1 **Whole genome sequencing of *Mycobacterium tuberculosis*: current standards**
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,
91    phylogenetic analyses, quality control processes, and standardised reporting.
92

93    1.  **Introduction**

94    *Mycobacterium tuberculosis* complex (Mtbc) pathogens are collectively the top infectious
95    disease killer globally, causing 10 million new tuberculosis (TB) cases annually[1]. Increasingly,
96    new TB cases are already resistant to rifampicin and isoniazid (termed multidrug resistance;
97    MDR-TB), the key first line drugs[1]. Tackling the spread and drug resistance burden of this
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[2,3].

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

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
116  surveillance[19] and is scheduled to evaluate sequencing technologies for routine genotypic drug
117  susceptibility testing in 2019[1]. As WGS-guided individualized treatment[20] and WGS-based
118  surveillance systems[15] are being implemented in several countries (e.g. the UK and the
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
121  output formats, with little standardisation amongst them[21], making cross-comparisons and
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**
134  The standard workflow for WGS analysis of Mtbc strains is outlined in Figure 2. It involves
135  culturing sputum specimens on solid (Löwenstein–Jensen) or liquid (Mycobacteria Growth
136  Indicator Tube) media, extracting DNA from Mtbc strains, library preparation, and sequencing
137  by short read technologies (e.g. Illumina platforms)[22]. The complete Mtbc WGS analysis pipeline
138  involves several key steps such as input data validation and quality control followed by mapping
139  to a reference genome (often H37Rv) and detection of genomic variants such as single
140  nucleotide polymorphisms (SNPs) and insertion/deletions (indels). Numerous resequencing

141 pipelines for the Mtbc currently exist with currently no single 'gold standard'. These pipelines
142 typically exclude about ~10% of the genome because erroneous mapping in certain regions
143 result in false variant calls (PE/PPE gene families, other repetitive genes, mobile elements[4]) and
144 apply various criteria, such as read depth, base quality, and strand bias to filter out false
145 positive variants. Finally, based on the detected variants, several tasks can be performed
146 including (but not limited to) prediction of drug resistance and susceptibility profiles, strain
147 typing, and identification of transmission clusters.

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[23,24]. Drug
151 resistance and susceptibility profiles can be determined with high accuracy for many drugs used
152 for the treatment of TB by comparing variant calls to lists of high-confidence resistance
153 conferring variants. These lists have been established primarily using genotype-phenotype
154 associations identified from statistical analyses of large sets of clinical WGS data[25,26] (Fig. 3). A
155 prime effort in the construction of these lists is the Relational Sequencing Tuberculosis Data
156 Platform (ReSeqTB, http://www.reseqtb.org), where researchers from around the world can
157 contribute data[27]. This database contains curated, aggregated genotypic and phenotypic
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
162 TB treatment[2]. By comparing the SNPs present in a sequenced isolate to these lists, WGS can
163 not only predict resistance but also 1st line pan-susceptibility under specific conditions[2],
164 replacing the need for phenotypic testing.

166 Similarly, strain classification of the seven major human-associated lineages, many of the
167 animal-associated lineages, and their sub-lineages, can be derived directly from variant calls
168 using lists of lineage-defining SNPs[7–9]. This is important for understanding population structure
169 and potential phenotypic differences between lineages[28] and comparing isolates on the global
170 level[18,29].
171 The genomic data for a set of isolates can also be used for surveillance and transmission
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[30,31].
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
179 calibrated in low incidence settings with a diverse strain population[32]. It is not yet clear if a
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:
184 the core genome (termed cgMLST; 2891 genes covering 2.86 million bases[31]) and an extended

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

188

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

194

195  Currently, two Mtbc-specific pipelines are available, which perform multiple core tasks in single
196  install set-up to produce genetic variant calls from raw Illumina sequence data (MTBseq[34] and
197  UVP-ReSeqTB[35]). Other pathogen-agnostic pipelines can be used with an Mtbc-specific
198  reference genome and drug resistance database to achieve similar results[33,36–38]. Numerous
199  custom-built pipelines also exist[8,39–46], often incorporating similar tools for mapping and variant
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
203  resistance prediction[25,47–51] or strain typing[7,50] are available and have been comprehensively
204  compared elsewhere[52–55].

