



Epidemiology

Resolving a clinical tuberculosis outbreak using palaeogenomic genome reconstruction methodologies



Rhys Jones^a, Marcela Sandoval Velasco^b, Llinos G. Harris^a, Sue Morgan^c, Mark Temple^c, Ruddy Michael C^{d,e}, Rhian Williams^d, Perry Michael D^d, Matt Hitchings^a, Thomas S. Wilkinson^a, Thomas Humphrey^a, M. Thomas P. Gilbert^{b,f}, Davies Angharad P^{a,e,*}

^a Swansea University Medical School, Institute of Life Science, Swansea University, Swansea, Wales, UK

^b Natural History Museum of Denmark, University of Copenhagen, Copenhagen K, Denmark

^c Health Protection Division (Mid and West Wales), Public Health Wales, Swansea, Wales, UK

^d Wales Centre for Mycobacteriology, Llandough Hospital, Cardiff, Wales, UK

^e Public Health Wales Microbiology, Swansea, Wales, UK

^f Norwegian University of Science and Technology, University Museum, Trondheim, Norway

ARTICLE INFO

Keywords:

Mycobacterium tuberculosis
Whole genome sequencing
Outbreak investigation
Ancient DNA library construction
Palaeogenomics

ABSTRACT

This study describes the analysis of DNA from heat-killed (boilate) isolates of *Mycobacterium tuberculosis* from two UK outbreaks where DNA was of sub-optimal quality for the standard methodologies routinely used in microbial genomics. An Illumina library construction method developed for sequencing ancient DNA was successfully used to obtain whole genome sequences, allowing analysis of the outbreak by gene-by-gene MLST, SNP mapping and phylogenetic analysis. All cases were spoligotyped to the same Haarlem H1 sub-lineage. This is the first described application of ancient DNA library construction protocols to allow whole genome sequencing of a clinical tuberculosis outbreak. Using this method it is possible to obtain epidemiologically meaningful data even when DNA is of insufficient quality for standard methods.

1. Introduction

In 2017, tuberculosis incidence in Wales was 3.4 cases per 100 000 population [1], lower than the 2017 overall UK incidence of 8.4 cases per 100 000 population [2]. However, the structure of the National Health Service Trust responsible for health protection in Wales, Public Health Wales (PHW), links microbiology, public health and epidemiology into one organizational team. This, coupled with the relatively stable population, enables detailed analysis of links between cases and makes Wales an attractive place to study tuberculosis transmission dynamics.

During 2003–2005 an outbreak of *M. tuberculosis* occurred in a small town in south Wales (town G). Seven cases (GO1–GO7) were associated with a public house in the town. Case GO3, a barman in the public house, and case GO4, both had direct contact with all the other *M. tuberculosis* cases within the outbreak, at least six of whom had visited the public house. Case GO4 was considered to have been highly infectious, having been symptomatic for about 16 months before diagnosis. Both GO3 and

GO4 were thought to be the public health team to represent potential super-spreaders within the outbreak. During this period, the standard typing method in use by Public Health Wales was MIRU-VNTR. All the cases had identical MIRU-VNTR patterns. They were all fully susceptible to standard therapy.

In 2008, another outbreak was identified in an area of a nearby town, approximately 6 km away (town T). There were close links between cases in the two outbreaks, with two of the cases in town T (TH1 and TH2) being direct contacts of GO3. TH2 was a regular at the public house at the centre of the outbreak in town G. However, MIRU-VNTR typing of the two groups of cases differed, with a polymorphism at a single MIRU-VNTR locus (MIRU16), suggesting the presence of two independent outbreaks within the area. Nonetheless the public health team felt it likely that all the cases formed part of one larger outbreak caused by the same strain of *M. tuberculosis* with divergence seen at locus MIRU16 due to a change in an endemic circulating strain. Epidemiologically linked isolates differing at fewer than two MIRU-VNTR loci have been suggested previously to be likely to be part of the same clonal complex

* Corresponding author. Swansea University Medical School, Institute of Life Science, Swansea University, Swansea, Wales, UK.

