



Swansea University
Prifysgol Abertawe



Cronfa - Swansea University Open Access Repository

This is an author produced version of a paper published in:

Tuberculosis

Cronfa URL for this paper:

<http://cronfa.swan.ac.uk/Record/cronfa51995>

Paper:

Jones, R., Velasco, M., Harris, L., Morgan, S., Temple, M., Ruddy, M., Williams, R., Perry, M., Hitchings, M., et. al. (2019). Resolving a clinical tuberculosis outbreak using palaeogenomic genome reconstruction methodologies.

Tuberculosis, 119, 101865

<http://dx.doi.org/10.1016/j.tube.2019.101865>

Released under the terms of a Creative Commons Attribution Non-Commercial No Derivatives License (CC-BY-NC-ND).

This item is brought to you by Swansea University. Any person downloading material is agreeing to abide by the terms of the repository licence. Copies of full text items may be used or reproduced in any format or medium, without prior permission for personal research or study, educational or non-commercial purposes only. The copyright for any work remains with the original author unless otherwise specified. The full-text must not be sold in any format or medium without the formal permission of the copyright holder.

Permission for multiple reproductions should be obtained from the original author.

Authors are personally responsible for adhering to copyright and publisher restrictions when uploading content to the repository.

<http://www.swansea.ac.uk/library/researchsupport/ris-support/>

1 **Resolving a clinical tuberculosis outbreak using palaeogenomic genome reconstruction**
2 **methodologies**

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

7 ^aSwansea University Medical School, Institute of Life Science, Swansea University,
8 Swansea, Wales, UK.

9 ^bNatural History Museum of Denmark, University of Copenhagen, Copenhagen K,
10 Denmark

11 ^cHealth Protection Division (Mid and West Wales), Public Health Wales, Swansea,
12 Wales, UK.

13 ^dWales Centre for Mycobacteriology, Llandough Hospital, Cardiff, Wales, UK

14 ^ePublic Health Wales Microbiology Swansea, Wales, UK

15 ^fNorwegian University of Science and Technology, University Museum, Trondheim,
16 Norway.

17

18 #Corresponding Author: Angharad Davies:

19 Tel: 00 44 792 295034

20 e-mail: angharad.p.davies@swansea.ac.uk

21

22 **Abstract**

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

33

34 **Keywords:** *Mycobacterium tuberculosis*, whole genome sequencing, outbreak investigation,
35 ancient DNA library construction, palaeogenomics.

36

37 **Introduction**

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

45 During 2003-2005 an outbreak of *M. tuberculosis* occurred in a small town in south Wales
46 (town G). Seven cases (GO1-GO7) were associated with a public house in the town. Case
47 GO3, a barman in the public house, and case GO4, both had direct contact with all the
48 other *M. tuberculosis* cases within the outbreak, at least six of whom had visited the public
49 house. Case GO4 was considered to have been highly infectious, having been symptomatic
50 for about 16 months before diagnosis. Both GO3 and GO4 were thought be the public
51 health team to represent potential super-spreaders within the outbreak. During this period,
52 the standard typing method in use by Public Health Wales was MIRU-VNTR. All the cases
53 had identical MIRU-VNTR patterns. They were all fully susceptible to standard therapy.

54 In 2008, another outbreak was identified in an area of a nearby town, approximately 6km
55 away (town T). There were close links between cases in the two outbreaks, with two of the
56 cases in town T (TH1 and TH2) being direct contacts of GO3. TH2 was a regular at the public
57 house at the centre of the outbreak in town G. However, MIRU-VNTR typing of the two
58 groups of cases differed, with a polymorphism at a single MIRU-VNTR locus (MIRU16),
59 suggesting the presence of two independent outbreaks within the area. Nonetheless the
60 public health team felt it likely that all the cases formed part of one larger outbreak caused
61 by the same strain of *M. tuberculosis* with divergence seen at locus MIRU16 due to a
62 change in an endemic circulating strain. Epidemiologically linked isolates differing at fewer
63 than two MIRU-VNTR loci have been suggested previously to be likely to be part of the
64 same clonal complex within an outbreak [3] and MIRU-VNTR typing has been found to be
65 unable to account for within-outbreak heterogeneity. It also provides limited information
66 on the direction of transmission, identification of super-spreaders or outbreak origin [4, 5,
67 6, 7, 8, 9].