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
214  have been undertaken to great effect[22,54]. Pipeline standardisation could be achieved through
215  the use of a single pipeline in all settings or through validation with rigorous testing and
216  convergence on a defined outcome for all pipelines developed. Since multiple pipelines have
217  already been implemented (e.g. MTBseq[34] for the EUSeqMyTB consortium and the Unified
218  Variant Pipeline[35] for ReSeqTB) (supplementary table 1), agreement on validation criteria seems
219  more realistic. Since WGS-based diagnostics present a potential paradigm shift for regulatory
220  approvals, there is an urgent need to understand how to validate and standardise these
221  multiple pipelines for clinical use[56]. In 2016, the US Food and Drug Administration (FDA)
222  released draft guidelines on sequencing-based infectious disease diagnostics and bodies such as
223  the WHO and ECDC are taking steps towards international standardisations of Mtbc WGS[15,22,57].

224

225  ### a.  Technical validation and external quality control of Mtbc WGS
226  First, the extraction of DNA needs to meet minimal standards as defined for a given WGS
227  instrument[22]. Next, the pipeline to convert the raw sequencing reads into accurate variant calls
228  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
231 (e.g. from long read sequencing)[58]. Mapping these read sets to  their respective assembled
232 genomes allows to calculate the rate of false positive and negative SNPs called by the pipeline
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
235 record all thresholds, steps and implementation arguments for a given pipeline is utilised[59].
236 Comparisons of BCOs from different pipelines can then be used to set acceptable lower limits
237 for the assessed parameters, refining technical validation criteria across pipelines[60].
238 A prime example of external quality control of bioinformatics pipelines is the efforts by the
239 National Institute for Public Health and the Environment (RIVM) to standardize the use of WGS
240 for Mtbc genotyping across the European Reference Laboratory Network for TB (ERLTB-Net)[21].
241 Panels of DNA extracted from selected Mtbc isolates are sent annually by RIVM to reference
242 laboratories to assess intra- and inter laboratory reproducibility of WGS. Similar efforts in high
243 burden settings are needed to monitor the reliability of Mtbc WGS outputs when used in these
244 settings.
245
246     b.   **Validation for core tasks: transmission, phylogeny and drug resistance**
247 Task validation is used to demonstrate that a given pipeline is verified for a specific analysis,
248 e.g. drug resistance profiling. For task validation, Mtbc bioinformatics pipelines should use
249 defined validation datasets, ideally with hundreds or thousands of well characterized clinical
250 Mtbc strains representing the diversity of a specific core task (e.g. different drug susceptibility
251 profiles for resistance detection, representatives of all Mtbc lineages and (sub)-species for
252 typing, or varying degrees of clustering for transmission analyses). The number of readily
253 available, well-curated validation datasets is currently limited.
254
255 ***Validation of transmission clustering***. The national public health institute of the Netherlands
256 (RIVM) has provided laboratories with sequenced reads from 535 Mtbc isolates for which
257 epidemiological links were known. Using this dataset, the EUSeqMyTB consortium showed that
258 existing pipelines could confidently distinguish linked from unlinked cases, especially when the
259 SNP distances are high, as is often the case in low burden settings[12]. This comparison was
260 undertaken as part of an effort to standardise WGS for monitoring MDR-TB cross border
261 transmission in Europe[15].
262
263 ***Validation of classification systems.*** The clonality of Mtbc strains means that lineage and strain
264 typing can be performed using only a handful of SNPs that are specific for strains of a particular
265 lineage. Several studies have demonstrated the reliability of specific SNPs to determine the
266 Mtbc (sub)lineage[8,9,61]. However, sub-lineage classifications are often less resolved, and parallel
267 nomenclatures for lineage 2 are being used[18,62,63]. As the diversity of the Mtbc is further
268 explored, especially for animal-associated and zoonotic TB, these under-described lineages can
269 also easily be typed using the same SNP-based approach[7].
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

273 genotypic predictions, no matter the sequencing platform used[19,54]. In the past two years,
274 major progress has been made in the linkage between genotype and resistance phenotype by
275 employing a standardized statistical approach[25,26]. The task of incrementally improving our
276 knowledge base on genetic resistance profiling is primarily being addressed by the two global
277 consortia outlined above: ReSeqTB's single platform for genotype-phenotype investigation of
278 drug resistance[27,35] and CRyPTIC's genotypic-phenotypic linking of over 10,000 isolates
279 demonstrating susceptibility prediction for rifampicin and isoniazid with 99% sensitivity and 93-
280 96% for ethambutol and pyrazinamide[2]. These results have led to some low burden countries
281 (Netherlands, UK) replacing phenotypic DST with WGS-based DST for first line drugs. Resistance
282 predictions for 2$^{nd}$ line drugs can also be undertaken with sensitivity often around 90%[25]. Large
283 comparative studies using phenotype-genotype associations are expanding the catalogues[64,65]
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

288     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.
291 While the need for plain language reporting of genomic results has been recognized[52,66], there
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
298 attention to important elements with shading, bolding, and other types of emphasis, and
299 incorporating summary statements to rapidly communicate key results[67,68].