E-mail address: angharad.p.davies@swansea.ac.uk (A.P. Davies).

<https://doi.org/10.1016/j.tube.2019.101865>

Received 3 May 2019; Received in revised form 19 September 2019; Accepted 22 September 2019

Available online 23 September 2019

1472-9792/© 2019 Elsevier Ltd. All rights reserved.

within an outbreak [3] and MIRU-VNTR typing has been found to be unable to account for within-outbreak heterogeneity. It also provides limited information on the direction of transmission, identification of super-spreaders or outbreak origin [4–9].

The development of affordable and accessible whole genome sequencing (WGS) protocols based around next generation sequencing (NGS) platforms such as the Illumina series, has provided an alternative method for the investigation of *M. tuberculosis* outbreaks. The quality of the DNA sample is critical to the success of WGS. The *M. tuberculosis* samples for this study were provided by the Wales Centre for Mycobacteriology (WCM) in boilate form.

Although boilate extraction does release DNA, it is crude, inconsistent and yields DNA of lower integrity, in low quantity and of poorer quality in comparison with other extraction methods [10]. As might be expected, poor sample quality has a negative effect on the standard Nextera XT library preparation [11]. Given the need sometimes to generate WGS data from such samples, lessons might be learnt from the field of palaeogenomics, which attempts routinely to generate genome-scale data from nucleic acids that are highly fragmented, contaminated with non-target DNA, and often contain residual chemical impurities [12–14]. Recent developments in palaeogenomic methodologies have allowed WGS data to be obtained from a wide range of samples, spanning ancient humans and hominids [15,16], mammals [17,18], plants [19] and even pathogens [20] – many of which contain DNA with fragment lengths of <80bp. We therefore hypothesised that the application of palaeogenomic sequencing protocols might overcome the challenge of retrieving genomic data from the low purity and low-quality DNA of crude *M. tuberculosis* boilate samples.

2. Materials and methods

2.1. Sample collection

Boilate samples from the *M. tuberculosis* isolates described above were obtained from the WCM, Public Health Wales, Llandough Hospital, Cardiff. The isolates had been cultured using the BACTEC™ MGIT™ 960 System (Becton Dickinson Diagnostic Systems, Sparks, MD) in containment level 3 facilities and then heat-killed by boiling for 35 min at 110 °C.

MIRU-VNTR typing had been performed at the PHW Molecular Unit, University Hospital of Wales, Cardiff, with typing based on 15 loci, namely: ETRA, ETRB, ETRC, ETRD, ETRE, MIRU2, MIRU10, MIRU16, MIRU20, MIRU23, MIRU24, MIRU26, MIRU27, MIRU39 and MIRU40. Epidemiological information for each isolate was obtained through face-to-face interviews with a senior public health nurse from the original PHW outbreak investigation team, and from the outbreak documentation. All the cases in these outbreaks and under consideration here were fully susceptible to standard anti-tuberculosis chemotherapy.

2.2. Sequencing attempt using conventional illumina protocols

Sequencing was first attempted directly from the boilates, following an ethanol precipitation. Indexed genomic DNA libraries were prepared for sequencing using the Illumina Nextera XT (V3) sample preparation protocol following the manufacturer's guidelines (2017 Illumina Inc., San Diego, CA, USA), size-selecting for fragments with an average size of 500bp. Bead-normalised sequencing libraries were pooled and sequenced on a MiSeq platform (2017 Illumina Inc., San Diego, CA, USA) using the V3 reagent kits and 600 cycles. The resulting paired-end reads were quality filtered with the Trimmomatic tool software [21] using a sliding window approach of 5 bases and a quality score of Q20 prior to contig assembly using the SPAdes genome assembler (Version 3.9.0) with K-mer sizes 33, 55, 77, 99 and 127 used [22]. In each case, this method failed to generate any sequence data at all, and unfortunately in this instance, replacement samples were not available. A method optimized for sequencing degraded DNA sources was therefore

explored.