68 The development of affordable and accessible whole genome sequencing (WGS) protocols
69 based around next generation sequencing (NGS) platforms such as the Illumina series, has

70 provided an alternative method for the investigation of *M. tuberculosis* outbreaks. The
71 quality of the DNA sample is critical to the success of WGS. The *M. tuberculosis* samples for
72 this study were provided by the Wales Centre for Mycobacteriology (WCM) in boilate form.

73 Although boilate extraction does release DNA, it is crude, inconsistent and yields DNA of
74 lower integrity, in low quantity and of poorer quality in comparison with other extraction
75 methods [10]. As might be expected, poor sample quality has a negative effect on the
76 standard Nextera XT library preparation [11]. Given the need sometimes to generate WGS
77 data from such samples, lessons might be learnt from the field of palaeogenomics, which
78 attempts routinely to generate genome-scale data from nucleic acids that are highly
79 fragmented, contaminated with non-target DNA, and often contain residual chemical
80 impurities [12, 13, 14]. Recent developments in palaeogenomic methodologies have
81 allowed WGS data to be obtained from a wide range of samples, spanning ancient humans
82 and hominids [15, 16], mammals [17, 18], plants [19] and even pathogens [20] – many of
83 which contain DNA with fragment lengths of <80bp. We therefore hypothesised that the
84 application of palaeogenomic sequencing protocols might overcome the challenge of
85 retrieving genomic data from the low purity and low-quality DNA of crude *M. tuberculosis*
86 boilate samples.

87

88 **Materials and Methods**

89 *Sample collection*

90 Boilate samples from the *M. tuberculosis* isolates described above were obtained from the
91 WCM, Public Health Wales, Llandough Hospital, Cardiff. The isolates had been cultured
92 using the BACTEC™ MGIT™ 960 System (Becton Dickinson Diagnostic Systems, Sparks, MD)
93 in containment level 3 facilities and then heat-killed by boiling for 35 minutes at 110°C.
94 MIRU-VNTR typing had been performed at the PHW Molecular Unit, University Hospital of
95 Wales, Cardiff, with typing based on 15 loci, namely: ETRA, ETRB, ETRC, ETRD, ETRE, MIRU2,
96 MIRU10, MIRU16, MIRU20, MIRU23, MIRU24, MIRU26, MIRU27, MIRU39 and MIRU40.
97 Epidemiological information for each isolate was obtained through face-to-face interviews
98 with a senior public health nurse from the original PHW outbreak investigation team, and
99 from the outbreak documentation. All the cases in these outbreaks and under
100 consideration here were fully susceptible to standard anti-tuberculosis chemotherapy.

101 *Sequencing attempt using conventional Illumina protocols*

102 Sequencing was first attempted directly from the boilates, following an ethanol
103 precipitation. Indexed genomic DNA libraries were prepared for sequencing using the
104 Illumina Nextera XT (V3) sample preparation protocol following the manufacturer's
105 guidelines (2017 Illumina Inc., San Diego, CA, USA), size-selecting for fragments with an
106 average size of 500bp. Bead-normalised sequencing libraries were pooled and sequenced
107 on a MiSeq platform (2017 Illumina Inc., San Diego, CA, USA) using the V3 reagent kits and
108 600 cycles. The resulting paired-end reads were quality filtered with the Trimmomatic tool
109 software [21] using a sliding window approach of 5 bases and a quality score of Q20 prior
110 to contig assembly using the SPAdes genome assembler (Version 3.9.0) with K-mer sizes 33,
111 55, 77, 99 and 127 used [22]. In each case, this method failed to generate any sequence
112 data at all, and unfortunately in this instance, replacement samples were not available. A
113 method optimized for sequencing degraded DNA sources was therefore explored.

