Whole genome sequencing of *Mycobacterium tuberculosis*: current standards 1

### 2 and open issues

- 3
- 4 Conor J Meehan, Galo A. Goig, Thomas Andreas Kohl, Lennert Verboven, Anzaan Dippenaar,
- 5 Matthew Ezewudo, Maha Farhat, Jennifer L. Guthrie, Kris Laukens, Paolo Miotto, Boatema
- 6 Ofori-Anyinam, Viola Dreyer, Philip Supply, Anita Suresh, Christian Utpatel, Dick van Soolingen,
- 7 Yang Zhou, Philip Ashton, Daniela Brites, Andrea M. Cabibbe, Bouke C. de Jong, Margaretha de
- 8 Vos, Fabrizio Menardo, Sebastien Gagneux, Qian Gao, Tim H Heupink, Qingyun Liu, Chloé
- 9 Loiseau, Leen Rigouts, Timothy C Rodwell, Elisa Tagliani, Timothy M. Walker, Robin Mark
- 10 Warren, Yanlin Zhao, Matteo Zignol, Marco Schito, Jennifer Gardy, Daniela Maria Cirillo, Stefan
- 11 Niemann, Inaki Comas\* and Annelies Van Rie\*
- 12

#### 13 Affiliations:

- 14 CJM Unit of Mycobacteriology, Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp,
- 15 Belgium
- 16 GAG Institute of Biomedicine of Valencia, CSIC and CIBER in Epidemiology and Public Health, Valencia, Spain
- 17 TAK Molecular and Experimental Mycobacteriology, Priority Area Infections, Research Center Borstel, D-23845 18 Borstel, Germany
- 19 LV Tuberculosis Omics Research Consortium, Department of Epidemiology and Social Medicine, Institute of Global
- 20 Health, Faculty of Medicine and Health Sciences, University of Antwerp
- 21 AD DST-NRF Centre of Excellence for Biomedical Tuberculosis Research; South African Medical Research Council
- 22 Centre for Tuberculosis Research; Division of Molecular Biology and Human Genetics, Faculty of Medicine and
- 23 Health Sciences, Stellenbosch University, Cape Town, South Africa
- 24 ME Critical Path Institute, Arizona, USA
- 25 MF Harvard Medical School, and Massachusetts General Hospital, Boston, MA, USA
- 26 JLG University of British Columbia, Vancouver, Canada
- 27 KL Adrem Data Lab, Department of Mathematics & Computer Science, University of Antwerp
- 28 PM Emerging Bacterial Pathogens Unit, Division of Immunology, Transplantation and Infectious Diseases, IRCCS
- 29 San Raffaele Scientific Institute, Milano, Italy
- 30 BO Center for Global Health Security and Diplomacy, Ottawa, Canada, Food and Drugs Authority, Accra, Ghana
- 31 VD Molecular and Experimental Mycobacteriology, Priority Area Infections, Research Center Borstel, D-23845
- 32 Borstel, Germany
- 33 PS University Lille, CNRS, INSERM, CHU Lille, Institut Pasteur de Lille, U1019, UMR 8204, CIIL, Centre d'Infection et
- 34 d'Immunité de Lille, Lille, France
- 35 AS Foundation for Innovative New Diagnostics (FIND), Geneva, Switzerland
- 36 CU Molecular and Experimental Mycobacteriology, Priority Area Infections, Research Center Borstel, D-23845
- 37 Borstel, Germany
- 38 DvS National Tuberculosis Reference Laboratory, Centre for Infectious Disease Control, National Institute for Public
- 39 Health and the Environment (RIVM), Bilthoven, The Netherlands
- 40 ZY National Center for Tuberculosis Control and Prevention, Chinese Center for Disease Control and Prevention,
- 41 Beijing, China
- 42 PA Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine, University of Oxford, UK
- 43 DB Swiss Tropical and Public Health Institute, Basel, Switzerland
- 44 AMC Emerging Bacterial Pathogens Unit, Division of Immunology, Transplantation and Infectious Diseases, IRCCS
- 45 San Raffaele Scientific Institute, Milano, Italy
- 46 BCdJ Unit of Mycobacteriology, Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp, 47 Belgium
- 48 MDV DST-NRF Centre of Excellence for Biomedical Tuberculosis Research; South African Medical Research Council
- 49 Centre for Tuberculosis Research; Division of Molecular Biology and Human Genetics, Faculty of Medicine and
- 50 Health Sciences, Stellenbosch University, Cape Town, South Africa

- 51 FM Swiss Tropical and Public Health Institute, Basel, Switzerland
- 52 SG Swiss Tropical and Public Health Institute, Basel, Switzerland
- 53 QG Key Laboratory of Medical Molecular Virology, Ministry of Education and Health, School of Basic Medical
- 54 Sciences, Fudan University, Shanghai, China.
- 55 THH Tuberculosis Omics Research Consortium, Department of Epidemiology and Social Medicine, Institute of
- 56 Global Health, Faculty of Medicine and Health Sciences, University of Antwerp
- 57 CL Swiss Tropical and Public Health Institute, Basel, Switzerland
- 58 QL Key Laboratory of Medical Molecular Virology, Ministry of Education and Health, School of Basic Medical
- 59 Sciences, Fudan University, Shanghai, China.
- 60 LR Unit of Mycobacteriology, Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp, Belgium
- 61 TCR Foundation for Innovative New Diagnostics (FIND), Geneva, Switzerland
- 62 ET Emerging Bacterial Pathogens Unit, Division of Immunology, Transplantation and Infectious Diseases, IRCCS San
- 63 Raffaele Scientific Institute, Milano, Italy
- 64 TMW Nuffield Department of Medicine, University of Oxford, John Radcliffe Hospital, Oxford, UK
- 65 RMW DST-NRF Centre of Excellence for Biomedical Tuberculosis Research; South African Medical Research Council
- 66 Centre for Tuberculosis Research; Division of Molecular Biology and Human Genetics, Faculty of Medicine and
- 67 Health Sciences, Stellenbosch University, Cape Town, South Africa
- 68 YZ National Center for Tuberculosis Control and Prevention, Chinese Center for Disease Control and Prevention, 69 Beijing, China
- 70 MZ Global Tuberculosis Programme, World Health Organization, Geneva, Switzerland
- 71 MS Critical Path Institute, Arizona, USA
- 72 JG University of British Columbia, Vancouver, Canada
- 73 DMC Emerging Bacterial Pathogens Unit, Division of Immunology, Transplantation and Infectious Diseases, IRCCS
- 74 San Raffaele Scientific Institute, Milano, Italy
- 75 SN Molecular and Experimental Mycobacteriology, Priority Area Infections, Research Center Borstel, D-23845 76 Borstel, Germany
- 77 IC Institute of Biomedicine of Valencia, CSIC and CIBER in Epidemiology and Public Health, Valencia, Spain
- 78 AVR Tuberculosis Omics Research Consortium, Department of Epidemiology and Social Medicine, Institute of
- 79 Global Health, Faculty of Medicine and Health Sciences, University of Antwerp
- 80

#### 81 Abstract

- 82 Whole genome sequencing (WGS) of *Mycobacterium tuberculosis* has rapidly evolved from a 83 research tool to a clinical application for the diagnosis and management of tuberculosis and in 84 public health surveillance. This evolution has been facilitated by the dramatic drop in costs,
- 85 advances in technology, and concerted efforts to translate sequencing data into actionable
- 86 information. There is however a risk that, in the absence of a consensus and international
- 87
- standards, the widespread use of WGS technology may result in data and processes that lack 88
- harmonisation, comparability and validation. In this review, we outline the current landscape of
- 89 WGS pipelines and applications and set out best practices for *M. tuberculosis* WGS, including 90
- standards for bioinformatics pipelines, curated repository of resistance-causing variants, phylogenetic analyses, quality control processes, and standardised reporting.
- 91
- 92
- 93 1. Introduction
- 94 Mycobacterium tuberculosis complex (Mtbc) pathogens are collectively the top infectious
- disease killer globally, causing 10 million new tuberculosis (TB) cases annually<sup>1</sup>. Increasingly, 95
- 96 new TB cases are already resistant to rifampicin and isoniazid (termed multidrug resistance;
- MDR-TB), the key first line drugs<sup>1</sup>. Tackling the spread and drug resistance burden of this 97
- 98 pathogen requires concerted global effort in prevention, diagnosis, treatment and surveillance.

99 Over the past decades, research and public health practices, including contact investigation and 100 phenotypic methods for drug susceptibility testing (DST), have been complemented by 101 molecular approaches. These can now provide rapid diagnosis, drug susceptibility profiling, and 102 an understanding of *Mtb* transmission dynamics<sup>2,3</sup>.

Whole genome sequencing (WGS) approaches use DNA sequencing platforms to reconstruct 103 104 the complement of DNA found inside a cell. The small (~4.4Mb), single chromosome genome of 105 Mtbc strains<sup>4</sup> lends itself well to WGS approaches. . Rapid, reliable, and increasingly affordable WGS technologies, can now guide all components of TB control: diagnosis, treatment, 106 surveillance and contact tracing<sup>5,6</sup> (Fig. 1). Individual (sub)species of human and animal Mtbc 107 lineages can be identified,<sup>7-9</sup> and drug resistance profiles can be predicted, especially well for 108 1<sup>st</sup> line drugs<sup>2</sup>, enabling prompt, appropriate initiation of treatment and monitoring the 109 acquisition of drug resistance<sup>10</sup>. TB outbreaks can be identified with high resolution<sup>11-13</sup>, 110 including across borders,<sup>14,15</sup> and diseases control measures implemented. The analysis of the 111 emergence, spread, genetic makeup, and evolution of particular outbreak strains, e.g. highly 112 resistant or highly virulent clones, can enable the development of targeted measures<sup>16–18</sup>. 113

114 WGS-based approaches are quickly moving from research-only to clinical care and public health 115 applications. The World Health Organization (WHO) is already using WGS for drug resistance surveillance<sup>19</sup> and is scheduled to evaluate sequencing technologies for routine genotypic drug 116 susceptibility testing in 2019<sup>1</sup>. As WGS-guided individualized treatment<sup>20</sup> and WGS-based 117 surveillance systems<sup>15</sup> are being implemented in several countries (e.g. the UK and the 118 119 Netherlands) with more to come, accurate methods and standardized reporting are vital. At 120 present, multiple WGS data analysis solutions exist that vary widely in scope, pipelines, and output formats, with little standardisation amongst them<sup>21</sup>, making cross-comparisons and 121 122 rigorous validation of these pipelines difficult. Because clinical decisions such as the effective 123 drugs that can be included in a patients' regimen may be influenced by differences in the 124 bioinformatic analysis, robustness of the pipeline used in clinically-relevant predictions tools is 125 critical.

126 In this review, we present the current state of the art for the three core Mtbc WGS tasks: drug 127 susceptibility profiling, transmission cluster detection and subspecies/lineage identification 128 (referred to as strain typing). We highlight those places where a general agreement in the 129 analysis parameters or interpretation of the results has been already reached by the 130 community. Alternatively, we discuss those items where there is still open discussion about the 131 best practices and will require more effort to reach a consensus in the future.