2.3. DNA extraction

The remaining *M. tuberculosis* boilates were transferred to tubes containing a 500 µL solution of digestion buffer (10 mM Tris-HCl pH8, 10 mM NaCl, 5 nM CaCl, 2.5 mM EDTA, 1% SDS, 1% Proteinase K, and DTT) and 500 µL of Phenol: Chloroform: Isoamyl alcohol solution (Sigma-Aldrich, St. Louis, MO, United States). Next, 0.6 g of Zirconia/Silica beads (Cat. No. 11079105z, Biospec Products Inc., Bartlesville, OK, USA) were added to the tubes and each sample was homogenized using a TissueLyser II (Qiagen, Valencia, California), for 4 rounds of 20 s bursts with cooling on ice for 30 s between rounds. After homogenization, samples were centrifuged for 10 min at 16 000×g in a bench centrifuge to separate the phases. The aqueous upper phase (around 500 µL) was gently transferred to a new low-bind 2 mL Eppendorf tube and two volumes (1 mL) of ice cold absolute ethanol (kept at –20C) were added to each sample. Samples were vortexed briefly and centrifuged again for 10 min at 16,000 g. The supernatant was discarded and the tubes washed carefully with 700 µL of 70% ethanol without disturbing the pellet. The ethanol wash was discarded and the pellet was left to dry for 2 min. Finally, the pellet was re-suspended and DNA eluted in buffer EB (Qiagen, Valencia, California). Extracted DNA was quantified on a Qubit fluorometer using a dsDNA high sensitivity assay (Life Technologies, Carlsbad, California) and an Agilent 2100 Bioanalyser (Santa Clara, California). Following extraction, the samples were fragmented for 20 cycles in 30 s cycles within a Diagenode bioruptor 300.

2.4. Ancient DNA library preparation protocol

Carøe et al. [14] have recently published a new library construction protocol developed specifically for use on low concentration and degraded nucleic acid extracts. Based around a single tube blunt-end adaptor ligation, this so-called 'BEST' protocol has been shown to yield more complex libraries than other methods, due to removal of intermediate purification steps that generally lead to loss of the DNA molecules within the extract [14]. Illumina compatible libraries were constructed in this way at the laboratories of the Natural History Museum of Denmark, using 32 µl of extracted (see DNA extraction above) DNA from each boilate per sample as input. Based on qPCR results, libraries were indexed through PCR amplification for 10, 12 or 15 cycles, prior to visualisation and quantification on an Agilent Bioanalyser using the High Sensitivity DNA assay (Agilent technologies, Cheshire, UK). Subsequently, the indexed libraries were pooled at equimolar concentrations and then sequenced on an Illumina HiSeq 2500 platform (Illumina sequencing platforms, 2017) in 80bp single read mode by the Danish National High-Throughput DNA Sequencing Centre. The resulting single-end reads were quality filtered with the Trimmomatic tool [21] using a sliding window approach of 5 bases and a quality score of Q20 prior to contig assembly using the SPAdes genome assembler [22]. Raw reads for all the isolates are publicly available (NCBI BioProject PRJNA556450).

2.5. Transmission chain

Cytoscape software [23] was used to generate a transmission tree.

2.6. Gene-by-gene MLST analysis

Gene-by-gene MLST analysis was carried out using Ridom SeqSphere Software [24]. A published core genome MLST (cgMLST) scheme [6,24] was used for the analysis, which was based on 2891 core genes.

2.7. Whole genome sequence SNP mapping

Single nucleotide polymorphisms (SNPs) were identified using the

standardised online CSI Phylogeny programme (Version 1.4; Call SNPs & Infer Phylogeny) of the Centre for Genomic Epidemiology (CGE) online tool [<https://cge.cbs.dtu.dk/services/CSIPhylogeny/>] [25]. A minimum spanning tree was constructed based on SNPs from 1123 sites from the WGS data using an adapted application of the Ridom SeqSphere software [24].