114 *DNA extraction*

115 The remaining *M. tuberculosis* boilates were transferred to tubes containing a 500µL
116 solution of digestion buffer (10mM Tris-HCl pH8, 10mM NaCl, 5nM CaCl, 2.5mM EDTA, 1%
117 SDS, 1% Proteinase K, and DTT) and 500µL of Phenol: Chloroform: Isoamyl alcohol solution
118 (Sigma-Aldrich, St. Louis, MO, United States). Next, 0.6g of Zirconia/Silica beads (Cat. No.
119 11079105z, Biospec Products Inc., Bartlesville, OK, USA) were added to the tubes and each
120 sample was homogenized using a TissueLyser II (Qiagen, Valencia, California), for 4 rounds

121 of 20 second bursts with cooling on ice for 30 seconds between rounds. After
122 homogenization, samples were centrifuged for 10 minutes at 16,000 X g in a bench
123 centrifuge to separate the phases. The aqueous upper phase (around 500µL) was gently
124 transferred to a new low-bind 2mL Eppendorf tube and two volumes (1mL) of ice cold
125 absolute ethanol (kept at -20C) were added to each sample. Samples were vortexed briefly
126 and centrifuged again for 10 min at 16,000g. The supernatant was discarded and the tubes
127 washed carefully with 700µL of 70% ethanol without disturbing the pellet. The ethanol
128 wash was discarded and the pellet was left to dry for two minutes. Finally, the pellet was
129 re-suspended and DNA eluted in buffer EB (Qiagen, Valencia, California), Extracted DNA
130 was quantified on a Qubit fluorometer using a dsDNA high sensitivity assay (Life
131 Technologies, Carlsbad, California) and an Agilent 2100 Bioanalyzer (Santa Clara,
132 California). Following extraction, the samples were fragmented for 20 cycles in 30 second
133 cycles within a Diagenode bioruptor 300.

134 *Ancient DNA library preparation protocol*

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

152 *Transmission chain*

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

154 *Gene-by-gene MLST analysis*

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

158 *Whole Genome Sequence SNP mapping*

159 Single nucleotide polymorphisms (SNPs) were identified using the standardised online CSI
160 Phylogeny programme (Version 1.4; Call SNPs & Infer Phylogeny) of the Centre for
161 Genomic Epidemiology (CGE) online tool [<https://cge.cbs.dtu.dk/services/CSIPhylogeny/>]
162 [25]. A minimum spanning tree was constructed based on SNPs from 1123 sites from the
163 WGS data using an adapted application of the Ridom SeqSphere software [24].

164 *In silico spoligotype*

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

174

175 **Results**

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

183 *cgMLST*

184 Gene-by-gene cgMLST analysis appeared to indicate the presence of two outbreaks within
185 the dataset, labelled outbreak 1 and outbreak 2, separated by 124 allelic differences.
186 Outbreak 1 included the outbreak G public house cases GO2, GO3, GO4, GO5, GO6 and
187 GO7 which all had fewer than 12 allelic differences between them. Within outbreak 1, two
188 isolates could not be distinguished from one another, with no allelic differences detected
189 (GO2 and GO6). A star-like topology is seen for outbreak 1, with GO2 and GO6 in the
190 central position. Outbreak 2 is represented by three isolates, two from outbreak TH (TH1
191 and TH2) and from outbreak G, GO1 (Figure 1a). TH1 and GO1 could not be distinguished.
192 Isolate TH2 diverged from GO1 and TH1 by only 4 allelic differences, well within the 12-
193 allele limit and thus representing a direct transmission event [6]. Isolates GO8 and GO9,
194 also from town G and so included in the figure for comparison, substantially exceeded the
195 12-allelic difference threshold for direct transmission with any other isolate within this
196 dataset, consistent with the results of MIRU-VNTR typing, which previously found that GO8
197 and GO9 were not outbreak-related strains.