132

## **133 2. State of the art**

The standard workflow for WGS analysis of Mtbc strains is outlined in Figure 2. It involves culturing sputum specimens on solid (Löwenstein–Jensen) or liquid (Mycobacteria Growth Indicator Tube) media, extracting DNA from Mtbc strains, library preparation, and sequencing by short read technologies (e.g. Illumina platforms)<sup>22</sup>. The complete Mtbc WGS analysis pipeline involves several key steps such as input data validation and quality control followed by mapping to a reference genome (often H37Rv) and detection of genomic variants such as single nucleotide polymorphisms (SNPs) and insertion/deletions (indels). Numerous resequencing pipelines for the Mtbc currently exist with currently no single 'gold standard'. These pipelines typically exclude about ~10% of the genome because erroneous mapping in certain regions result in false variant calls (PE/PPE gene families, other repetitive genes, mobile elements<sup>4</sup>) and apply various criteria, such as read depth, base quality, and strand bias to filter out false positive variants. Finally, based on the detected variants, several tasks can be performed including (but not limited to) prediction of drug resistance and susceptibility profiles, strain typing, and identification of transmission clusters.

148

149 Due to the clonality of their genomes and their inability to undergo lateral gene transfer, Mtbc 150 strains acquire drug resistance primarily through variants in core genes or promoters<sup>23,24</sup>. Drug 151 resistance and susceptibility profiles can be determined with high accuracy for many drugs used for the treatment of TB by comparing variant calls to lists of high-confidence resistance 152 153 conferring variants. These lists have been established primarily using genotype-phenotype associations identified from statistical analyses of large sets of clinical WGS data<sup>25,26</sup> (Fig. 3). A 154 155 prime effort in the construction of these lists is the Relational Sequencing Tuberculosis Data 156 Platform (ReSegTB, http://www.resegtb.org), where researchers from around the world can contribute data<sup>27</sup>. This database contains curated, aggregated genotypic and phenotypic 157 158 information on global Mtbc isolates accompanied by metadata including clinical outcome. 159 Another important initiative is the Comprehensive Resistance Prediction for Tuberculosis: an 160 International Consortium (CRyPTIC) project. CRyPTIC aims to better understand the relationship 161 between genetic variants and minimum inhibitory concentrations (MIC) for most drugs used for TB treatment<sup>2</sup>. By comparing the SNPs present in a sequenced isolate to these lists, WGS can 162 not only predict resistance but also 1<sup>st</sup> line pan-susceptibility under specific conditions<sup>2</sup>, 163 164 replacing the need for phenotypic testing.

165

Similarly, strain classification of the seven major human-associated lineages, many of the animal-associated lineages, and their sub-lineages, can be derived directly from variant calls using lists of lineage-defining SNPs<sup>7–9</sup>. This is important for understanding population structure and potential phenotypic differences between lineages<sup>28</sup> and comparing isolates on the global level<sup>18,29</sup>.

The genomic data for a set of isolates can also be used for surveillance and transmission 171 172 investigations. For this, the most common approach is to use a SNP cut-off-based clustering 173 although genome-based multi locus sequence typing (MLST) has shown comparable results<sup>30,31</sup>. 174 The SNP cut-off approach starts by constructing a list of high-confidence, unambiguous SNPs 175 found in each isolate, often excluding indels and drug resistance related sites. This filtering is 176 important when predefined SNP distance thresholds are used to cluster strains and define 177 recent transmission chains. Given the very low genetic diversity of the Mtbc, thresholds of 5 or 178 12 SNPs are frequently used to suggest epidemiological links, although these thresholds were calibrated in low incidence settings with a diverse strain population<sup>32</sup>. It is not yet clear if a 179 180 single threshold can be employed to detect epidemiologically linked cases in all timeframes and 181 contexts. The MLST approach employs a predefined set of shared genes and assigns a number 182 to each allele sequence identified for each gene. Coded allele combinations can be compared 183 between strains to detect potential transmission clusters. Two schema exist for this approach: the core genome (termed cgMLST; 2891 genes covering 2.86 million bases<sup>31</sup>) and an extended 184

pan-genome including 1141 accessory loci<sup>11</sup> (termed wgMLST). These WGS-based approaches
 have been shown to perform better than contact tracing and with higher resolution than
 classical approaches such as MIRU-VNTR<sup>12,13,30,31,33</sup>.

188

This currently recommended data processing workflow (Fig. 2) leading to SNP-based drug resistance profiling, transmission clustering at a given SNP cut-off and strain profiling using lineage-defining SNPs is often robust and reliable. However, steps towards standardisation and validation of this workflow are required to ease integration into current clinical and public health initiatives.

194

195 Currently, two Mtbc-specific pipelines are available, which perform multiple core tasks in single install set-up to produce genetic variant calls from raw Illumina sequence data (MTBseq<sup>34</sup> and 196 UVP-ReSeqTB<sup>35</sup>). Other pathogen-agnostic pipelines can be used with an Mtbc-specific 197 reference genome and drug resistance database to achieve similar results<sup>33,36–38</sup>. Numerous 198 custom-built pipelines also exist<sup>8,39–46</sup>, often incorporating similar tools for mapping and variant 199 200 calling with additional accessory tools and in-house scripts to parse and refine outputs. A non-201 exhaustive list of such pipelines is given in supplementary table 1 to demonstrate the range of 202 tools and settings routinely implemented. Lastly, pipelines specific for a single task such as drug resistance prediction<sup>25,47–51</sup> or strain typing<sup>7,50</sup> are available and have been comprehensively 203 204 compared elsewhere<sup>52–55</sup>.

205

# 206 3. Mtbc WGS validation and standardisation

207 Before a workflow can become a gold standard, the validity of that workflow needs to be 208 ensured for its intended uses. For Mtbc WGS workflows, this essentially means ensuring 209 virtually every variant that is reported is truly present in the isolate (validation) and each 210 pipeline calls the same variants (standardisation). Ideally, all steps of the workflow, from DNA 211 extraction to sequencing, data analysis and reporting, should be standardised (or at least 212 comparable) and well documented, and an external quality assessment (EQA) program should 213 be in place. Efforts to standardise and validate the upstream (pre-bioinformatics pipeline) steps have been undertaken to great effect<sup>22,54</sup>. Pipeline standardisation could be achieved through 214 the use of a single pipeline in all settings or through validation with rigorous testing and 215 216 convergence on a defined outcome for all pipelines developed. Since multiple pipelines have already been implemented (e.g. MTBseq<sup>34</sup> for the EUSeqMyTB consortium and the Unified 217 Variant Pipeline<sup>35</sup> for ReSegTB) (supplementary table 1), agreement on validation criteria seems 218 219 more realistic. Since WGS-based diagnostics present a potential paradigm shift for regulatory approvals, there is an urgent need to understand how to validate and standardise these 220 multiple pipelines for clinical use<sup>56</sup>. In 2016, the US Food and Drug Administration (FDA) 221 222 released draft guidelines on sequencing-based infectious disease diagnostics and bodies such as the WHO and ECDC are taking steps towards international standardisations of Mtbc WGS<sup>15,22,57</sup>. 223

224 225

## a. Technical validation and external quality control of Mtbc WGS

First, the extraction of DNA needs to meet minimal standards as defined for a given WGS instrument<sup>22</sup>. Next, the pipeline to convert the raw sequencing reads into accurate variant calls should be technically valid, i.e. call the correct variantss. While there is much debate about the

229 reference standard to be used for technical validation of WGS pipelines, currently this is best 230 undertaken by using short read datasets derived from isolates with known complete genomes (e.g. from long read sequencing)<sup>58</sup>. Mapping these read sets to their respective assembled 231 genomes allows to calculate the rate of false positive and negative SNPs called by the pipeline 232 233 under consideration. Ideally, to promote interoperability and ease the verification of 234 bioinformatics protocols, a standard reporting format such as a BioCompute Object (BCO) to record all thresholds, steps and implementation arguments for a given pipeline is utilised<sup>59</sup>. 235 236 Comparisons of BCOs from different pipelines can then be used to set acceptable lower limits for the assessed parameters, refining technical validation criteria across pipelines<sup>60</sup>. 237

A prime example of external quality control of bioinformatics pipelines is the efforts by the National Institute for Public Health and the Environment (RIVM) to standardize the use of WGS for Mtbc genotyping across the European Reference Laboratory Network for TB (ERLTB-Net)<sup>21</sup>. Panels of DNA extracted from selected Mtbc isolates are sent annually by RIVM to reference laboratories to assess intra- and inter laboratory reproducibility of WGS. Similar efforts in high burden settings are needed to monitor the reliability of Mtbc WGS outputs when used in these

244 settings.245

## b. Validation for core tasks: transmission, phylogeny and drug resistance

Task validation is used to demonstrate that a given pipeline is verified for a specific analysis, e.g. drug resistance profiling. For task validation, Mtbc bioinformatics pipelines should use defined validation datasets, ideally with hundreds or thousands of well characterized clinical Mtbc strains representing the diversity of a specific core task (e.g. different drug susceptibility profiles for resistance detection, representatives of all Mtbc lineages and (sub)-species for typing, or varying degrees of clustering for transmission analyses). The number of readily available, well-curated validation datasets is currently limited.

254

**Validation of transmission clustering**. The national public health institute of the Netherlands (RIVM) has provided laboratories with sequenced reads from 535 Mtbc isolates for which epidemiological links were known. Using this dataset, the EUSeqMyTB consortium showed that existing pipelines could confidently distinguish linked from unlinked cases, especially when the SNP distances are high, as is often the case in low burden settings<sup>12</sup>. This comparison was undertaken as part of an effort to standardise WGS for monitoring MDR-TB cross border transmission in Europe<sup>15</sup>.

262

Validation of classification systems. The clonality of Mtbc strains means that lineage and strain typing can be performed using only a handful of SNPs that are specific for strains of a particular lineage. Several studies have demonstrated the reliability of specific SNPs to determine the Mtbc (sub)lineage<sup>8,9,61</sup>. However, sub-lineage classifications are often less resolved, and parallel nomenclatures for lineage 2 are being used<sup>18,62,63</sup>. As the diversity of the Mtbc is further explored, especially for animal-associated and zoonotic TB, these under-described lineages can also easily be typed using the same SNP-based approach<sup>7</sup>.

270

271 *Validation of drug resistance profiling.* Validation of WGS for TB resistance is the most 272 advanced of all the core tasks. Phelan *et al* showed high concordance between phenotypic and

genotypic predictions, no matter the sequencing platform used<sup>19,54</sup>. In the past two years, 273 major progress has been made in the linkage between genotype and resistance phenotype by 274 employing a standardized statistical approach<sup>25,26</sup>. The task of incrementally improving our 275 knowledge base on genetic resistance profiling is primarily being addressed by the two global 276 consortia outlined above: ReSegTB's single platform for genotype-phenotype investigation of 277 drug resistance<sup>27,35</sup> and CRyPTIC's genotypic-phenotypic linking of over 10,000 isolates 278 demonstrating susceptibility prediction for rifampicin and isoniazid with 99% sensitivity and 93-279 96% for ethambutol and pyrazinamide<sup>2</sup>. These results have led to some low burden countries 280 (Netherlands, UK) replacing phenotypic DST with WGS-based DST for first line drugs. Resistance 281 predictions for 2<sup>nd</sup> line drugs can also be undertaken with sensitivity often around 90%<sup>25</sup>. Large 282 comparative studies using phenotype-genotype associations are expanding the catalogues<sup>64,65</sup> 283 284 and will help to increase the sensitivity for drugs used to treat MDR-TB. Efforts are now 285 directed towards increasing the diversity of isolates and including accompanying high quality 286 phenotypic and clinical data, especially for new anti-TB drugs.