2.8. *In silico* spoligotype

Each isolate sequence was submitted to the Python-based SpolTyping [26] *in silico* software for prediction of spoligotype pattern. Resulting octal and binary patterns were then submitted to the SitVit database for determination of international typing assignment and assignment to globally recognised spoligotype clades [27]. Additionally, spoligotype patterns were submitted to the TB-insight online server for identification of the corresponding *M. tuberculosis* lineage [27]. Isolates were assigned to major lineages based on the Conformal Bayesian Network (CBN) parameters, which employ a hierarchical Bayesian network based on PCR based biomarkers such as spoligotypes to classify isolates into given lineages [28].

3. Results

All the isolates were sequenced successfully using the aDNA sequencing protocol (Table 1). Sequence data for all isolates covered >99% of the core genes used in the cgMLST scheme. In addition, in all cases, sequence data covered >98% of the specified reference genome according to the CGE CSI phylogeny software. The minimum spanning tree (Fig. 1) shows the genomic distances, based on allelic differences across 2891 core genes, between each of the isolates included in this dataset. A >12 allele difference was used as the threshold for exclusion from the outbreak [6].

3.1. cgMLST

Gene-by-gene cgMLST analysis appeared to indicate the presence of two outbreaks within the dataset, labelled outbreak 1 and outbreak 2, separated by 124 allelic differences. Outbreak 1 included the outbreak G public house cases GO2, GO3, GO4, GO5, GO6 and GO7 which all had fewer than 12 allelic differences between them. Within outbreak 1, two isolates could not be distinguished from one another, with no allelic differences detected (GO2 and GO6). A star-like topology is seen for outbreak 1, with GO2 and GO6 in the central position. Outbreak 2 is represented by three isolates, two from outbreak TH (TH1 and TH2) and from outbreak G, GO1 (Fig. 1a). TH1 and GO1 could not be distinguished. Isolate TH2 diverged from GO1 and TH1 by only 4 allelic differences, well within the 12-allele limit and thus representing a direct transmission event [6]. Isolates GO8 and GO9, also from town G and so included in the figure for comparison, substantially exceeded the 12-allelic difference threshold for direct transmission with any other isolate within this dataset, consistent with the results of MIRU-VNTR

typing, which previously found that GO8 and GO9 were not outbreak-related strains.

3.2. SNP mapping

SNP mapping also highlighted the presence of the same two outbreaks within this dataset (Fig. 1b). For outbreak 1, a star-like topology was again present including isolates GO2, GO3, GO4, GO5, GO6 and GO7, as seen in the cgMLST analysis in Fig. 1a. However, SNP mapping was able to distinguish GO2 and GO6 and found GO3 to be the central isolate. Outbreak 2 contained only two isolates, GO1 and TH1, with TH2 being excluded due to exceeding the threshold of 12 SNPs for outbreak inclusion. SNP mapping supported cgMLST and MIRU-VNTR in excluding GO8 and GO9 from either outbreak.

3.3. *In silico* spoligotyping

All the isolates were successfully predicted a spoligotype *in silico* (Table 2). The isolates could all be assigned to the same lineage and correlating spoligotype clades, the Euro-American lineage and Haarlem H1 spoligotype clades respectively. Three different spoligotypes and correlating international types were found: 47, 46 and 742. Isolates GO2 and GO4 had a different international type (742) from isolates GO3, GO5, GO6 and GO7 (international type 47) despite being directly linked within the same outbreak by both cgMLST and SNP mapping. In line with WGS, GO1 showed a closer association with the TH outbreak cases compared to the other town G cases, having the same spoligotype pattern (international type 46).