198 *SNP mapping*

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

206 *In silico spoligotyping*

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

216 **Discussion**

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

224 WGS provided extensive information on the outbreak. The application of gene-by-gene
225 cgMLST provided similar, but not identical, results to those achieved by the SNP mapping
226 procedure. Both agreed in assigning two separate outbreaks and both supported the
227 epidemiological suspicion that a super-spreader was present, although the cgMLST and SNP
228 mapping analyses conflicted in their assignment of the likely super-spreader. cgMLST
229 indicated that the super-spreader was either GO2 or GO6, while traditional SNP mapping
230 supported the assumption of the public health team that GO3 was the super-spreader
231 responsible for multiple secondary cases within this outbreak. A possible reason for the
232 discrepancy between the two methods could be the presence of gene families of a
233 repetitive nature being included in the analyses, such as those for PE_PPE which show
234 disproportionately high amount of divergence [29]. cgMLST removed these regions and
235 this could at least partly explain the discrepancies. It is important to note that most SNP-
236 calling procedures in *M. tuberculosis* epidemiology filter out repetitive regions, and that
237 therefore the procedure of mapping and SNP-calling without filtering, described here, is
238 not directly comparable with that described in other studies.

239 *In silico* application of SpoTyping software was achieved without further laboratory work
240 and at no extra cost to initial sequencing data [26]. Each isolate in this outbreak had a
241 spoligotype pattern corresponding to the Haarlem H1 clade. It has been demonstrated
242 previously that the value of inferring a recent common ancestor of isolates within
243 potentially related outbreaks, with identification of a causative circulating strain being a
244 common feature [8, 30]. Previous studies of *M. tuberculosis* in the Inuit and Greenland
245 populations of North America, which have stable populations, have also documented the
246 cause of multiple outbreaks within the region as the ongoing spread of an evolving founder
247 strain that has continually spread across the area over decades [8, 30].

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

257 The work described here demonstrates that when necessary, clinically useful data can be
258 obtained from sub-optimal quality samples by applying an ancient DNA library construction
259 protocol to overcome the need for DNA of high purity and quality. The success of the 'BEST'
260 aDNA library construction protocol here highlights a clinical application for a method
261 previously associated only with palaeogenomic studies, and shows how the transfer of such
262 techniques in defined circumstances could provide clinical benefit.

263 **Acknowledgments**

264 This work was funded by St. David's Medical Foundation & Coleg Cenedlaethol Cymraeg
265 funding to APD & RJ, a EMBO Short-term Fellowship to RJ, and an ERC Consolidator Grant
266 (681396- Extinction Genomics) to MTPG. The authors would like to thank Christian Carøe
267 and other members of the Gilbert laboratory at the Natural History Museum, Denmark for
268 technical guidance.