287

## a. Standardization of communication of Mtbc WGS results and data sharing

289 Communication to end users: Effective communication of WGS-based results to a diverse 290 audience of end-users is key to positively impacting patient care and TB control programs. While the need for plain language reporting of genomic results has been recognized<sup>52,66</sup>, there 291 292 are no international standards yet. Reporting standards should be flexible enough to address 293 the varying levels of familiarity of end-users with genomic data interpretation and allow 294 customization to region-specific treatment guidelines and formatting requirements. For 295 example, the ISO15189:2012 standard mandates information such as patient identifiers, assay 296 details, and the testing laboratory be reported. Recommendations from Mtbc WGS report 297 design validation studies included the use of complete terms instead of abbreviations, drawing attention to important elements with shading, bolding, and other types of emphasis, and 298 incorporating summary statements to rapidly communicate key results<sup>67,68</sup>. 299

300

301 Communication to the research community. In peer-reviewed publications, the parameters 302 used at each step of a bioinformatics pipeline must be stated in a way that makes it 303 reproducible and understandable to non-bioinformaticians (e.g. using a BCO as outlined above). 304 Custom code used in the analysis should be made available through a public repository 305 (e.g.GitHub), ensuring ease of installation elsewhere. Pipelines should report the outcome of 306 technical validations, at least for the core tasks they aim to address (e.g. lineage-defining SNPs 307 for a typing pipeline). Examples of standard reporting include the MIABi (Minimum Information About a Bioinformatics investigation)<sup>69</sup> and the STROME-ID (Strengthening the Reporting of 308 Molecular Epidemiology for Infectious Diseases)<sup>70</sup> guidelines. In supplementary table 2, we 309 310 suggest data elements to include according to intended use, but note that a report may need to 311 include elements from more than one use case.

312 **Data sharing** will be crucial as incremental knowledge improves drug resistance predictions and 313 strain tracking relies on the number and diversity of strain genome data available. This can 314 come in the form of sharing coded strain identifiers such as MLST patterns or raw sequence 315 data not yet processed by a pipeline. Indeed data sharing has been shown already to be

- 316 invaluable for detecting cross-Europe transmission clusters<sup>14</sup>. Data sharing should encompass
- 317 data produced by research and collected in public health laboratories and surveillance efforts<sup>71</sup>,
- 318 similar to the GenomeTrakr network for foodborne pathogens<sup>72</sup>, while still safeguarding patient
- 319 data and appropriately acknowledging contributions. This setup would be of great value for
- 320 moving the field of Mtbc WGS forward.
- 321

The crucial next step for fully utilising Mtbc WGS data is implementation of validations, both technical and task oriented, for all pipelines. Once undertaken, the agreed upon pipeline(s) can then be widely implemented, once infrastructure and usability is accounted for.

325

## 326 2. Implementation of WGS in routine clinical practice

While the use of WGS is rapidly expanding in research, minimal progress has been made in programmatic use of WGS. Some reasons include the lack of standardised end-to-end solutions, the required wet-lab and computing infrastructure, need for sufficient internet connectivity and bandwidth, and training deficits in genomics and bioinformatics<sup>73–75</sup>. Efforts are thus needed to expand accessibility to perform analysis by non-experts. How these factors are addressed will depend a country's income and public health sector strength.

333

334 High-income countries will probably use a mixture of closed (end-to-end) solutions and more 335 complex pipelines as they likely will have on-site bioinformatics support. Ideally, routine analysis of WGS will require little to no bioinformatics knowledge by the end user. 336 337 Implementation of these pipelines can be undertaken by either local set-ups with supporting infrastructure or a cloud/web-based approach with easy, affordable access<sup>76</sup>. Many large 338 339 healthcare facilities such as referral hospitals are already incorporating bioinformatics units into 340 their support services as part of the push towards personalized medicine, something TB 341 treatment can take advantage of. These services should mediate the implementation of 342 complex pipelines and make all required software readily available without a requirement to install additional software tools, as is done with certain existing pipelines<sup>34,77</sup>. 343

344

345 Giving the heterogeneity of pipelines already in place (e.g. supplementary table 1) it is 346 conceivable that something similar will happen when implementation is done in hundreds of 347 care services. Some will opt-in for end-to-end solutions, perhaps integrated with the 348 sequencing platform, or others for task-specific, such as resistance prediction only. Those 349 implementing their own pipeline should be aware of the limitations, cautions and recommendations detailed by expert consensus here and elsewhere<sup>6,76</sup>. In order to evaluate 350 new pipelines it is preferable to develop inside 'containers', such as Docker or Singularity<sup>78,79</sup>, or 351 one-command installation wrappers like Bioconda or Homebrew<sup>80,81</sup>. Creating a container for 352 353 each step (Figure 2) also allows for easy updating of a specific step without the need to install a 354 whole new pipeline and allows for tasks (e.g. resistance profiling) to be added to the pipeline as 355 needed. To allow usability by a range of end-users, fine-grained access to the individual steps 356 should be available for advanced users with functionality layers abstracted away for users with 357 limited bioinformatics expertise. The pipelines should be open source and user-friendly, by 358 employing intuitive and well-documented command line and graphical user interfaces with 359 relevant and validated default parameters.

360

361 The situation in LMIC countries, especially those with a high burden of TB is currently totally 362 different. End-to-end solutions based on cloud computing are the most logical step forward similar to the roll-out of qPCR systems (Box 1). Centralized web-based analysis platforms have 363 recently emerged and promise to aid in computational efficiency, access and usability<sup>47,51</sup>. Roll-364 365 out of such initiatives to more countries would greatly improve the potential for large-scale 366 WGS implementation. The primary barrier to this is usually unstable internet connectivity with limited bandwidth, although using methods that can effectively handle connection 367 interruptions, such as BioTorrents<sup>82</sup>, or direct transfer from sequencing centres to cloud storage 368 and/or web-based pipelines may help circumvent these issues. 369

370

371 The use of end-to-end, cloud-based solutions is likely to play an important role in LMICs. It is, however, advisable to build in those countries human capacity for WGS of Mtbc strains<sup>83,84</sup>. 372 373 While standardised, immutable pipelines are optimal for global implementation of WGS, there 374 are several reasons why local bioinformatics knowledge is required, such as the necessity to 375 adapt analyses to the country-specific epidemiological profiles and public health ecosystems or 376 regulatory laws that do not allow storage beyond country borders. Such customised, yet 377 reproducible solutions are being supported by capacity building initiatives (e.g. the Human, 378 Heredity and Health in Africa Consortium (https://h3abionet.org) and the TORCH consortium 379 (https://torch-consortium.com/vliruos)). TB supranational reference laboratories should also play an important coordinating role, as is currently done for phenotypic workflows<sup>19,85</sup>. 380 381 Ultimately, expanding education curricula to include bioinformatics are needed to generate 382 sufficient capacity<sup>86</sup>.

383

Finally, supportive policy and political commitment will be essential for sustainable implementation of WGS, especially in TB endemic LMICs<sup>74,83,87</sup>. This implementation will benefit from the lessons learned during the step-wise approach used to roll-out line probe assays and GeneXpert (Box 1)<sup>88</sup>.

- 388
- 389

## **390 3.** Extensions of the current standard

391

While current pipelines (Fig. 2) appear to be highly accurate for many aspects of the three core tasks, multiple important issues remain open and should be part of future research and evaluation.

395

# **396 a. Input data validation and quality control**

Most current pipelines do not routinely filter out reads that do not come from the sequenced Mtbc strains. However, sequencing files can contain reads from other organisms and these contaminants can introduce errors during the variant calling process, modifying both the variants identified and their respective frequencies<sup>89</sup>. Additionally, any host DNA sequencing reads should be removed especially if the data is shared online for legal/ethical reasons. Computationally removing non-Mtbc strain reads prior to mapping is an efficient strategy to

implement contamination-proof analysis pipelines<sup>40</sup>, but requires a taxonomic classification of 403 404 individual reads. Using taxonomic classification methods, where reads are assigned to the 405 closest matched species, allows for quick and efficient removal of contaminating reads but requires comprehensive genome databases, often making their implementation extremely 406 memory consuming<sup>90,91</sup>. Additionally, elimination of reads from highly conserved core bacterial 407 genes of heterologous sources still remains a problem. Proposed alternatives include masking 408 genomic regions known to accumulate artefactual polymorphisms<sup>89</sup>, filtering the alignments 409 produced by contaminant reads, or fine-tuning the read aligners such that only the Mtbc strains 410 411 sequences are mapped to the reference genome. Any methodology will require thorough 412 technical validation to ensure that contaminant reads are removed without eliminating true 413 Mtbc sequences, e.g. through in silico generation of datasets with varying levels of reads from 414 other organisms.

415

## 416 **b.** Sequence read mapping and reference genomes

417 The use of a single reference genome for mapping all Mtbc strains is the ideal approach for comparable and standard variant calling. While most pipelines use the H37Rv genome<sup>4,92</sup> as the 418 419 reference genome, several alternative approaches should be explored. Since H37Rv is a lineage 420 4 strain, its use as a reference for other lineages may be insufficient due to gene content differences between lineages<sup>93–96</sup>. Additionally, H37Rv contains many variants not found in any 421 other strain<sup>97</sup>, including in genes related to drug resistance (e.g. gyrAS95T), creating confusion 422 423 in SNP interpretations. Any replacement of H37Rv as the reference genome should be assessed 424 by in-silico studies across datasets and clinical settings. An example of such a study tested seven 425 different references against sequence reads from lineage 4 isolates showing that very limited 426 variation occurred, and that reference choice should be based on criteria other than matching lineage<sup>98</sup>. 427

428 One alternative to the H37Rv genome is a pan-genome which incorporates the entire gene pool 429 of Mtbc lineages. Previous studies have found small but notable differences in gene content between lineages, often affecting genes involved in pathogenesis<sup>93–96</sup>. While these differences 430 431 are unlikely to affect drug resistance profiling (since associated mutations are in the core 432 genome), they may impact delineation of transmission clusters if additional SNPs are found in 433 these genes that would push strain comparisons over the predetermined thresholds. Building a 434 *Mtbc* pan-genome should be straightforward due to the close relationship between different 435 strains (average nucleotide identity between any two strains ≥ 99.8%) and the lack of horizontal 436 gene transfers events. So far this approach has not been effectively explored.

A second alternative is the use of an inferred ancestral genome representative of the Mtbc population and diversity<sup>29,40</sup>. From an evolutionary perspective, this approach addresses the H37Rv-specific variants outlined above. In addition, because all extant strains are equidistant to a common ancestor, the number SNPs called for any Mtbc strain will be similar (normalized) regardless of its lineage. This expected SNP range is useful for quality control, as deviations may indicate poor quality sequencing, co/super-infections and contaminations<sup>40</sup>.

A third approach is to use ad-hoc reference genomes, depending on the study being conducted.
 For instance, lineage-specific ancestral genomes or high-quality, closed, outbreak-specific

- 445 reference genome<sup>99–101</sup> could be used as reference to reduce mapping errors<sup>10</sup>. A disadvantage
- of this approach is that it hampers comparison of results between pipelines and standardizedreporting of results.
- 448 A completely different alternative involves de-novo assembly, using a reference-free approach,
- 449 which has been successfully applied for human population genomics data<sup>102</sup>.
- 450

Independent of the selection of the reference genome, other steps such as mapping and filtering are not consistent between different pipelines, yet might greatly affect the analysis outcome. For instance, removal of duplicates, both PCR and optical, may have a large impact in the variants identified and the allele frequencies. Similarly, local assembly/realignment around indels, reducing false positive SNPs derived from mapping artefacts, is rarely used in Mtbc WGS pipelines<sup>58</sup> but is known to affect variant calling<sup>47</sup>. The question of whether these steps have a relevant effect on the final outcome should be incorporated into future technical validations.