4. Discussion

This study describes the first application of an aDNA library construction protocol for the investigation of a clinical outbreak of *M. tuberculosis*. The 'BEST' ancient DNA library preparation protocol and subsequent sequencing were able to provide extensive sequence data from sub-optimal quality DNA from *M. tuberculosis* outbreak samples, where the standard protocol had failed. The aDNA library preparation protocol was able to circumvent the issue of short DNA fragment sizes and yield WGS data from *M. tuberculosis* outbreak samples that would otherwise have been lost.

WGS provided extensive information on the outbreak. The application of gene-by-gene cgMLST provided similar, but not identical, results to those achieved by the SNP mapping procedure. Both agreed in assigning two separate outbreaks and both supported the epidemiological suspicion that a super-spreader was present, although the cgMLST and SNP mapping analyses conflicted in their assignment of the likely super-spreader. cgMLST indicated that the super-spreader was either GO2 or GO6, while traditional SNP mapping supported the assumption of the public health team that GO3 was the super-spreader responsible for multiple secondary cases within this outbreak. A possible reason for the discrepancy between the two methods could be the presence of gene

Table 1

Ancient DNA protocol sequencing of nine outbreak isolates. The percentage of the core genome MLST genes present is shown as well as the percentage of the reference *M. tuberculosis* H37Rv genome covered according to the CSI phylogeny algorithm provided by the Centre for Genomic Epidemiology [7].

Sample ID	No. contigs	Largest contig	Total length	N50	Average depth of coverage	% cgMLST	% reference genome covered	GenBank Accession	% Error rate
TH1	151	174895	4372963	64226	177.0	99.52	98.27	VOGE000000000	0.45
TH2	181	171547	4347932	61791	86.9	99.34	98.18	VOGF000000000	0.46
GO1	150	174836	4360310	69175	159.0	99.52	98.45	VOGD000000000	0.44
GO2	163	174750	4353499	63534	90.1	99.52	98.07	VOGG000000000	0.44
GO3	162	211047	4401955	75381	221.2	99.52	98.26	VOGH000000000	0.45
GO4	143	257745	4357032	69538	129.0	99.41	98.15	VOGI000000000	0.42
GO5	134	228152	4363387	76515	158.8	99.45	98.50	VOGJ000000000	0.43
GO6	148	210603	4362517	72538	125.0	99.48	98.36	VOGK000000000	0.40
GO7	161	174721	4349894	64061	76.3	99.52	98.50	VOGL000000000	0.41