269

270

271

272

273 **References**

- 274 1. Public Health Wales Communicable Disease Surveillance Centre. 2018. Tuberculosis
275 in Wales Annual Report.
- 276 2. Public Health England 2018. Reports of cases of tuberculosis to enhanced
277 tuberculosis surveillance systems: United Kingdom, 2000 to 2017.
- 278 3. Allix-Béguec C, Harmsen D, Weniger T, Supply P, Niemann S. 2008. Evaluation and
279 strategy for use of MIRU-VNTRplus, a multifunctional database for online analysis
280 of genotyping data and phylogenetic identification of *Mycobacterium tuberculosis*
281 complex isolates. J Clin Microbiol. 46:2692-2699.
- 282 4. Walker TM, Ip CLC, Harrell RH, Evans JT, Kapatai G, Dedicoat M, Eyre DW, Wilson
283 DJ, Hawkey PM, Crook DW, Parkhill J, Harris D, Walker AS, Bowden R, Monk P,
284 Smith EG, Peto TE. 2013. Whole-genome sequencing to delineate *Mycobacterium*
285 *tuberculosis* outbreaks: a retrospective observational study. Lancet Inf Dis. 13:137-
286 146.
- 287 5. Walker TM, Kohl TA, Omar SV, Hedge J, Elias CDO, Bradley P, Iqbal Z, Feuerriegel S,
288 Niehaus KE, Wilson DJ. 2015. Whole-genome sequencing for prediction of
289 *Mycobacterium tuberculosis* drug susceptibility and resistance: a retrospective
290 cohort study. Lancet Inf Dis.15:1193-1202.
- 291 6. Kohl TA, Diel R, Harmsen D, Rothgänger J, Walter KM, Merker M, Weniger T,
292 Niemann S. 2014. Whole-genome-based *Mycobacterium tuberculosis* surveillance:
293 a standardized, portable, and expandable approach. J Clin Microbiol. 52:2479-2486.
- 294 7. Takiff HE & Feo O. 2015. Clinical value of whole-genome sequencing of
295 *Mycobacterium tuberculosis*. Lancet Inf Dis. 15:1077-1090.
- 296 8. Bjorn-Mortensen K, Soborg B, Koch A, Ladefoged K, Merker M, Lillebaek T,
297 Andersen A, Niemann S, Kohl T. 2016. Tracing *Mycobacterium tuberculosis*
298 transmission by whole genome sequencing in a high incidence setting: a
299 retrospective population-based study in East Greenland. Sci Rep. 6:33180.
- 300 9. Gardy JL, Johnston JC, Sui SJH, Cook VJ, Shah L, Brodtkin E, Rempel S, Moore R, Zhao
301 Y, Holt R, Varhol R, Birol I, Lem M, Sharma M, Elwood K, Jones S, Brinkman F,
302 Brunham R, Tang P. 2011. Whole-Genome Sequencing and Social-Network Analysis
303 of a Tuberculosis Outbreak. NEJM. 364:730-739.
- 304 10. Aldous WK, Pounder JI, Cloud JL, Woods GL. 2005. Comparison of six methods of
305 extracting *Mycobacterium tuberculosis* DNA from processed sputum for testing by
306 quantitative real-time PCR. J Clin Microbiol. 43:2471-2473.

- 307 11. Tyler AD, Christianson S, Knox NC, Mabon P, Wolfe J, Van Domselaar G, Graham
308 MR, Sharma MK. 2016. Comparison of sample preparation methods used for the
309 next-generation sequencing of *Mycobacterium tuberculosis*. PLoS ONE.
310 11:e0148676.
- 311 12. Schubert M, Ermini L, Der Sarkissian C, Jónsson H, Ginolhac A, Schaefer R, Martin
312 MD, Fernandez R, Kircher M, McCue M, Willersely E, Orlando L. 2014.
313 Characterization of ancient and modern genomes by SNP detection and
314 phylogenomic and metagenomic analysis using PALEOMIX. Nat Protoc. 9:1056-82.
- 315 13. Orlando L, Gilbert MT, Willerslev E. 2015. Reconstructing ancient genomes and
316 epigenomes. Nat Rev Genet. 16:395-408.
- 317 14. Carøe CGS, Gopalakrishnan S, Vinner L, Mak SS, Sinding MH, Samaniego JA, Wales
318 N, Sicheritz-Ponten T, Gilbert MTP. 2017. Single-tube library preparation for
319 degraded DNA. Methods Ecol Evol. doi10.1111/2041-210X.12871
- 320 15. Rasmussen M, Li Y, Lindgreen S, Pedersen JS, Albrechtsen A, Moltke I, Metspalu
321 M, Metspalu E, Kivisild T, Gupta R, Bertalan M, Nielsen K, Gilbert MT, Wang
322 Y, Raghavan M, Campos PF, Kamp HM, Wilson AS, Gledhill A, Tridico S, Bunce
323 M, Lorenzen ED, Binladen J, Guo X, Zhao J, Zhang X, Zhang H, Li Z, Chen M, Orlando
324 L, Kristiansen K, Bak M, Tommerup N, Bendixen C, Pierre TL, Grønnow
325 B, Meldgaard M, Andreasen C, Fedorova SA, Osipova LP, Higham TF, Ramsey
326 CB, Hansen TV, Nielsen FC, Crawford MH, Brunak S, Sicheritz-Pontén T, Villems
327 R, Nielsen R, Krogh A, Wang J, Willerslev E. 2010. Ancient human genome sequence
328 of an extinct Palaeo-Eskimo. Nature. 463:757-62.
- 329 16. Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud
330 G, Sudmant PH, de Filippo C, Li H, Mallick S, Dannemann M, Fu Q, Kircher
331 M, Kuhlwilm M, Lachmann M, Meyer M, Ongyerth M, Siebauer M, Theunert
332 C, Tandon A, Moorjani P, Pickrell J, Mullikin JC, Vohr SH, Green RE, Hellmann
333 I, Johnson PL, Blanche H, Cann H, Kitzman JO, Shendure J, Eichler EE, Lein
334 ES, Bakken TE, Golovanova LV, Doronichev VB, Shunkov MV, Derevianko AP, Viola
335 B, Slatkin M, Reich D, Kelso J, Pääbo S. 2014. The complete genome sequence of a
336 Neanderthal from the Altai Mountains. Nature. 505,43-9.
- 337 17. Orlando L, Ginolhac A, Zhang G, Froese D, Albrechtsen A, Stiller M, Schubert M,
338 Cappellini E, Petersen B, Moltke I, Johnson PL, Fumagalli M, Vilstrup JT, Raghavan
339 M, Korneliusen T, Malaspinas AS, Vogt J, Szklarczyk D, Kelstrup CD, Vinther J,
340 Dolocan A, Stenderup J, Velazquez AM, Cahill J, Rasmussen M, Wang X, Min J,