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
308 About a Bioinformatics investigation)[69] and the STROME-ID (Strengthening the Reporting of
309 Molecular Epidemiology for Infectious Diseases)[70] guidelines. In supplementary table 2, we
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[14]. Data sharing should encompass
317  data produced by research and collected in public health laboratories and surveillance efforts[71],
318  similar to the GenomeTrakr network for foodborne pathogens[72], 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

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

325

326  2.  **Implementation of WGS in routine clinical practice**
327  While the use of WGS is rapidly expanding in research, minimal progress has been made in
328  programmatic use of WGS. Some reasons include the lack of standardised end-to-end solutions,
329  the required wet-lab and computing infrastructure, need for sufficient internet connectivity and
330  bandwidth, and training deficits in genomics and bioinformatics[73–75]. Efforts are thus needed to
331  expand accessibility to perform analysis by non-experts. How these factors are addressed will
332  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
336  analysis of WGS will require little to no bioinformatics knowledge by the end user.
337  Implementation of these pipelines can be undertaken by either local set-ups with supporting
338  infrastructure or a cloud/web-based approach with easy, affordable access[76]. Many large
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
343  install additional software tools, as is done with certain existing pipelines[34,77].

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
350  recommendations detailed by expert consensus here and elsewhere[6,76]. In order to evaluate
351  new pipelines it is preferable to develop inside 'containers', such as Docker or Singularity[78,79], or
352  one-command installation wrappers like Bioconda or Homebrew[80,81]. Creating a container for
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
363 similar to the roll-out of qPCR systems (Box 1). Centralized web-based analysis platforms have
364 recently emerged and promise to aid in computational efficiency, access and usability[47,51]. Roll-
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
367 limited bandwidth, although using methods that can effectively handle connection
368 interruptions, such as BioTorrents[82], or direct transfer from sequencing centres to cloud storage
369 and/or web-based pipelines may help circumvent these issues.

370

371 The use of end-to-end, cloud-based solutions is likely to play an important role in LMICs. It is,
372 however, advisable to build in those countries human capacity for WGS of Mtbc strains[83,84].
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
380 play an important coordinating role, as is currently done for phenotypic workflows[19,85].
381 Ultimately, expanding education curricula to include bioinformatics are needed to generate
382 sufficient capacity[86].

383

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

388
389

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

391

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

395

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

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

403   implement contamination-proof analysis pipelines[40], but requires a taxonomic classification of
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
406   requires comprehensive genome databases, often making their implementation extremely
407   memory consuming[90,91]. Additionally, elimination of reads from highly conserved core bacterial
408   genes of heterologous sources still remains a problem. Proposed alternatives include masking
409   genomic regions known to accumulate artefactual polymorphisms[89], filtering the alignments
410   produced by contaminant reads, or fine-tuning the read aligners such that only the Mtbc strains
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.

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
418   comparable and standard variant calling. While most pipelines use the H37Rv genome[4,92] as the
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
421   differences between lineages[93–96]. Additionally, H37Rv contains many variants not found in any
422   other strain[97], including in genes related to drug resistance (e.g. *gyrA*S95T), creating confusion
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
427   lineage[98].

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
430   between lineages, often affecting genes involved in pathogenesis[93–96]. While these differences
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.