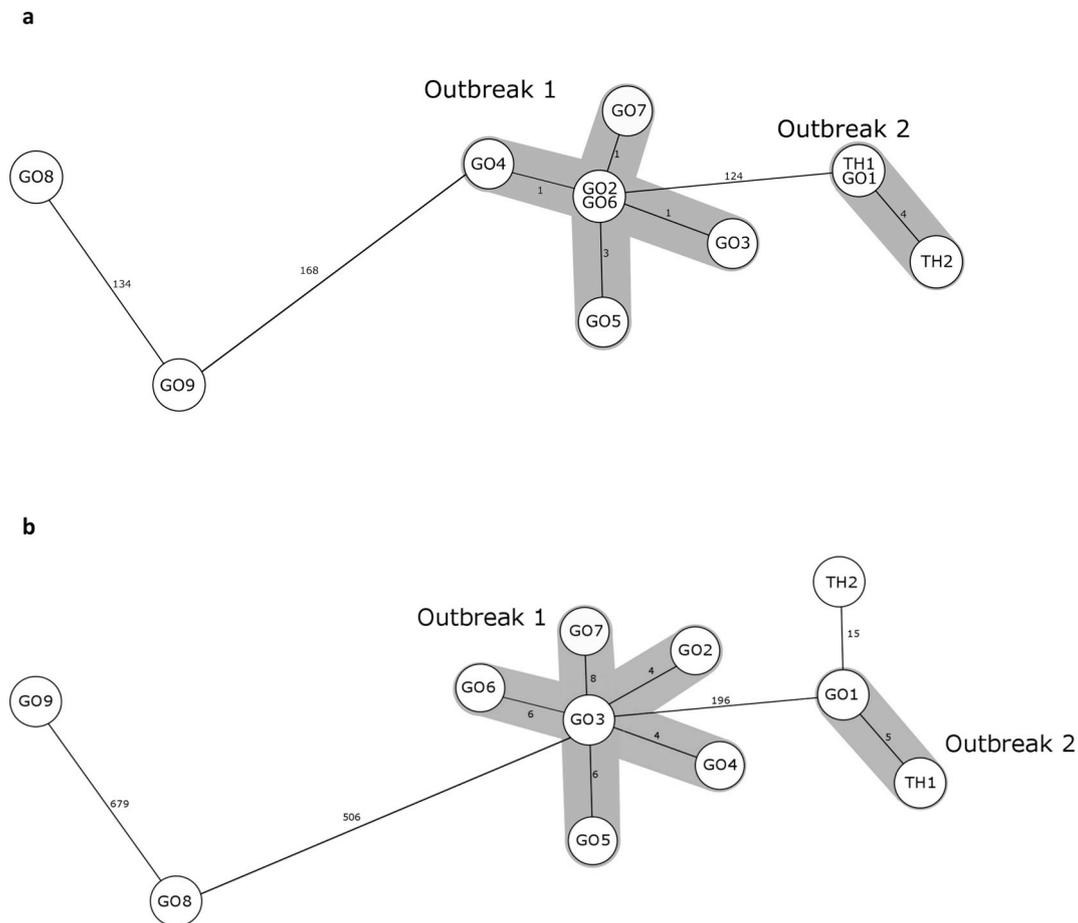


Fig. 1. Results from analysing the nine isolates, with the addition of 2 isolates (GO8 & GO9) from the same geographical area isolated during the same time-period, previously sequenced by a standard method. a: cgMLST minimum spanning tree of the 11 isolates with numbers representing the number of allelic differences between isolates; and b: CSI phylogeny based minimum spanning tree of the 11 isolates on a total of 1123 SNPs, with numbers representing the number of SNPs between isolates. Branches between isolates are not to scale. Shaded areas represent those associated with the labelled outbreak. Isolates with prefix TH were from Town T; prefix GO indicates isolates from Town G.

Table 2

in silico spoligotyping results for each of the outbreak isolates. Results include the predicted spoligotype (produced by SpolDB4), international spoligotype (SITVIT database), lineage assignment and clade assignment (both outputted from the TB-insight online server).

Isolate	Predicted Spoligotype	International type	Lineage	Clade
GO3	77777774020771	47	Euro-American	H1
TH1	77777770000000	46	Euro-American	H1
TH2	77777770000000	46	Euro-American	H1
GO1	77777770000000	46	Euro-American	H1
GO2	77777770020771	742	Euro-American	H1
GO4	77777770020771	742	Euro-American	H1
GO5	77777774020771	47	Euro-American	H1
GO6	77777774020771	47	Euro-American	H1
GO7	77777774020771	47	Euro-American	H1

families of a repetitive nature being included in the analyses, such as those for PE_PPE which show disproportionately high amount of divergence [29]. cgMLST removed these regions and this could at least partly explain the discrepancies. It is important to note that most SNP-calling procedures in *M. tuberculosis* epidemiology filter out repetitive regions, and that therefore the procedure of mapping and SNP-calling without filtering, described here, is not directly comparable with that described in other studies.

In silico application of SpoTyping software was achieved without further laboratory work and at no extra cost to initial sequencing data

[26]. Each isolate in this outbreak had a spoligotype pattern corresponding to the Haarlem H1 clade. It has been demonstrated previously that the value of inferring a recent common ancestor of isolates within potentially related outbreaks, with identification of a causative circulating strain being a common feature [8,30]. Previous studies of *M. tuberculosis* in the Inuit and Greenland populations of North America, which have stable populations, have also documented the cause of multiple outbreaks within the region as the ongoing spread of an evolving founder strain that has continually spread across the area over decades [8,30].