- 341 Zazula GD, Seguin-Orlando A, Mortensen C, Magnussen K, Thompson JF, Weinstock
342 J, Gregersen K, Røed KH, Eisenmann V, Rubin CJ, Miller DC, Antczak DF, Bertelsen
343 MF, Brunak S, Al-Rasheid KA, Ryder O, Andersson L, Mundy J, Krogh A, Gilbert MT,
344 Kjær K, Sicheritz-Ponten T, Jensen LJ, Olsen JV, Hofreiter M, Nielsen R, Shapiro B,
345 Wang J, Willerslev E. 2013. Recalibrating Equus evolution using the genome
346 sequence of an early Middle Pleistocene horse. *Nature*. 499:74-8.
- 347 18. Miller W, Drautz DI, Ratan A, Pusey B, Qi J, Lesk AM, Tomsho LP, Packard MD, Zhao
348 F, Sher A, Tikhonov A, Raney B, Patterson N, Lindblad-Toh K, Lander ES, Knight
349 JR, Irzyk GP, Fredrikson KM, Harkins TT, Sheridan S, Pringle T, Schuster SC. 2008.
350 Sequencing the nuclear genome of the extinct woolly mammoth. *Nature*. 456:387-
351 90.
- 352 19. Mascher M., Schuenemann VJ, Davidovich U, Marom N, Himmelbach A, Hübner
353 S, Korol A, David M, Reiter E, Riehl S, Schreiber M, Vohr SH, Green RE, Dawson
354 IK, Russell J, Kilian B, Muehlbauer GJ, Waugh R, Fahima T, Krause J, Weiss E, Stein
355 N. 2016. Genomic analysis of 6,000-year-old cultivated grain illuminates the
356 domestication history of barley. *Nat Genet*. 48:1089-93.
- 357 20. Schuenemann VJ, Singh P, Mendum TA, Krause-Kyora B, Jäger G, Bos KI, Herbig
358 A, Economou C, Benjak A, Busso P, Nebel A, Boldsen JL, Kjellström A, Wu H, Stewart
359 GR, Taylor GM, Bauer P, Lee OY, Wu HH, Minnikin DE, Besra GS, Tucker K, Roffey
360 S, Sow SO, Cole ST, Nieselt K, Krause J. 2013. Genome-wide comparison of medieval
361 and modern *Mycobacterium leprae*. *Science*. 341:179-83.
- 362 21. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina
363 sequence data. *Bioinformatics*. 30:2114-2120
- 364 22. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin
365 VM, Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler
366 G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm
367 and its applications to single-cell sequencing. *J Comput Biol*. 19:455-477.
- 368 23. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin
369 N, Schwikowski B, Ideker T. 2003. Cytoscape: a software environment for
370 integrated models of biomolecular interaction networks. *Genome Res*. 13:2498-
371 2504.
- 372 24. Junemann S, Sedlazeck FJ, Prior K, Albersmeier A, John U, Kalinowski J, Mellmann
373 A, Goesmann A, von Haeseler A, Stoye J, Harmsen D. 2013. Updating benchtop
374 sequencing performance comparison. *Nat Biotechnol*. 31:294-6.