458

## 459 c. Interpretation of drug resistance results and predictions

460 Currently, the bulk of routine drug resistance testing is undertaken using pDST. While this 461 approach will still be required for a subset of difficult to interpret drug resistance patterns, the 462 overarching goal is to detect all variants associated with resistance for comprehensive genome-463 based resistance profiling. While the current statistical approach to calling resistance-464 associated variants using WGS data is an important step forward for clinical use, a weakness is 465 that phenotype predictions of rare and/or novel genetic variants cannot be assessed (Fig. 3). 466 This problem is especially relevant for new and repurposed drugs, or drugs such as 467 pyrazinamide and ethionamide for which mutations are not limited to hotspots but appear 468 across genes (pncA and ethA) and in promoter regions. For these drugs, the standard statistical approach could be complemented by experimental data, comprehensive single nucleotide 469 mutagenesis<sup>103</sup> followed by systematic phenotypic screening, multi-omics studies, and machine 470 learning approaches to predict the resistance phenotype of uncommon or novel genomic 471 variants<sup>104,105</sup>. With the final aim of replacing the majority of phenotypic DST by sequence-472 473 based testing, it will also be essential to catalogue "benign" variants that are not associated to 474 resistance, i.e. phylogenetic markers or other neutral variants<sup>2</sup>. New statistical approaches like large-scale GWAS<sup>64,65</sup>, protein structure modelling<sup>44,106</sup> and machine learning<sup>104,105,107</sup> will likely 475 476 play a key role identifying causative versus benign variants. Comprehensive databases of WGS 477 data linked with phenotypic and clinical outcome data (e.g. CRyPTIC or ReSeqTB) are key to 478 moving towards this goal.

479

480 Once established, endorsement of a single standardised variant list by the WHO or other 481 regulating bodies, with regular updating should be favoured.

482

## 483 d. Variant calling for other purposes

484

Accurate variant calling has major implications on downstream interpretation of the results for evolutionary, epidemiological and clinical applications. Because of the low level diversity and the slow substitution rate of Mtbc genomes<sup>32,42,100,108</sup>, a few falsely called SNPs can affect the interpretation of transmission events, impact the classification of a second episode of TB as
 relapse versus re-infection, or influence the interpretation of sub-populations within a patient
 (Fig. 4).

A primary use of Mtbc WGS is the identification of recent transmission chains and its direction 491 at high resolution. While some studies have used thresholds from 0 to <50 SNPs<sup>109-111</sup>, a 492 threshold of 5- or 12-SNP genetic distances is most frequently used to identify possible 493 epidemiological links and recent transmission<sup>30,32</sup>. For WGS-based distinction of relapse versus 494 reinfection, studies have used often arbitrary thresholds of < 6 or <10 SNPs to define 495 reactivation, and >100 to >1306 to define re-infection<sup>46,112,113</sup>. Any threshold selection can be 496 497 problematic as inferences based on relatedness must include possible underlying methodological bias (culture, sampling and pipeline). In addition, genetic distances may be 498 impacted by biological factors such as potential mutational bursts<sup>42,114</sup>, clonal variants in 499 different lesions<sup>10,115</sup>, the impact of strain type (lineage or subspecies) or drug resistance on 500 substitution rates<sup>108,116</sup>, and genome stability/instability during latency<sup>116,117</sup>. For example, 501 identifying transmission from unrelated cases or distinguishing relapse and reinfection in low 502 503 burden countries is relatively easy, where the distribution of SNP distances is bimodal, separating linked from unlinked cases<sup>12,14</sup>. Conversely, inferring transmission clusters within the 504 505 context of institutional- or household settings or in high TB-incidence scenarios where the SNP 506 distance distribution is continuous remains difficult especially if epidemiological links in large clusters of patients with seemingly identical strains are lacking<sup>118–120</sup>. 507

508

509 Other approaches have meanwhile been developed to improve the identification of 510 epidemiological links and outbreak reconstruction beyond SNP-based clustering. These either use transmission event thresholds<sup>121</sup> and/or often combine genomic and epidemiological data 511 to identify the most probable transmission trees for infectious diseases<sup>122,123</sup>. Of particular 512 513 importance when reconstructing Mtbc outbreaks is that phylogeny and transmission events do not necessarily coincide as a results of genetic diversification during latency and long 514 generation times<sup>124</sup>; it is thus necessary to model the within-host genetic dynamics<sup>125–127</sup>. 515 Besides transmission reconstruction, phylodynamic approaches also allow for the inference of 516 517 epidemiological relevant parameters such as the effective reproduction number as well as the timing and geographic origin of an outbreak<sup>128,129</sup>. 518

519

520 Unravelling within-host dynamics in terms of subpopulation detection remains even more 521 challenging. Low frequency variants that are not due to technical artefacts can indicate the 522 presence of mixed infections (two distinct Mtbc strains co-circulating in a host), or microevolution leading to closely related subpopulations or heteroresistance (subpopulations 523 that differ in drug resistance-related variants)<sup>10,115,130</sup>. Proposed sub-population detection limits 524 525 in different pipelines vary considerably from 10% to <75% (supplementary table 1) and are 526 strongly influenced by factors such as read depth. While the presence of a subpopulation of at least 1% resistant bacilli is considered clinically relevant<sup>131</sup>, the chain reaction of selection bias 527 means that what is observed in sequencing data may not be representative of what is present 528 529 in the culture isolate, which in turn is likely not representative of the diversity in the sputum sample, which is known to not represent the entirety of the within-patient diversity<sup>115,132</sup>.
Mathematical modelling approaches have been developed to identify mixed infections<sup>133,134</sup>.
However with the current approaches the detection of mixed infections is limited by the
relative ratio of the two strains and the number of differing SNPs between both. Future
research and methodological improvements are needed to better understand and interpret this
within-host diversity.

536

## 537 **4. Beyond the current standards**

538

539 As current culture-based approaches require time for Mtbc strain growth, culture-free WGS, 540 directly from clinical samples (e.g.sputum), would be transformative for clinical and public 541 health applications of WGS. This approach would not only eliminate the culture delay but also 542 remove culture selection biases. While studies have shown some success, this approach is still 543 mired with problems such as contamination by human and commensal microbial reads, preventing sufficient coverage depth of the Mtbc genomes and thus reliable variant calling, 544 even in samples with high bacterial loads<sup>135–137</sup>. Improvements in cell lysis or capture coupled 545 with selective DNA enrichment or depletion could reduce this technical complexity and cost. 546 547 Additionally, downstream bioinformatic filtering could be used to control for and remove 548 possible remaining false variants.

549

550 Much is expected from the development of highly portable sequencing devices (e.g. the 551 MinION). Such technology offers the capacity to detect variants in real-time during sample 552 acquisition, potentially giving results from sputum within hours if mycobacterial loads are high. 553 Their portability and ability to work in resource limited settings also favours direct sequencing 554 of clinical samples, even in LMICs. Moreover, although progress has been made in analysis of variants in repeat-rich genome regions (e.g. PE/PPE family genes) or structural changes 555 (duplications, large indels, etc.) by short read mapping<sup>112,138</sup>, long read sequencing will make 556 this more robust<sup>101,135</sup>. Unfortunately, application of this technology is currently limited by high 557 558 error rates (although new dual sequence reading systems promise substantial improvement) 559 and, specifically for mycobacteria, difficulty in cell lysis without over-shearing DNA.

560 561 F

# 561 5. **Conclusion**

562 A decade after first proof-of-principle studies, the community consensus is that Mtbc WGS is 563 now mature enough to inform clinical decisions and public health. This is evident as WGS has 564 already replaced phenotypic testing for first line drugs in some settings, has become the basis 565 of drug resistance surveillance surveys supported by the WHO, and has become the standard 566 for Mtbc molecular epidemiology and strain typing studies. Before its full-scale implementation, 567 we call for extensive standardisation and validation efforts. This will require political 568 commitment, and involvement of supranational laboratories and regulatory authorities. There 569 also remains an important role for the research community at large to continue to improve the 570 technical and analytical aspects of WGS. Consideration is also needed towards the ethical 571 implications and consequences of routine WGS sequencing and the information it provides. 572 There is thus a need now to commit resources to ensure access to standardized and validated 573 WGS approaches, especially in high burden countries where WGS will have the greatest impact.

574

## 575 Acknowledgements

576 CJM and LR are also affiliated with BCCM/ITM Mycobacterial Culture Collection, Institute of 577 Tropical Medicine, Antwerp, Belgium; JG is also affiliated with BC Centre for Disease Control, 578 Vancouver, Canada. SG is also affiliated with the University of Basel, Switzerland. BOA is also 579 affiliated with Center for Global Health Security and Diplomacy, Ottawa, Canada. TAK, CU, VD 580 and SN are also affiliated with the German Center for Infection Research, Partner Site Hamburg-581 Lübeck-Borstel-Riems, D-23845 Borstel, Germany. MS is also affiliated with the University of 582 Arizona, Tuscon, USA. IC is also affiliated to the CIBER in Epidemiology and Public Health, Spain. 583 Funding: 584 CJM, BO, LR and BCdJ are supported by an ERC grant [INTERRUPTB; no. 311725]. IC and GAG

585 are supported by an ERC grant (TB-ACCELERATE ; no. 638553).TR receives salary support from 586 the not-for-profit organization Foundation for Innovative New Diagnostics (FIND). The terms of 587 this arrangement have been reviewed and approved by UCSD. TMW is an NIHR Academic 588 Clinical Lecturer. JLG and JG receive funding from the University of British Columbia, Vancouver, 589 Canada. TAK, CU, VD and SN receive funding from the German Center for Infection Research 590 (DZIF) and funded by the Deutsche Forschungsgemeinschaft (DFG, German Research 591 Foundation) under Germany's Excellence Strategy – EXC 22167-390884018. LV, THH and AVR 592 are funded by FWO Oddyseus G0F8316N. MF is supported by NIH BD2K K01 (MRF ES026835). 593 PS is supported by the Agence Nationale de la Recherche (ANR-16-CE35-0009), and was a 594 consultant for Genoscreen.