The results of *in silico* spoligotyping, together with the strong epidemiological links between the two outbreaks, lends support to the public health team’s case that both G and T outbreak cases were part of the same on-going outbreak. The apparent existence of two outbreaks may be due to the absence of intermediate isolates not included in this dataset coupled with recent minor genomic changes [8], a conclusion consistent with the presence of a polymorphism at only one MIRU-VNTR locus between the two genotypes. Such problems with MIRU-VNTR typing have been described before [4,8,30]. Approximately one-third of reported tuberculosis cases are culture-negative, so that there is no isolate for analysis.

The work described here demonstrates that when necessary, clinically useful data can be obtained from sub-optimal quality samples by applying an ancient DNA library construction protocol to overcome the need for DNA of high purity and quality. The success of the ‘BEST’ aDNA library construction protocol here highlights a clinical application for a

method previously associated only with palaeogenomic studies, and shows how the transfer of such techniques in defined circumstances could provide clinical benefit.

Acknowledgments

This work was funded by St. David's Medical Foundation & Coleg Cenedlaethol Cymraeg funding to APD & RJ, a EMBO Short-term Fellowship to RJ, and an ERC Consolidator Grant (681396- Extinction Genomics) to MTPG. The authors would like to thank Christian Carøe and other members of the Gilbert laboratory at the Natural History Museum, Denmark for technical guidance.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.tube.2019.101865>.