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

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

445    reference genome[99–101] could be used as reference to reduce mapping errors[10]. A disadvantage
446    of this approach is that it hampers comparison of results between pipelines and standardized
447    reporting 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[102].
450
451    Independent of the selection of the reference genome, other steps such as mapping and
452    filtering are not consistent between different pipelines, yet might greatly affect the analysis
453    outcome. For instance, removal of duplicates, both PCR and optical, may have a large impact in
454    the variants identified and the allele frequencies. Similarly, local assembly/realignment around
455    indels, reducing false positive SNPs derived from mapping artefacts, is rarely used in Mtbc WGS
456    pipelines[58] but is known to affect variant calling[47]. The question of whether these steps have a
457    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
469    approach could be complemented by experimental data, comprehensive single nucleotide
470    mutagenesis[103] followed by systematic phenotypic screening, multi-omics studies, and machine
471    learning approaches to predict the resistance phenotype of uncommon or novel genomic
472    variants[104,105]. With the final aim of replacing the majority of phenotypic DST by sequence-
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[2]. New statistical approaches like
475    large-scale GWAS[64,65], protein structure modelling[44,106] and machine learning[104,105,107] will likely
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
485    Accurate variant calling has major implications on downstream interpretation of the results for
486    evolutionary, epidemiological and clinical applications. Because of the low level diversity and
487    the slow substitution rate of Mtbc genomes[32,42,100,108], a few falsely called SNPs can affect the

488  interpretation of transmission events, impact the classification of a second episode of TB as
489  relapse versus re-infection, or influence the interpretation of sub-populations within a patient
490  (Fig. 4).

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

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

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
523  microevolution leading to closely related subpopulations or heteroresistance (subpopulations
524  that differ in drug resistance-related variants)[10,115,130]. Proposed sub-population detection limits
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
527  least 1% resistant bacilli is considered clinically relevant[131] , the chain reaction of selection bias
528  means that what is observed in sequencing data may not be representative of what is present
529  in the culture isolate, which in turn is likely not representative of the diversity in the sputum

530  sample, which is known to not represent the entirety of the within-patient diversity[115,132].
531  Mathematical modelling approaches have been developed to identify mixed infections[133,134].
532  However with the current approaches the detection of mixed infections is limited by the
533  relative ratio of the two strains and the number of differing SNPs between both. Future
534  research and methodological improvements are needed to better understand and interpret this
535  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,
544  preventing sufficient coverage depth of the Mtbc genomes and thus reliable variant calling,
545  even in samples with high bacterial loads[135–137]. Improvements in cell lysis or capture coupled
546  with selective DNA enrichment or depletion could reduce this technical complexity and cost.
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
555  variants in repeat-rich genome regions (e.g. PE/PPE family genes) or structural changes
556  (duplications, large indels, etc.) by short read mapping[112,138], long read sequencing will make
557  this more robust[101,135]. Unfortunately, application of this technology is currently limited by high
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  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
595
596

597 **Box 1 : Primary Mtbc diagnostics**
598 Solid or liquid culture (e.g. MGIT, Beckton Dickinson, USA[139]) are the conventional diagnostics
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
603 rely on hybridization of amplified mycobacterial DNA with nucleotide probes on strips to detect
604 selected drug resistance-associated mutations or their wild-type alleles. MTBDRplus[140,141], TB
605 NTM+MDR[141,142] and MTBDRsl[143,144] were endorsed by WHO. The two former assays target
606 mutations associated with resistance to rifampicin (in *rpoB*) and isoniazid (*katG*, *inhA*), i.e.
607 detect MDR-TB. The MTBDRsl[143,144] assay targets mutations associated with resistance to
608 fluoroquinolones (*gyrA*, *gyrB*) and injectables (*rrs*, *eis*), i.e. detect XDR-TB. Other tests use
609 cartridge-based real-time PCR (GeneXpert MTB-Rif[88,145] (and updated Ultra[146,147]); Anyplex II
610 MDR/XDR[148]; FluoroType MTBDR[149], Hain) or PCR melt-curve (Meltpro[150]) for mutation
611 detection. The FluoroType as well as the WHO-endorsed and globally deployed GeneXpert both
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
616 mutations[151,152]. To circumvent these limitations, newer assays use targeted amplicon
617 sequencing. The Next Gen-RDST[153,154] and Deeplex-MycTB[155,156] assays are directly applicable

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**
629  **region of Mycobacterium tuberculosis complex. It is used for genotyping 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**
641  **antibiotics either by testing the inhibition in culture or by identifying drug resistance**
642  **associated 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**
651  1.  World Health Organization. The use of next-generation sequencing technologies for the
652      detection of mutations associated with drug resistance in Mycobacterium tuberculosis
653      complex: technical guide. (2018).
654  2.  Starks, A. M. et al. Collaborative Effort for a Centralized Worldwide Tuberculosis
655      Relational Sequencing Data Platform. Clin. Infect. Dis. 61, S141–S146 (2015).
656  3.  The CRyPTIC Consortium and the 100000 Genomes project. Prediction of Susceptibility to
657      First-Line Tuberculosis Drugs by DNA Sequencing. N. Engl. J. Med. NEJMoa1800474
658      (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).