- 595
- 596

## 597 Box 1 : Primary Mtbc diagnostics

Solid or liquid culture (e.g. MGIT, Beckton Dickinson, USA<sup>139</sup>) are the conventional diagnostics 598 599 for Mtbc identification and drug susceptibility testing. However such phenotypic tests can take 600 weeks to months to obtain results, require high-level biosafety infrastructure, and are 601 considered unreliable for certain drugs (e.g. pyrazinamide). Therefore, several molecular tests 602 (besides WGS) directly applicable on clinical samples have been developed. Line probe assays rely on hybridization of amplified mycobacterial DNA with nucleotide probes on strips to detect 603 604 selected drug resistance-associated mutations or their wild-type alleles. MTBDRplus<sup>140,141</sup>, TB NTM+MDR<sup>141,142</sup> and MTBDRsl<sup>143,144</sup> were endorsed by WHO. The two former assays target 605 606 mutations associated with resistance to rifampicin (in rpoB) and isoniazid (katG, inhA), i.e. detect MDR-TB. The MTBDRsl<sup>143,144</sup> assay targets mutations associated with resistance to 607 608 fluoroquinolones (qyrA, qyrB) and injectables (rrs, eis), i.e. detect XDR-TB. Other tests use cartridge-based real-time PCR (GeneXpert MTB-Rif<sup>88,145</sup> (and updated Ultra<sup>146,147</sup>); Anyplex II 609 MDR/XDR<sup>148</sup>; FluoroType MTBDR<sup>149</sup>, Hain) or PCR melt-curve (Meltpro<sup>150</sup>) for mutation 610 detection. The FluoroType as well as the WHO-endorsed and globally deployed GeneXpert both 611 612 detect rifampicin-associated mutations in rpoB, plus in the first case, isoniazid resistance 613 mutations (katG, inhA, ahpC). Because all aforementioned molecular tests use indirect 614 sequencing technologies, they are intrinsically limited to the detection of common pre-selected 615 mutations and are prone to false positive results due to indiscriminate detection of unrelated mutations<sup>151,152</sup>. To circumvent these limitations, newer assays use targeted amplicon 616 sequencing. The Next Gen-RDST<sup>153,154</sup> and Deeplex-MycTB<sup>155,156</sup> assays are directly applicable 617

618 on clinical samples and sequences (some with promoter regions) of 6 or 18 genes associated 619 with resistance to 7 or 13 anti-tuberculosis drugs, respectively. Deeplex-MycTB additionally 620 includes mycobacterial species and spoligotyping. The large coverage depths that can be

621 achieved enables high confidence mutation calls, including those born by minor subpopulations 622 in case of heteroresistance. Nevertheless, accessible targets are inherently fewer than with

- 623 WGS.
- 624

## 625 Glossary terms

- 626 *Mycobacterium tuberculosis* complex (Mtbc): the genetically related group of organisms 627 within the mycobacterium genus that cause tuberculosis in humans or animals.
- 628 Spoligotyping: a PCR-based approach based on the amplification of spacers in the CRISPR
- region of Mycobacterium tuberculosis complex. It is used for genotyping Mtbc strains.
   MIRU-VNTR: Mtbc-specific variable tandem repeats loci used to genotype Mtbc strains
- 630 MIRU-VNTR: Mtbc-specific variable tandem repeats loci used to genotype Mtbc strains
- 631 cgMLST: core genome multi-locus sequence typing; a scheme that converts genome-wide SNP
- 632 data into an allele-numbering system using a pre-selected set of core genes
- 633 wgMLST: whole genome multi-locus sequence typing; a scheme that converts genome-wide
- 634 SNP data into an allele-numbering system using a pre-selected set of core genes and 635 additional accessory genes
- 636 Löwenstein-Jensen: is a selective culture media in Mycobacteria and commonly used to 637 isolate Mtbc strains
- 638 MGIT: the Mycobacteria Growth Indicator Tube is tube that contains mycobacteria selective
- 639 culture media and which is usually coupled to automated instrument to read the results
- 640 Drug susceptibility testing: a procedure to determine if clinical isolates are resistant to
- antibiotics either by testing the inhibition in culture or by identifying drug resistanceassociated mutations
- 643 SNPs: Single nucleotide polymorphisms; differences in the nucleotide composition of a strain, 644 often compared to a reference (e.g. H37Rv).
- 645 WGS workflow: all steps involved (from culturing to SNP calling and analyses) for whole 646 genome sequencing of an isolate
- 647 WGS pipeline: the bioinformatics section of the WGS workflow, starting from fastQ files 648 through to SNP calling and analyses
- 649

# 650 Highlighted references

- World Health Organization. The use of next-generation sequencing technologies for the
   detection of mutations associated with drug resistance in Mycobacterium tuberculosis
   complex: technical guide. (2018).
- 6542.Starks, A. M. et al. Collaborative Effort for a Centralized Worldwide Tuberculosis655Relational Sequencing Data Platform. Clin. Infect. Dis. 61, S141–S146 (2015).
- 6563.The CRyPTIC Consortium and the 100000 Genomes project. Prediction of Susceptibility to657First-Line Tuberculosis Drugs by DNA Sequencing. N. Engl. J. Med. NEJMoa1800474658(2018). doi:10.1056/NEJMoa1800474
- 659 4. Coll, F. et al. Rapid determination of anti-tuberculosis drug resistance from whole-660 genome sequences. Genome Med. 7, 51 (2015).

6615.Tagliani, E. et al. EUSeqMyTB to set standards and build capacity for whole genome662sequencing for tuberculosis in the EU. Lancet Infect. Dis. 18, 377 (2018).

- 663
  6. Jajou, R. et al. Epidemiological links between tuberculosis cases identified twice as
  664 efficiently by whole genome sequencing than conventional molecular typing: A
  665 population-based study. PLoS One 13, e0195413 (2018).1
- 666 7. Coll, F. et al. A robust SNP barcode for typing Mycobacterium tuberculosis complex
  667 strains. Nat. Commun. 5, 4812 (2014).
- 6688.Satta, G. et al. Mycobacterium tuberculosis and whole-genome sequencing: how close are669we to unleashing its full potential? Clin. Microbiol. Infect. 24, 604–609 (2018).
- 670 9. Crisan, A., McKee, G., Munzner, T. & Gardy, J. L. Evidence-Based Design and Evaluation of
  671 a Whole Genome Sequencing Clinical Report for the Reference Microbiology Laboratory.
  672 doi.org 199570 (2017). doi:10.1101/199570
- 67310. Hatherell, H.-A. et al. Interpreting whole genome sequencing for investigating674tuberculosis transmission: a systematic review. BMC Med. 14, 21 (2016).

# 675676 References

- 677678 1. WHO. Global tuberculosis report 2018. (2018).
- The CRyPTIC Consortium and the 100000 Genomes project. Prediction of Susceptibility to
  First-Line Tuberculosis Drugs by DNA Sequencing. *N. Engl. J. Med.* NEJMoa1800474
  (2018). doi:10.1056/NEJMoa1800474
- 682 3. Gardy, J. L. *et al.* Whole-Genome Sequencing and Social-Network Analysis of a
  683 Tuberculosis Outbreak. *N. Engl. J. Med.* 364, 730–739 (2011).
- 6844.Cole, S. T. *et al.* Deciphering the biology of Mycobacterium tuberculosis from the685complete genome sequence. *Nature* **393**, 537–544 (1998).
- 686 5. Cabibbe, A. M., Walker, T. M., Niemann, S. & Cirillo, D. M. Whole genome sequencing of
  687 Mycobacterium tuberculosis. *Eur. Respir. J.* 52, 1801163 (2018).
- 6886.Satta, G. *et al.* Mycobacterium tuberculosis and whole-genome sequencing: how close689are we to unleashing its full potential? *Clin. Microbiol. Infect.* **24,** 604–609 (2018).
- 690 7. Lipworth, S. *et al.* SNP-IT Tool for Identifying Subspecies and Associated Lineages of
  691 Mycobacterium tuberculosis Complex. *Emerg. Infect. Dis.* 25, 482–488 (2019).
- 6928.Coll, F. *et al.* A robust SNP barcode for typing Mycobacterium tuberculosis complex693strains. *Nat. Commun.* **5,** 4812 (2014).
- 6949.Homolka, S. *et al.* High Resolution Discrimination of Clinical Mycobacterium tuberculosis695Complex Strains Based on Single Nucleotide Polymorphisms. *PLoS One* 7, e39855 (2012).
- 69610.Trauner, A. *et al.* The within-host population dynamics of Mycobacterium tuberculosis697vary with treatment efficacy. *Genome Biol.* 18, 71 (2017).
- Merker, M., Kohl, T. A., Niemann, S. & Supply, P. in *Advances in experimental medicine and biology* **1019**, 43–78 (2017).
- Jajou, R. *et al.* Epidemiological links between tuberculosis cases identified twice as
   efficiently by whole genome sequencing than conventional molecular typing: A
   population-based study. *PLoS One* 13, e0195413 (2018).
- Wyllie, D. H. *et al.* A Quantitative Evaluation of MIRU-VNTR Typing Against Whole Genome Sequencing for Identifying Mycobacterium tuberculosis Transmission: A

705		Prospective Observational Cohort Study. EBioMedicine 34, 122–130 (2018).
706	14.	Walker, T. M. et al. A cluster of multidrug-resistant Mycobacterium tuberculosis among
707		patients arriving in Europe from the Horn of Africa: a molecular epidemiological study.
708		Lancet Infect. Dis. (2018). doi:10.1016/S1473-3099(18)30004-5
709	15.	Tagliani, E. <i>et al.</i> EUSeqMyTB to set standards and build capacity for whole genome
710		sequencing for tuberculosis in the EU. <i>Lancet Infect. Dis.</i> <b>18,</b> 377 (2018).
711	16.	Cohen, K. A. et al. Evolution of Extensively Drug-Resistant Tuberculosis over Four
712		Decades: Whole Genome Sequencing and Dating Analysis of Mycobacterium tuberculosis
713		Isolates from KwaZulu-Natal. PLoS Med. 12, e1001880 (2015).
714	17.	Eldholm, V. et al. Four decades of transmission of a multidrug-resistant Mycobacterium
715		tuberculosis outbreak strain. <i>Nat. Commun.</i> <b>6,</b> 7119 (2015).
716	18.	Merker, M. et al. Evolutionary history and global spread of the Mycobacterium
717		tuberculosis Beijing lineage. <i>Nat. Genet.</i> <b>47,</b> 242–249 (2015).
718	19.	Zignol, M. et al. Genetic sequencing for surveillance of drug resistance in tuberculosis in
719		highly endemic countries: a multi-country population-based surveillance study. Lancet
720		Infect. Dis. (2018). doi:10.1016/S1473-3099(18)30073-2
721	20.	Gröschel, M. I. et al. Pathogen-based precision medicine for drug-resistant tuberculosis.
722		<i>PLOS Pathog.</i> <b>14,</b> e1007297 (2018).
723	21.	Anthony, R., Kamst, M., Nikolayevskyy, V. & van Soolingen, D. External Quality
724		Assessment of Mycobacterium Interspersed Repetitive Units - Variable Number of
725		Tandem Repeats (MIRU-VNTR) typing and Whole Genome Sequencing analysis of
726		Mycobacterium tuberculosis complex isolates across the European Reference Laboratory .
727		(2018).
728	22.	World Health Organization. The use of next-generation sequencing technologies for the
729		detection of mutations associated with drug resistance in Mycobacterium tuberculosis
730		complex: technical guide. (2018).
731	23.	Nebenzahl-Guimaraes, H., Jacobson, K. R., Farhat, M. R. & Murray, M. B. Systematic
732		review of allelic exchange experiments aimed at identifying mutations that confer drug
733		resistance in Mycobacterium tuberculosis. J. Antimicrob. Chemother. 69, 331–42 (2014).
734	24.	Sandgren, A. et al. Tuberculosis Drug Resistance Mutation Database. PLoS Med. 6,
735		e1000002 (2009).
736	25.	Coll, F. et al. Rapid determination of anti-tuberculosis drug resistance from whole-
737		genome sequences. Genome Med. 7, 51 (2015).
738	26.	Miotto, P. et al. A standardised method for interpreting the association between
739		mutations and phenotypic drug resistance in Mycobacterium tuberculosis. Eur. Respir. J.
740		<b>50,</b> 1701354 (2017).
741	27.	Starks, A. M. et al. Collaborative Effort for a Centralized Worldwide Tuberculosis
742		Relational Sequencing Data Platform. Clin. Infect. Dis. 61, S141–S146 (2015).
743	28.	Brown, T., Nikolayevskyy, V., Velji, P. & Drobniewski, F. Associations between
744		Mycobacterium tuberculosis strains and phenotypes. Emerg. Infect. Dis. 16, 272–80
745		(2010).
746	29.	Comas, I. et al. Out-of-Africa migration and Neolithic coexpansion of Mycobacterium
747		tuberculosis with modern humans. Nat. Genet. 45, 1176–82 (2013).
748	30.	Meehan, C. J. et al. The relationship between transmission time and clustering methods