References

- [1] Public Health Wales Communicable Disease Surveillance Centre. Tuberculosis in Wales annual report. 2018.
- [2] Public Health England. Reports of cases of tuberculosis to enhanced tuberculosis surveillance systems. 2018. United Kingdom, 2000 to 2017.
- [3] Allix-Béguec C, Harmsen D, Weniger T, Supply P, Niemann S. Evaluation and strategy for use of MIRU-VNTRplus, a multifunctional database for online analysis of genotyping data and phylogenetic identification of *Mycobacterium tuberculosis* complex isolates. *J Clin Microbiol* 2008;46:2692–9.
- [4] Walker TM, Ip CLC, Harrell RH, Evans JT, Kapatai G, Dedicoat M, et al. Whole-genome sequencing to delineate *Mycobacterium tuberculosis* outbreaks: a retrospective observational study. *Lancet Infect Dis* 2013;13:137–46.
- [5] Walker TM, Kohl TA, Omar SV, Hedge J, Elias CDO, Bradley P, et al. Whole-genome sequencing for prediction of *Mycobacterium tuberculosis* drug susceptibility and resistance: a retrospective cohort study. *Lancet Infect Dis* 2015;15:1193–202.
- [6] Kohl TA, Diel R, Harmsen D, Rothgänger J, Walter KM, Merker M, et al. Whole-genome-based *Mycobacterium tuberculosis* surveillance: a standardized, portable, and expandable approach. *J Clin Microbiol* 2014;52:2479–86.
- [7] Takiff HE, Feo O. Clinical value of whole-genome sequencing of *Mycobacterium tuberculosis*. *Lancet Infect Dis* 2015;15:1077–90.
- [8] Bjorn-Mortensen K, Soborg B, Koch A, Ladefoged K, Merker M, Lillebaek T, et al. Tracing *Mycobacterium tuberculosis* transmission by whole genome sequencing in a high incidence setting: a retrospective population-based study in East Greenland. *Sci Rep* 2016;6:33180.
- [9] Gardy JL, Johnston JC, Sui SJH, Cook VJ, Shah L, Brodtkin E, et al. Whole-genome sequencing and social-network analysis of a tuberculosis outbreak. *N Engl J Med* 2011;364:730–9.
- [10] Aldous WK, Pounder JI, Cloud JL, Woods GL. Comparison of six methods of extracting *Mycobacterium tuberculosis* DNA from processed sputum for testing by quantitative real-time PCR. *J Clin Microbiol* 2005;43:2471–3.
- [11] Tyler AD, Christianson S, Knox NC, Mabon P, Wolfe J, Van Domselaar G, et al. Comparison of sample preparation methods used for the next-generation sequencing of *Mycobacterium tuberculosis*. *PLoS One* 2016;11:e0148676.
- [12] Schubert M, Ermini L, Der Sarkissian C, Jónsson H, Ginolhac A, Schaefer R, et al. Characterization of ancient and modern genomes by SNP detection and phylogenomic and metagenomic analysis using PALEOMIX. *Nat Protoc* 2014;9:1056–82.
- [13] Orlando L, Gilbert MT, Willerslev E. Reconstructing ancient genomes and epigenomes. *Nat Rev Genet* 2015;16:395–408.
- [14] Carøe CGS, Gopalakrishnan S, Vinner L, Mak SS, Sinding MH, Samaniego JA, et al. Single-tube library preparation for degraded DNA. *Methods Ecol Evol* 2017. <https://doi.org/10.1111/2041-210X.12871>.
- [15] Rasmussen M, Li Y, Lindgreen S, Pedersen JS, Albrechtsen A, Moltke I, et al. Ancient human genome sequence of an extinct Palaeo-Eskimo. *Nature* 2010;463:757–62.
- [16] Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, et al. The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* 2014;505:43–9.
- [17] Orlando L, Ginolhac A, Zhang G, Froese D, Albrechtsen A, Stiller M, et al. Recalibrating Equus evolution using the genome sequence of an early Middle Pleistocene horse. *Nature* 2013;499:74–8.
- [18] Miller W, Drautz DI, Ratan A, Pusey B, Qi J, Lesk AM, et al. Sequencing the nuclear genome of the extinct woolly mammoth. *Nature* 2008;456:387–90.
- [19] Mascher M, Schuenemann VJ, Davidovich U, Marom N, Himmelbach A, Hübner S, et al. Genome-wide analysis of 6,000-year-old cultivated grain illuminates the domestication history of barley. *Nat Genet* 2016;48:1089–93.
- [20] Schuenemann VJ, Singh P, Mendum TA, Krause-Kyora B, Jäger G, Bos KI, et al. Genome-wide comparison of medieval and modern *Mycobacterium leprae*. *Science* 2013;341:179–83.
- [21] Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 2014;30:2114–20.
- [22] Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 2012;19:455–77.
- [23] Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* 2003;13:2498–504.
- [24] Junemann S, Sedlazeck FJ, Prior K, Albersmeier A, John U, Kalinowski J, et al. Updating benchtop sequencing performance comparison. *Nat Biotechnol* 2013;31:294–6.
- [25] Kaas RS, Leekitcharoenphon P, Aarestrup FM, Lund O. Solving the problem of comparing whole bacterial genomes across different sequencing platforms. *PLoS One* 2014;9:e104984.
- [26] Xia E, Teo YY, Ong RTH. SpoTyping: fast and accurate *in-silico* *Mycobacterium* spoligotyping from sequence reads. *Genome Med* 2016;8:19.
- [27] Shabbeer A, Cowan LS, Ozcaglar C, Rastogi N, Vandenberg SL, Yener B, et al. TB-Lineage: an online tool for classification and analysis of strains of *Mycobacterium tuberculosis* complex. *Infect Genet Evol* 2012;12:789–97.
- [28] Aminian M, Shabbeer A, Bennett K. A conformational Bayesian network for classification of *Mycobacterium tuberculosis* complex lineages. *BMC Bioinf* 2010;11: S4.
- [29] Roetzer A, Diel R, Kohl TA, Ruckert C, Nubel U, Blom J, et al. Whole genome sequencing versus traditional genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a longitudinal molecular epidemiological study. *PLoS Med* 2013;10:e1001387.
- [30] Lee RS, Radomski N, Proulx JF, Levade I, Shapiro BJ, McIntosh F, et al. Population genomics of *Mycobacterium tuberculosis* in the Inuit. *Proc Natl Acad Sci* 2015;112:13609–14.