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

6.  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).1

7.  Coll, F. et al. A robust SNP barcode for typing Mycobacterium tuberculosis complex strains. Nat. Commun. 5, 4812 (2014).

8.  Satta, G. et al. Mycobacterium tuberculosis and whole-genome sequencing: how close are we to unleashing its full potential? Clin. Microbiol. Infect. 24, 604–609 (2018).

9.  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

10. Hatherell, H.-A. et al. Interpreting whole genome sequencing for investigating tuberculosis transmission: a systematic review. BMC Med. 14, 21 (2016).

**References**

1.  WHO. *Global tuberculosis report 2018*. (2018).

2.  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

3.  Gardy, J. L. *et al.* Whole-Genome Sequencing and Social-Network Analysis of a Tuberculosis Outbreak. *N. Engl. J. Med.* **364,** 730–739 (2011).

4.  Cole, S. T. *et al.* Deciphering the biology of Mycobacterium tuberculosis from the complete genome sequence. *Nature* **393,** 537–544 (1998).

5.  Cabibbe, A. M., Walker, T. M., Niemann, S. & Cirillo, D. M. Whole genome sequencing of Mycobacterium tuberculosis. *Eur. Respir. J.* **52,** 1801163 (2018).

6.  Satta, G. *et al.* Mycobacterium tuberculosis and whole-genome sequencing: how close are we to unleashing its full potential? *Clin. Microbiol. Infect.* **24,** 604–609 (2018).

7.  Lipworth, S. *et al.* SNP-IT Tool for Identifying Subspecies and Associated Lineages of Mycobacterium tuberculosis Complex. *Emerg. Infect. Dis.* **25,** 482–488 (2019).

8.  Coll, F. *et al.* A robust SNP barcode for typing Mycobacterium tuberculosis complex strains. *Nat. Commun.* **5,** 4812 (2014).

9.  Homolka, S. *et al.* High Resolution Discrimination of Clinical Mycobacterium tuberculosis Complex Strains Based on Single Nucleotide Polymorphisms. *PLoS One* **7,** e39855 (2012).

10. Trauner, A. *et al.* The within-host population dynamics of Mycobacterium tuberculosis vary with treatment efficacy. *Genome Biol.* **18,** 71 (2017).

11. Merker, M., Kohl, T. A., Niemann, S. & Supply, P. in *Advances in experimental medicine and biology* **1019,** 43–78 (2017).

12. 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).

13. 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. *et al.* EUSeqMyTB to set standards and build capacity for whole genome
710            sequencing for tuberculosis in the EU. *Lancet Infect. Dis.* **18,** 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. *Nat. Commun.* **6,** 7119 (2015).

716    18.    Merker, M. *et al.* Evolutionary history and global spread of the Mycobacterium
717            tuberculosis Beijing lineage. *Nat. Genet.* **47,** 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            *PLOS Pathog.* **14,** 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            **50,** 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. *EBioMedicine* **37,** 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. *EBioMedicine* **0,** (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. *BMC Infect. Dis.* **18,** 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.* **15,** 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. *bioRxiv* 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         *Clin. Microbiol.* **55,** 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         **14,** 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.* **53,** 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. *Tuberculosis (Edinb).* **95,** 843–844 (2015).
797   50.   Steiner, A., Stucki, D., Coscolla, M., Borrell, S. & Gagneux, S. KvarQ: targeted and direct
798         variant calling from fastq reads of bacterial genomes. *BMC Genomics* **15,** 881 (2014).
799   51.   Farhat, M. *et al.* 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. *Sci. Rep.* **7,** 46327 (2017).
803   53.   Ngo, T.-M. & Teo, Y.-Y. Genomic prediction of tuberculosis drug-resistance:
804         benchmarking existing databases and prediction algorithms. *BMC Bioinformatics* **20,** 68
805         (2019).
806   54.   Phelan, J. *et al.* The variability and reproducibility of whole genome sequencing
807         technology for detecting resistance to anti-tuberculous drugs. *Genome Med.* **8,** 132
808         (2016).
809   55.   Macedo, R. *et al.* Dissecting whole-genome sequencing-based online tools for predicting
810         resistance in Mycobacterium tuberculosis: can we use them for clinical decision
811         guidance? *Tuberculosis* **110,** 44–51 (2018).
812   56.   Angers-Loustau, A. *et al.* The challenges of designing a benchmark strategy for
813         bioinformatics pipelines in the identification of antimicrobial resistance determinants
814         using next generation sequencing technologies. *F1000Research* **7,** (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. *Methods Microbiol.* **42,** 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. *Nat. Ecol. Evol.* **2,** 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. *Nat. Genet.* **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         Commun.* (2019). doi:10.1101/429159