749		in Mycobacterium tuberculosis epidemiology. <i>EBioMedicine</i> <b>37,</b> 410–416 (2018).
750	31.	Kohl, T. A. et al. Harmonized Genome Wide Typing of Tubercle Bacilli Using a Web-Based
751		Gene-By-Gene Nomenclature System. <i>EBioMedicine</i> <b>0,</b> (2018).
752	32.	Walker, T. M. et al. Whole-genome sequencing to delineate Mycobacterium tuberculosis
753		outbreaks: a retrospective observational study. Lancet Infect. Dis. 13, 137–46 (2013).
754	33.	Koster, K. J. et al. Genomic sequencing is required for identification of tuberculosis
755		transmission in Hawaii. <i>BMC Infect. Dis.</i> <b>18,</b> 608 (2018).
756	34.	Kohl, T. A. et al. MTBseq: a comprehensive pipeline for whole genome sequence analysis
757		of Mycobacterium tuberculosis complex isolates. PeerJ 6, e5895 (2018).
758	35.	Ezewudo, M. et al. Integrating standardized whole genome sequence analysis with a
759		global Mycobacterium tuberculosis antibiotic resistance knowledgebase. Sci. Rep. 8,
760		15382 (2018).
761	36.	Brynildsrud, O. B. et al. Global expansion of Mycobacterium tuberculosis lineage 4 shaped
762		by colonial migration and local adaptation. Sci. Adv. 4, eaat5869 (2018).
763	37.	Brown, A. C. et al. Rapid Whole-Genome Sequencing of Mycobacterium tuberculosis
764		Isolates Directly from Clinical Samples. J. Clin. Microbiol. 53, 2230–2237 (2015).
765	38.	Conceição, E. C. et al. Analysis of potential household transmission events of tuberculosis
766		in the city of Belem, Brazil. Tuberculosis 113, 125–129 (2018).
767	39.	Walker, T. M. et al. Whole-genome sequencing for prediction of Mycobacterium
768		tuberculosis drug susceptibility and resistance: a retrospective cohort study. Lancet
769		Infect. Dis. <b>15,</b> 1193–1202 (2015).
770	40.	Goig, G. A., Blanco, S., Garcia-Basteiro, A. & Comas, I. Pervasive contaminations in
771		sequencing experiments are a major source of false genetic variability: a Mycobacterium
772		tuberculosis meta-analysis. <i>bioRxiv</i> 403824 (2018). doi:10.1101/403824
773	41.	Menardo, F. et al. Treemmer: a tool to reduce large phylogenetic datasets with minimal
774		loss of diversity. BMC Bioinformatics 19, 164 (2018).
775	42.	Bryant, J. M. et al. Inferring patient to patient transmission of Mycobacterium
776		tuberculosis from whole genome sequencing data. BMC Infect. Dis. 13, 110 (2013).
777	43.	Shea, J. et al. Comprehensive Whole-Genome Sequencing and Reporting of Drug
778		Resistance Profiles on Clinical Cases of Mycobacterium tuberculosis in New York State. J.
779		<i>Clin. Microbiol.</i> <b>55,</b> 1871–1882 (2017).
780	44.	Phelan, J. et al. Mycobacterium tuberculosis whole genome sequencing and protein
781		structure modelling provides insights into anti-tuberculosis drug resistance. BMC Med.
782		<b>14,</b> 31 (2016).
783	45.	Witney, A. A. et al. Use of whole-genome sequencing to distinguish relapse from
784		reinfection in a completed tuberculosis clinical trial. BMC Med. 15, 71 (2017).
785	46.	Casali, N. et al. Whole Genome Sequence Analysis of a Large Isoniazid-Resistant
786		Tuberculosis Outbreak in London: A Retrospective Observational Study. PLOS Med. 13,
787		e1002137 (2016).
788	47.	Feuerriegel, S. et al. PhyResSE: a Web Tool Delineating Mycobacterium tuberculosis
789		Antibiotic Resistance and Lineage from Whole-Genome Sequencing Data. J. Clin.
790		Microbiol. <b>53,</b> 1908–1914 (2015).
791	48.	Bradley, P. et al. Rapid antibiotic-resistance predictions from genome sequence data for
792		Staphylococcus aureus and Mycobacterium tuberculosis. Nat. Commun. 6, 10063 (2015).

793	49.	Iwai, H., Kato-Miyazawa, M., Kirikae, T. & Miyoshi-Akiyama, T. CASTB (the comprehensive
794		analysis server for the Mycobacterium tuberculosis complex): A publicly accessible web
795		server for epidemiological analyses, drug-resistance prediction and phylogenetic
796		comparison of clinical isolates. <i>Tuberculosis (Edinb).</i> <b>95,</b> 843–844 (2015).
797	50.	Steiner, A., Stucki, D., Coscolla, M., Borrell, S. & Gagneux, S. KvarQ: targeted and direct
798		variant calling from fastg reads of bacterial genomes. BMC Genomics 15, 881 (2014).
799	51.	Farhat, M. <i>et al.</i> genTB: Translational Genomics of Tuberculosis. (2015).
800	52.	Schleusener, V., Köser, C. U., Beckert, P., Niemann, S. & Feuerriegel, S. Mycobacterium
801		tuberculosis resistance prediction and lineage classification from genome sequencing:
802		comparison of automated analysis tools. <i>Sci. Rep.</i> <b>7.</b> 46327 (2017).
803	53.	Ngo, TM. & Teo, YY. Genomic prediction of tuberculosis drug-resistance:
804		benchmarking existing databases and prediction algorithms. <i>BMC Bioinformatics</i> <b>20.</b> 68
805		(2019).
806	54.	Phelan, J. <i>et al.</i> The variability and reproducibility of whole genome sequencing
807	-	technology for detecting resistance to anti-tuberculous drugs. <i>Genome Med.</i> 8. 132
808		(2016).
809	55.	Macedo, R. <i>et al.</i> Dissecting whole-genome sequencing-based online tools for predicting
810		resistance in Mycobacterium tuberculosis: can we use them for clinical decision
811		guidance? <i>Tuberculosis</i> <b>110.</b> 44–51 (2018).
812	56.	Angers-Loustau, A. <i>et al.</i> The challenges of designing a benchmark strategy for
813		bioinformatics pipelines in the identification of antimicrobial resistance determinants
814		using next generation sequencing technologies. <i>F1000Research</i> <b>7</b> , (2018).
815	57.	FDA. Infectious Disease Next Generation Sequencing Based Diagnostic Devices: Microbial
816		Identification and Detection of Antimicrobial Resistance and Virulence Markers. (2016).
817	58.	Pouseele, H. & Supply, P. Accurate Whole-Genome Sequencing-Based Epidemiological
818		Surveillance of Mycobacterium Tuberculosis. <i>Methods Microbiol.</i> <b>42,</b> 359–394 (2015).
819	59.	Simonyan, V., Goecks, J. & Mazumder, R. Biocompute Objects-A Step towards Evaluation
820		and Validation of Biomedical Scientific Computations. PDA J. Pharm. Sci. Technol. 71,
821		136–146 (2017).
822	60.	Alterovitz, G. et al. Enabling precision medicine via standard communication of HTS
823		provenance, analysis, and results. PLOS Biol. 16, e3000099 (2018).
824	61.	Stucki, D. et al. Standard Genotyping Overestimates Transmission of Mycobacterium
825		tuberculosis among Immigrants in a Low-Incidence Country. J. Clin. Microbiol. 54, 1862–
826		70 (2016).
827	62.	Liu, Q. et al. China's tuberculosis epidemic stems from historical expansion of four strains
828		of Mycobacterium tuberculosis. <i>Nat. Ecol. Evol.</i> <b>2,</b> 1982–1992 (2018).
829	63.	Holt, K. E. et al. Frequent transmission of the Mycobacterium tuberculosis Beijing lineage
830		and positive selection for the EsxW Beijing variant in Vietnam. Nat. Genet. 50, 849–856
831		(2018).
832	64.	Coll, F. et al. Genome-wide analysis of multi- and extensively drug-resistant
833		Mycobacterium tuberculosis. <i>Nat. Genet.</i> 50, 307–316 (2018).
834	65.	Farhat, M. R. et al. Genome wide association with quantitative resistance phenotypes in
835		Mycobacterium tuberculosis reveals novel resistance genes and regulatory regions. Nat.
836		<i>Commun.</i> (2019). doi:10.1101/429159

- 83766.Kwong, J. C., Mccallum, N., Sintchenko, V. & Howden, B. P. Whole genome sequencing in838clinical and public health microbiology. *Pathology* **47**, 199–210 (2015).
- 67. Crisan, A., McKee, G., Munzner, T. & Gardy, J. L. Evidence-Based Design and Evaluation of
  a Whole Genome Sequencing Clinical Report for the Reference Microbiology Laboratory. *doi.org* 199570 (2017). doi:10.1101/199570
- Tornheim, J. A. *et al.* Building the framework for standardized clinical laboratory
  reporting of next generation sequencing data for resistance-associated mutations in *Mycobacterium tuberculosis* complex. *Clin. Infect. Dis.* (2019). doi:10.1093/cid/ciz219
- 845 69. Tan, T. W. *et al.* Advancing standards for bioinformatics activities: persistence,
  846 reproducibility, disambiguation and Minimum Information About a Bioinformatics
  847 investigation (MIABi). *BMC Genomics* **11 Suppl 4**, S27 (2010).
- Field, N. *et al.* Strengthening the Reporting of Molecular Epidemiology for Infectious
  Diseases (STROME-ID): an extension of the STROBE statement. *Lancet Infect. Dis.* 14,
  341–352 (2014).
- 851 71. World Health Organization. WHO's code of conduct for open and timely sharing of
  852 pathogen genetic sequence data during outbreaks of infectious disease. (2019).
- Allard, M. W. *et al.* Practical Value of Food Pathogen Traceability through Building a
  Whole-Genome Sequencing Network and Database. *J. Clin. Microbiol.* 54, 1975–1983
  (2016).
- 856 73. Karikari, T. K. Bioinformatics in Africa: The Rise of Ghana? *PLoS Comput. Biol.* 11,
  857 e1004308 (2015).
- Tekola-Ayele, F. & Rotimi, C. N. Translational Genomics in Low- and Middle-Income
   Countries: Opportunities and Challenges. *Public Health Genomics* 18, 242–247 (2015).
- 860 75. Helmy, M., Awad, M. & Mosa, K. A. Limited resources of genome sequencing in
  861 developing countries: Challenges and solutions. *Appl. Transl. Genomics* 9, 15–19 (2016).
- Satta, G., Atzeni, A. & McHugh, T. D. Mycobacterium tuberculosis and whole genome
  sequencing: a practical guide and online tools available for the clinical microbiologist. *Clin. Microbiol. Infect.* 23, 69–72 (2017).
- 865 77. Bradley, P. *et al.* Rapid antibiotic-resistance predictions from genome sequence data for
  866 Staphylococcus aureus and Mycobacterium tuberculosis. *Nat. Commun.* 6, 10063 (2015).
- 86778.Kurtzer, G. M., Sochat, V. & Bauer, M. W. Singularity: Scientific containers for mobility of868compute. *PLoS One* **12**, e0177459 (2017).
- 869 79. Merkel, D. Docker: lightweight linux containers for consistent development and
  870 deployment. *Linux J.* **2014**, 2 (2014).
- 87180.Grüning, B. *et al.* Bioconda: sustainable and comprehensive software distribution for the872life sciences. *Nat. Methods* **15**, 475–476 (2018).
- 87381.Jackman, S., Birol, I., Jackman, S. & Birol, I. Linuxbrew and Homebrew for cross-platform874package management. *F1000Research* 5, (2016).
- 875 82. Langille, M. G. I. & Eisen, J. A. BioTorrents: a file sharing service for scientific data. *PLoS*876 *One* 5, e10071 (2010).
- 83. Karikari, T. K., Quansah, E. & Mohamed, W. M. Y. Widening participation would be key in
  enhancing bioinformatics and genomics research in Africa. *Appl. Transl. genomics* 6, 35–
  41 (2015).
- 880 84. Bah, S. Y., Morang'a, C. M., Kengne-Ouafo, J. A., Amenga–Etego, L. & Awandare, G. A.