- 375 25. Kaas RS, Leekitcharoenphon P, Aarestrup FM, Lund, O. 2014. Solving the problem
376 of comparing whole bacterial genomes across different sequencing platforms. PLoS
377 ONE. 9:e104984.
- 378 26. Xia E, Teo YY, Ong RTH. 2016. SpoTyping: fast and accurate *in-silico Mycobacterium*
379 spoligotyping from sequence reads. Genome Med. 8:19.
- 380 27. Shabbeer A, Cowan LS, Ozcaglar C, Rastogi N, Vandenberg SL, Yener B, Bennett KP.
381 2012. TB-Lineage: an online tool for classification and analysis of strains of
382 *Mycobacterium tuberculosis* complex. Inf Genet Evol. 12:789-797.
- 383 28. Aminian M, Shabbeer A, Bennett K. 2010. A conformal Bayesian network for
384 classification of *Mycobacterium tuberculosis* complex lineages. BMC Bioinformatics.
385 11:S4
- 386 29. Roetzer A, Diel R, Kohl TA, Ruckert C, Nubel U, Blom J, Wirth T, Jaenicke S, Schuback
387 S, Rusch-Gerdes S, Supply P, Kalinowski J, Niemann S. 2013. Whole Genome
388 Sequencing versus Traditional Genotyping for Investigation of a *Mycobacterium*
389 *tuberculosis* Outbreak: A Longitudinal Molecular Epidemiological Study. PLOS Med.
390 10: e1001387
- 391 30. Lee RS, Radomski N, Proulx JF, Levade I, Shapiro BJ, McIntosh F, Soualhiné
392 H, Menzies D, Behr MA. 2015. Population genomics of *Mycobacterium tuberculosis*
393 in the Inuit. Proc Nat Acad Sci. 112:13609-13614.
- 394

395 **Table 1:** Ancient DNA protocol sequencing of nine outbreak isolates. The percentage of the
 396 core genome MLST genes present is shown as well as the percentage of the reference *M.*
 397 *tuberculosis* H37Rv genome covered according to the CSI phylogeny algorithm provided by
 398 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

399

400

401 **Table 2:** *in silico* spoligotyping results for each of the outbreak isolates. Results include the
 402 predicted spoligotype (produced by SpolDB4), international spoligotype (SITVIT database),
 403 lineage assignment and clade assignment (both outputted from the TB-insight online
 404 server).

Isolate	Predicted Spoligotype	International type	Lineage	Clade
GO3	777777774020771	47	Euro-American	H1
TH1	777777770000000	46	Euro-American	H1
TH2	777777770000000	46	Euro-American	H1
GO1	777777770000000	46	Euro-American	H1
GO2	777777770020771	742	Euro-American	H1
GO4	777777770020771	742	Euro-American	H1
GO5	777777774020771	47	Euro-American	H1
GO6	777777774020771	47	Euro-American	H1
GO7	777777774020771	47	Euro-American	H1

405

406 **Figure Legend**

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

415

416

Figure 1