837   66.   Kwong, J. C., Mccallum, N., Sintchenko, V. & Howden, B. P. Whole genome sequencing in
838          clinical and public health microbiology. *Pathology* **47,** 199–210 (2015).

839   67.   Crisan, A., McKee, G., Munzner, T. & Gardy, J. L. Evidence-Based Design and Evaluation of
840          a Whole Genome Sequencing Clinical Report for the Reference Microbiology Laboratory.
841          *doi.org* 199570 (2017). doi:10.1101/199570

842   68.   Tornheim, J. A. *et al.* Building the framework for standardized clinical laboratory
843          reporting of next generation sequencing data for resistance-associated mutations in
844          *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).

848   70.   Field, N. *et al.* Strengthening the Reporting of Molecular Epidemiology for Infectious
849          Diseases (STROME-ID): an extension of the STROBE statement. *Lancet Infect. Dis.* **14,**
850          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).

853   72.   Allard, M. W. *et al.* Practical Value of Food Pathogen Traceability through Building a
854          Whole-Genome Sequencing Network and Database. *J. Clin. Microbiol.* **54,** 1975–1983
855          (2016).

856   73.   Karikari, T. K. Bioinformatics in Africa: The Rise of Ghana? *PLoS Comput. Biol.* **11,**
857          e1004308 (2015).

858   74.   Tekola-Ayele, F. & Rotimi, C. N. Translational Genomics in Low- and Middle-Income
859          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).

862   76.   Satta, G., Atzeni, A. & McHugh, T. D. Mycobacterium tuberculosis and whole genome
863          sequencing: a practical guide and online tools available for the clinical microbiologist.
864          *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).

867   78.   Kurtzer, G. M., Sochat, V. & Bauer, M. W. Singularity: Scientific containers for mobility of
868          compute. *PLoS One* **12,** e0177459 (2017).

869   79.   Merkel, D. Docker: lightweight linux containers for consistent development and
870          deployment. *Linux J.* **2014,** 2 (2014).

871   80.   Grüning, B. *et al.* Bioconda: sustainable and comprehensive software distribution for the
872          life sciences. *Nat. Methods* **15,** 475–476 (2018).

873   81.   Jackman, S., Birol, I., Jackman, S. & Birol, I. Linuxbrew and Homebrew for cross-platform
874          package 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).

877   83.   Karikari, T. K., Quansah, E. & Mohamed, W. M. Y. Widening participation would be key in
878          enhancing bioinformatics and genomics research in Africa. *Appl. Transl. genomics* **6,** 35–
879          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. *Malawi Med. J.* **29,** 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. *Genome Biol.* **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, J.-C. & Pryor, M. J. Re-annotation of the genome
901               sequence of Mycobacterium tuberculosis H37Rv. *Microbiology* **148,** 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. *Microbiology* **151,** 5–14 (2005).

907   95.    Kato-Maeda, M. *et al.* Comparing genomes within the species Mycobacterium
908               tuberculosis. *Genome Res.* **11,** 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.    Ioerger, T. R. *et al.* Variation among genome sequences of H37Rv strains of
913               Mycobacterium tuberculosis from multiple laboratories. *J. Bacteriol.* **192,** 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               *Genome Announc.* **8,** 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               epidemiological study. *PLoS Med.* **10,** 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.    Iqbal, 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     105.    Chen, M. L. *et al.* Deep learning predicts tuberculosis drug resistance status from genome
933            sequencing data. *bioRxiv* 275628 (2018). doi:10.1101/275628

934     106.    Rajendran, V. & Sethumadhavan, R. Drug resistance mechanism of PncA in
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     115.    Lieberman, T. D. *et al.* Genomic diversity in autopsy samples reveals within-host
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     118.    Hatherell, H.-A. *et al.* Interpreting whole genome sequencing for investigating
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).

970    120.    Bjorn-Mortensen, K. *et al.* Tracing Mycobacterium tuberculosis transmission by whole
971          genome sequencing in a high incidence setting: a retrospective population-based study
972          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).