881		Highlights on the Application of Genomics and Bioinformatics in the Fight Against
882		Infectious Diseases: Challenges and Opportunities in Africa. Front. Genet. 9, 575 (2018).
883	85.	Zignol, M. et al. Population-based resistance of Mycobacterium tuberculosis isolates to
884		pyrazinamide and fluoroquinolones: results from a multicountry surveillance project.
885		Lancet Infect. Dis. (2016). doi:10.1016/S1473-3099(16)30190-6
886	86.	Kumwenda, S. et al. Challenges facing young African scientists in their research careers: A
887		qualitative exploratory study. <i>Malawi Med. J.</i> <b>29,</b> 1–4 (2017).
888	87.	Rabbani, F. et al. Schools of public health in low and middle-income countries: an
889		imperative investment for improving the health of populations? BMC Public Health 16,
890		941 (2016).
891	88.	Helb, D. et al. Rapid detection of Mycobacterium tuberculosis and rifampin resistance by
892		use of on-demand, near-patient technology. J. Clin. Microbiol. 48, 229–37 (2010).
893	89.	Wyllie, D. H. et al. Control of Artifactual Variation in Reported Intersample Relatedness
894		during Clinical Use of a Mycobacterium tuberculosis Sequencing Pipeline. J. Clin.
895		Microbiol. 56, e00104-18 (2018).
896	90.	Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using
897		exact alignments. <i>Genome Biol.</i> 15, R46 (2014).
898	91.	Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: rapid and sensitive
899		classification of metagenomic sequences. Genome Res. 26, 1721–1729 (2016).
900	92.	Médigue, C., Cole, S. T., Camus, JC. & Pryor, M. J. Re-annotation of the genome
901		sequence of Mycobacterium tuberculosis H37Rv. <i>Microbiology</i> <b>148</b> , 2967–2973 (2002).
902	93.	Periwal, V. et al. Comparative whole-genome analysis of clinical isolates reveals
903		characteristic architecture of Mycobacterium tuberculosis pangenome. PLoS One 10,
904		e0122979 (2015).
905	94.	Gao, Q. et al. Gene expression diversity among Mycobacterium tuberculosis clinical
906		isolates. <i>Microbiology</i> <b>151,</b> 5–14 (2005).
907	95.	Kato-Maeda, M. et al. Comparing genomes within the species Mycobacterium
908		tuberculosis. <i>Genome Res.</i> <b>11,</b> 547–54 (2001).
909	96.	Alland, D. et al. Role of large sequence polymorphisms (LSPs) in generating genomic
910		diversity among clinical isolates of Mycobacterium tuberculosis and the utility of LSPs in
911		phylogenetic analysis. J. Clin. Microbiol. 45, 39–46 (2007).
912	97.	loerger, T. R. et al. Variation among genome sequences of H37Rv strains of
913		Mycobacterium tuberculosis from multiple laboratories. <i>J. Bacteriol.</i> <b>192,</b> 3645–53
914		(2010).
915	98.	Lee, R. S. & Behr, M. A. Does Choice Matter? Reference-Based Alignment for Molecular
916		Epidemiology of Tuberculosis. J. Clin. Microbiol. 54, 1891–1895 (2016).
917	99.	Norman, A., Folkvardsen, D. B., Overballe-Petersen, S. & Lillebaek, T. Complete genome
918		sequence of Mycobacterium tuberculosis DKC2, the predominant Danish outbreak strain.
919	400	Genome Announc. <b>8,</b> e01554-18 (2019).
920	100.	Roetzer, A. et al. Whole genome sequencing versus traditional genotyping for
921		investigation of a Mycobacterium tuberculosis outbreak: a longitudinal molecular
922	404	epidemiological study. <i>PLoS Med.</i> <b>10</b> , e1001387 (2013).
923	101.	Bainomugisa, A. et al. A complete high-quality MinION nanopore assembly of an
924		extensively drug-resistant Mycobacterium tuberculosis Beijing lineage strain identifies

925 novel variation in repetitive PE/PPE gene regions. *Microb. Genomics* 4, 256719 (2018). 926 102. Igbal, Z., Caccamo, M., Turner, I., Flicek, P. & McVean, G. De novo assembly and 927 genotyping of variants using colored de Bruijn graphs. Nat. Genet. 44, 226–232 (2012). 928 103. Yadon, A. N. et al. A comprehensive characterization of PncA polymorphisms that confer 929 resistance to pyrazinamide. Nat. Commun. 8, 588 (2017). 930 104. Yang, Y. et al. Machine learning for classifying tuberculosis drug-resistance from DNA 931 sequencing data. Bioinformatics 34, 1666–1671 (2018). 932 Chen, M. L. et al. Deep learning predicts tuberculosis drug resistance status from genome 105. 933 sequencing data. bioRxiv 275628 (2018). doi:10.1101/275628 934 Rajendran, V. & Sethumadhavan, R. Drug resistance mechanism of PncA in 106. 935 Mycobacterium tuberculosis. J. Biomol. Struct. Dyn. (2013). 936 107. Kavvas, E. S. et al. Machine learning and structural analysis of Mycobacterium 937 tuberculosis pan-genome identifies genetic signatures of antibiotic resistance. Nat. 938 Commun. 9, 4306 (2018). 939 108. Duchêne, S. et al. Genome-scale rates of evolutionary change in bacteria. Microb. 940 genomics 2, e000094 (2016). 941 109. Lee, R. S. et al. Reemergence and Amplification of Tuberculosis in the Canadian Arctic. J. 942 Infect. Dis. 211, 1905–1914 (2015). 943 110. Clark, T. G. et al. Elucidating Emergence and Transmission of Multidrug-Resistant 944 Tuberculosis in Treatment Experienced Patients by Whole Genome Sequencing. PLoS One 945 8, e83012 (2013). 946 111. Guthrie, J. L. et al. Genotyping and Whole-Genome Sequencing to Identify Tuberculosis 947 Transmission to Pediatric Patients in British Columbia, Canada, 2005-2014. J. Infect. Dis. 948 **218,** 1155–1163 (2018). 949 112. Bryant, J. M. et al. Whole-genome sequencing to establish relapse or re-infection with 950 Mycobacterium tuberculosis: a retrospective observational study. Lancet. Respir. Med. 1, 951 786-92 (2013). 952 113. Guerra-Assunção, J. A. et al. Recurrence due to Relapse or Reinfection With 953 Mycobacterium tuberculosis : A Whole-Genome Sequencing Approach in a Large, 954 Population-Based Cohort With a High HIV Infection Prevalence and Active Follow-up. J. 955 Infect. Dis. 211, 1154–1163 (2015). 956 114. Schürch, A. C. et al. The tempo and mode of molecular evolution of Mycobacterium 957 tuberculosis at patient-to-patient scale. Infect. Genet. Evol. 10, 108–114 (2010). 958 Lieberman, T. D. et al. Genomic diversity in autopsy samples reveals within-host 115. 959 dissemination of HIV-associated Mycobacterium tuberculosis. Nat. Med. (2016). 960 doi:10.1038/nm.4205 961 116. Ford, C. B. et al. Mycobacterium tuberculosis mutation rate estimates from different 962 lineages predict substantial differences in the emergence of drug-resistant tuberculosis. 963 Nat. Genet. 45, 784-90 (2013). 964 117. Ford, C. B. et al. Use of whole genome sequencing to estimate the mutation rate of 965 Mycobacterium tuberculosis during latent infection. Nat. Genet. 43, 482–6 (2011). 966 Hatherell, H.-A. et al. Interpreting whole genome sequencing for investigating 118. 967 tuberculosis transmission: a systematic review. BMC Med. 14, 21 (2016). 968 119. Verver, S. et al. Transmission of tuberculosis in a high incidence urban community in