975    122.    Biek, R., Pybus, O. G., Lloyd-Smith, J. O. & Didelot, X. Measurably evolving pathogens in
976          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).

981    125.    Didelot, X., Fraser, C., Gardy, J. & Colijn, C. Genomic infectious disease epidemiology in
982          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).

986    127.    Klinkenberg, D., Backer, J. A., Didelot, X., Colijn, C. & Wallinga, J. Simultaneous inference
987          of phylogenetic and transmission trees in infectious disease outbreaks. *PLOS Comput.*
988          *Biol.* **13,** e1005495 (2017).

989    128.    Kühnert, D. *et al.* Tuberculosis outbreak investigation using phylodynamic analysis.
990          *Epidemics* **25,** 47–53 (2018).

991    129.    Eldholm, V. *et al.* Armed conflict and population displacement as drivers of the evolution
992          and dispersal of Mycobacterium tuberculosis. *Proc. Natl. Acad. Sci. U. S. A.* **113,** 13881–
993          13886 (2016).

994    130.    Streicher, E. M. *et al.* Mycobacterium tuberculosis population structure determines the
995          outcome of genetics-based second-line drug resistance testing. *Antimicrob. Agents*
996          *Chemother.* **56,** 2420–7 (2012).

997    131.    Folkvardsen, D. B. *et al.* Rifampin heteroresistance in Mycobacterium tuberculosis
998          cultures as detected by phenotypic and genotypic drug susceptibility test methods. *J.*
999          *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).

1003   133.    Sobkowiak, B. *et al.* Identifying mixed Mycobacterium tuberculosis infections from whole
1004         genome 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* **2,** e585 (2014).
1017  138.  Phelan, J. E. *et al.* 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.* **22,** 325–9 (1995).
1022  140.  Strydom, K. *et al.* 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. *et al.* 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. *et al.* Comprehensive Multicenter Evaluation of a New Line Probe Assay Kit for
1029        Identification of Mycobacterium Species and Detection of Drug-Resistant Mycobacterium
1030        tuberculosis. *J. Clin. Microbiol.* **50,** 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. *et al.* Diagnostic Performance of the New Version (v2.0) of GenoType MTBDR
1036        *sl* Assay 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. *et al.* 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. *et al.* 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. *MBio* **8,** e00812-17 (2017).
1044  147.  Ng, K. C. S. *et al.* Xpert Ultra Can Unambiguously Identify Specific Rifampin Resistance-
1045        Conferring Mutations. *J. Clin. Microbiol.* **56,** e00686-18 (2018).
1046  148.  Molina-Moya, B. *et al.* Diagnostic accuracy study of multiplex PCR for detecting
1047        tuberculosis drug resistance. *J. Infect.* **71,** 220–230 (2015).
1048  149.  Hillemann, D., Haasis, C., Andres, S., Behn, T. & Kranzer, K. Validation of the FluoroType
1049        MTBDR Assay for Detection of Rifampin and Isoniazid Resistance in Mycobacterium
1050        tuberculosis Complex Isolates. *J. Clin. Microbiol.* **56,** e00072-18 (2018).
1051  150.  Pang, Y. *et al.* Rapid diagnosis of MDR and XDR tuberculosis with the MeltPro TB assay in
1052        China. *Sci. Rep.* **6,** 25330 (2016).
1053  151.  Kaswa, M. K. *et al.* Pseudo-outbreak of pre-extensively drug-resistant (Pre-XDR)
1054        tuberculosis in Kinshasa: collateral damage caused by false detection of fluoroquinolone
1055        resistance by GenoType MTBDRsl. *J. Clin. Microbiol.* **52,** 2876–80 (2014).
1056  152.  Ajileye, A. *et al.* Some Synonymous and Nonsynonymous gyrA 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

**Display item legends**

**Figure 1:** The primary tasks for whole genome sequencing in public health. Assessing the epidemiology (surveillance and clustering/outbreaks) and determining the strain type or resistance profile to specific drugs can all be undertaken using the genomic variant calls derived from Mtbc WGS pipelines.

**Figure 2:** Common workflow for whole genome sequencing for Mtbc isolates. A clinical sample (often sputum) is first cultured for up to 6 weeks followed by gDNA extraction and sequencing. 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 calling. The resulting SNP lists can then be used for a variety of analyses, each which then can be reported to the end user.