- 969 South Africa. *Int. J. Epidemiol.* **33,** 351–357 (2004).
- Bjorn-Mortensen, K. *et al.* Tracing Mycobacterium tuberculosis transmission by whole
  genome sequencing in a high incidence setting: a retrospective population-based study
  in East Greenland. *Sci. Rep.* 6, 33180 (2016).
- 973 121. Stimson, J. *et al.* Beyond the SNP Threshold: Identifying Outbreak Clusters Using Inferred
  974 Transmissions. *Mol. Biol. Evol.* **36**, 587–603 (2019).
- Biek, R., Pybus, O. G., Lloyd-Smith, J. O. & Didelot, X. Measurably evolving pathogens in
  the genomic era. *Trends Ecol. Evol.* **30**, 306–313 (2015).
- 977 123. Campbell, F. *et al.* outbreaker2: a modular platform for outbreak reconstruction. *BMC*978 *Bioinformatics* 19, 363 (2018).
- 979 124. Didelot, X., Gardy, J. & Colijn, C. Bayesian inference of infectious disease transmission
  980 from whole-genome sequence data. *Mol. Biol. Evol.* **31**, 1869–79 (2014).
- 125. Didelot, X., Fraser, C., Gardy, J. & Colijn, C. Genomic infectious disease epidemiology in
  partially sampled and ongoing outbreaks. *Mol. Biol. Evol.* 34, msw075 (2017).
- 983 126. De Maio, N., Worby, C. J., Wilson, D. J. & Stoesser, N. Bayesian reconstruction of
  984 transmission within outbreaks using genomic variants. *PLOS Comput. Biol.* 14, e1006117
  985 (2018).
- 127. Klinkenberg, D., Backer, J. A., Didelot, X., Colijn, C. & Wallinga, J. Simultaneous inference
  of phylogenetic and transmission trees in infectious disease outbreaks. *PLOS Comput. Biol.* 13, e1005495 (2017).
- 128. Kühnert, D. *et al.* Tuberculosis outbreak investigation using phylodynamic analysis.
   *Epidemics* 25, 47–53 (2018).
- 129. Eldholm, V. *et al.* Armed conflict and population displacement as drivers of the evolution
  and dispersal of Mycobacterium tuberculosis. *Proc. Natl. Acad. Sci. U. S. A.* 113, 13881–
  13886 (2016).
- 130. Streicher, E. M. *et al.* Mycobacterium tuberculosis population structure determines the
  outcome of genetics-based second-line drug resistance testing. *Antimicrob. Agents Chemother.* 56, 2420–7 (2012).
- Folkvardsen, D. B. *et al.* Rifampin heteroresistance in Mycobacterium tuberculosis
  cultures as detected by phenotypic and genotypic drug susceptibility test methods. *J. Clin. Microbiol.* 51, 4220–2 (2013).
- 1000 132. Shamputa, I. C. *et al.* Mixed infection and clonal representativeness of a single sputum
  1001 sample in tuberculosis patients from a penitentiary hospital in Georgia. *Respir. Res.* 7, 99
  1002 (2006).
- 1003133.Sobkowiak, B. *et al.* Identifying mixed Mycobacterium tuberculosis infections from whole1004genome sequence data. *BMC Genomics* **19**, 613 (2018).
- 1005 134. Gan, M., Liu, Q., Yang, C., Gao, Q. & Luo, T. Deep Whole-Genome Sequencing to Detect
   1006 Mixed Infection of Mycobacterium tuberculosis. *PLoS One* **11**, e0159029 (2016).
- 1007 135. Votintseva, A. A. *et al.* Same-Day Diagnostic and Surveillance Data for Tuberculosis via
  1008 Whole-Genome Sequencing of Direct Respiratory Samples. *J. Clin. Microbiol.* 55, 1285–
  1009 1298 (2017).
- 1010 136. Doyle, R. M. *et al.* Direct Whole-Genome Sequencing of Sputum Accurately Identifies
   1011 Drug-Resistant Mycobacterium tuberculosis Faster than MGIT Culture Sequencing. *J. Clin.* 1012 *Microbiol.* 56, e00666-18 (2018).

1013	137.	Doughty, E. L., Sergeant, M. J., Adetifa, I., Antonio, M. & Pallen, M. J. Culture-
1014		independent detection and characterisation of Mycobacterium tuberculosis and M.
1015		africanum in sputum samples using shotgun metagenomics on a benchtop sequencer.
1016		PeerJ <b>2,</b> e585 (2014).
1017	138.	Phelan, J. E. <i>et al.</i> Recombination in pe/ppe genes contributes to genetic variation in
1018		Mycobacterium tuberculosis lineages. BMC Genomics 17, 151 (2016).
1019	139.	Reisner, B. S., Gatson, A. M. & Woods, G. L. Evaluation of mycobacteria growth indicator
1020		tubes for susceptibility testing of Mycobacterium tuberculosis to isoniazid and rifampin.
1021		Diagn. Microbiol. Infect. Dis. <b>22,</b> 325–9 (1995).
1022	140.	Strydom, K. <i>et al.</i> Comparison of Three Commercial Molecular Assays for Detection of
1023		Rifampin and Isoniazid Resistance among Mycobacterium tuberculosis Isolates in a High-
1024		HIV-Prevalence Setting. J. Clin. Microbiol. 53, 3032–4 (2015).
1025	141.	Nathavitharana, R. R. <i>et al.</i> Multicenter Noninferiority Evaluation of Hain GenoType
1026		MTBDRplus Version 2 and Nipro NTM+MDRTB Line Probe Assays for Detection of
1027		Rifampin and Isoniazid Resistance. J. Clin. Microbiol. 54, 1624–1630 (2016).
1028	142.	Mitarai, S. <i>et al.</i> Comprehensive Multicenter Evaluation of a New Line Probe Assav Kit for
1029		Identification of Mycobacterium Species and Detection of Drug-Resistant Mycobacterium
1030		tuberculosis. J. Clin. Microbiol. <b>50.</b> 884–890 (2012).
1031	143.	Hillemann, D., Rüsch-Gerdes, S. & Richter, E. Feasibility of the GenoType MTBDRsl assay
1032		for fluoroquinolone, amikacin-capreomycin, and ethambutol resistance testing of
1033		Mycobacterium tuberculosis strains and clinical specimens. J. Clin. Microbiol. 47, 1767–
1034		72 (2009).
1035	144.	Tagliani, E. <i>et al.</i> Diagnostic Performance of the New Version (v2.0) of GenoType MTBDR
1036		sl Assav for Detection of Resistance to Fluoroquinolones and Second-Line Injectable
1037		Drugs: a Multicenter Study. J. Clin. Microbiol. 53. 2961–2969 (2015).
1038	145.	Ng. K. C. <i>et al.</i> Potential Application of Digitally Linked Tuberculosis Diagnostics for Real-
1039		Time Surveillance of Drug-Resistant Tuberculosis Transmission: Validation and Analysis of
1040		Test Results. JMIR Med. informatics 6. e12 (2018).
1041	146.	Chakravorty, S. <i>et al.</i> The New Xpert MTB/RIF Ultra: Improving Detection of
1042		Mycobacterium tuberculosis and Resistance to Rifampin in an Assay Suitable for Point-of-
1043		Care Testing, <i>MBio</i> <b>8</b> , e00812-17 (2017).
1044	147.	Ng. K. C. S. <i>et al.</i> Xpert Ultra Can Unambiguously Identify Specific Rifampin Resistance-
1045		Conferring Mutations, J. Clin. Microbiol. <b>56.</b> e00686-18 (2018).
1046	148	Molina-Mova B <i>et al.</i> Diagnostic accuracy study of multiplex PCB for detecting
1047	110.	tuberculosis drug resistance / Infect <b>71</b> , 220–230 (2015)
1048	149	Hillemann D. Haasis C. Andres S. Behn T. & Kranzer K. Validation of the EluoroType
1049	145.	MTBDR Assay for Detection of Rifamnin and Isoniazid Resistance in Mycobacterium
1050		tuberculosis Complex Isolates 1 Clin Microbiol 56, e00072-18 (2018)
1050	150	Pang Y et al. Rapid diagnosis of MDR and XDR tuberculosis with the MeltPro TB assay in
1051	130.	China Sci Ren 6 25330 (2016)
1052	151	Kaswa M K et al Pseudo-outbreak of pre-extensively drug-resistant (Pre-XDR)
1054	191.	tuberculosis in Kinshasa: collateral damage caused by false detection of fluoroquinolone
1055		resistance by GenoType MTBDRsl / Clin Microbiol <b>52</b> , 2876–80 (2014)
1056	152	Ailleve A et al Some Synonymous and Nonsynonymous gyrA Mutations in
1050	192.	Agine ye, A et al. Some Synonymous and Nonsynonymous gyr A Mutations in

- 1057 Mycobacterium tuberculosis Lead to Systematic False-Positive Fluoroquinolone
   1058 Resistance Results with the Hain GenoType MTBDRsl Assays. *Antimicrob. Agents* 1059 *Chemother.* 61, e02169-16 (2017).
- 1060 153. Colman, R. E. *et al.* Detection of Low-Level Mixed-Population Drug Resistance in
  1061 Mycobacterium tuberculosis Using High Fidelity Amplicon Sequencing. *PLoS One* 10,
  1062 e0126626 (2015).
- 1063 154. Colman, R. E. *et al.* Rapid Drug Susceptibility Testing of Drug-Resistant Mycobacterium
   1064 tuberculosis Isolates Directly from Clinical Samples by Use of Amplicon Sequencing: a
   1065 Proof-of-Concept Study. J. Clin. Microbiol. 54, 2058–2067 (2016).
- 1066 155. Makhado, N. A. *et al.* Outbreak of multidrug-resistant tuberculosis in South Africa
  1067 undetected by WHO-endorsed commercial tests: an observational study. *Lancet Infect.*1068 *Dis.* 18, 1350–1359 (2018).
- 1069 156. Tagliani, E. *et al.* Culture and Next-generation sequencing-based drug susceptibility
  1070 testing unveil high levels of drug-resistant-TB in Djibouti: results from the first national
  1071 survey. *Sci. Rep.* 7, 17672 (2017).
- 1072
- 1073

1074

## 1075 Display item legends

- 1076 **Figure 1:** The primary tasks for whole genome sequencing in public health. Assessing the
- 1077 epidemiology (surveillance and clustering/outbreaks) and determining the strain type or
- 1078 resistance profile to specific drugs can all be undertaken using the genomic variant calls derived
- 1079 from Mtbc WGS pipelines.
- 1080 **Figure 2:** Common workflow for whole genome sequencing for Mtbc isolates. A clinical sample
- 1081 (often sputum) is first cultured for up to 6 weeks followed by gDNA extraction and sequencing.
- 1082 The resulting sequencing output (fastq files) can be deposited online to public repositories and
- also run through standard SNP-calling pipelines which will undertake read mapping and variant
- 1084 calling. The resulting SNP lists can then be used for a variety of analyses, each which then can 1085 be reported to the end user
- 1085 be reported to the end user.
- 1086 **Figure 3:** Current and potential future approach for determining resistance-related
- 1087 polymorphisms. In the current approach (green box), lists of resistance-related SNPs are
- 1088 primarily built using a statistical approach, often a likelihood ratio. This uses linked
- 1089 phenotypic/genotypic data derived from a variety of strains across the diversity of the Mtbc to
- 1090 create lists of known SNPs that cause drug resistance. The suggested extension (blue box)
- 1091 would complement this procedure with additional information from targeted mutagenesis etc.
- 1092 to detect drug resistance causing SNPs too rare to be detected using a statistical approach.
- 1093 **Figure 4:** Epidemiological and within-host applications of SNP-based comparisons between
- 1094 Mtbc isolates. At a population level, SNP-based phylogenetics can be used to recreate local
- 1095 diversity. These phylogenies are then sub-divided into transmission clusters using pre-defined
- 1096 SNP or allele cut-offs. At the individual level, within-host diversity can be generated either
- 1097 through sub-population divergence or infection with multiple concurrent strains.
- 1098 **Supplementary table 1:** A non-exhaustive list of common bioinformatics pipelines and their
- 1099 settings for SNP calling of Mtbc isolates. This list contains only a small portion of the available
- 1100~ pipelines but demonstrates the variability and breadth of the field.
- 1101 Supplementary table 2: Suggested elements and attributes for standardised reporting of Mtbc1102 WGS result





Input data				
Strain 1 Strain 2 Strain 3 Strain				Strain 4
Lineage	4	3	1	M. bovis
rpoB	Ser450Leu	Ser450Leu	Ser450Leu	Gln429Ala
pncA	Val130Gly	Val130Gly	Arg123Gly	Gly108Ser
Dhanatunic	RIF R	RIF R	RIF R	RIF S
Phenotypic	PZA R	PZA R	PZA S	PZA R

Genotypic

Phenotypic

Mutagenesis Multi-omics Machine learning



# Extended knowledge base

Gene	Gene Change		Outcome	Support	
rpoB	Ser450Leu	RIF	R	Statistical	
rpoB	Gln429Ala	RIF	S	Mutagenesis	
pncA	Val130Gly	PZA	R	Statistical; Mutagenesis	
pncA	Arg123Gly	PZA	S	Machine learning	
pncA	Val128Gly	PZA	R	Transcriptomics	