**Figure 3:** Current and potential future approach for determining resistance-related polymorphisms. In the current approach (green box), lists of resistance-related SNPs are primarily built using a statistical approach, often a likelihood ratio. This uses linked phenotypic/genotypic data derived from a variety of strains across the diversity of the Mtbc to create lists of known SNPs that cause drug resistance. The suggested extension (blue box) would complement this procedure with additional information from targeted mutagenesis etc. to detect drug resistance causing SNPs too rare to be detected using a statistical approach.

**Figure 4:** Epidemiological and within-host applications of SNP-based comparisons between Mtbc isolates. At a population level, SNP-based phylogenetics can be used to recreate local diversity. These phylogenies are then sub-divided into transmission clusters using pre-defined SNP or allele cut-offs. At the individual level, within-host diversity can be generated either through sub-population divergence or infection with multiple concurrent strains.
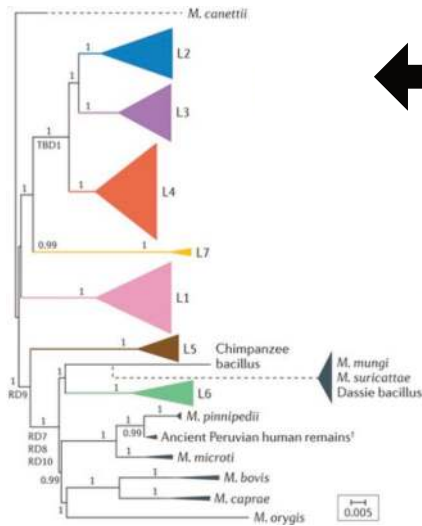
**Supplementary table 1:** A non-exhaustive list of common bioinformatics pipelines and their settings for SNP calling of Mtbc isolates. This list contains only a small portion of the available pipelines but demonstrates the variability and breadth of the field.

**Supplementary table 2:** Suggested elements and attributes for standardised reporting of Mtbc WGS result
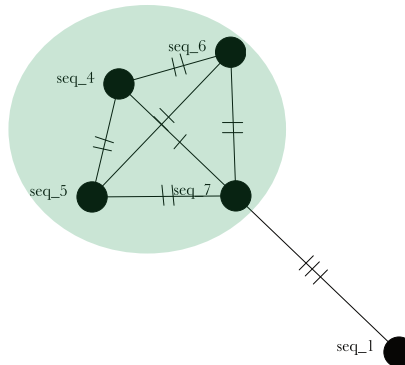
# Surveillance

**The Gambia**
MDR-TB prevalence
Clustering rate
...

**Indonesia**
MDR-TB prevalence
Clustering rate
...

**South Africa**
MDR-TB prevalence
Clustering rate
...

# Strain and subspecies typing

# Whole genome sequencing

# Drug susceptibility

| Drug | Resistance |
|------|------------|
| RIF | R |
| INH | R |
| ETH | S |
| PZA | S |

# Clustering and outbreaks

**Sample** → **Culture** → **Sequencing** → **Fastq file** → **Online deposit**

**Assembly pipeline**

**Read mapping**

**Variant Calling**
G C A
A G T

**SNP table**

| Position | SNP |
|----------|-----|
| 18253 | C |
| 23007 | T |
| ... | ... |

**Drug resistance**

| Position | SNP | Gene | Change | Drug |
|----------|-----|------|--------|------|
| 761109 | T | rpoB | Asp435Tyr | RIF |
| 2155168 | C | katG | Ser315Thr | INH |
| ... | ... | ... | ... | ... |

**Clustering**

seq_6
seq_4
seq_5
seq_3
seq_1

**Strain typing**

| Position | SNP | Gene | Lineage |
|----------|-----|------|---------|
| 5520 | T | gyrB | 4,3,2,1 |
| 1473079 | A | rrs | M. microti |
| ... | ... | ... | ... |

**Reporting**

Mycobacterium tuberculosis sequencing report
Patient:                                Cluster detection
Birth date:                             Cluster 48
...

Drug susceptibility                     Strain information
RIF        R                            Lineage 4
INH        S
...        ...

**Clonal diversification**

Minority variant

Hetero-resistance

Majority variant

Resistance

TRANSMISSION CLUSTER

- Distance threshold
- Phylogenetic modelling

**Mixed infection**

Majority variant strain 1

Hetero-resistance

Strain1

Strain2

Majority variant strain 2

Resistance