# GigaScience

## Screening methods for detection of ancient Mycobacterium tuberculosis complex fingerprints in NGS data derived from skeletal samples

### --Manuscript Draft--

| Manuscript Number: | GIGA-D-19-00014R1 |
|---|---|
| Full Title: | Screening methods for detection of ancient Mycobacterium tuberculosis complex fingerprints in NGS data derived from skeletal samples |
| Article Type: | Research |

| Abstract: | Background<br><br>Recent advances in ancient DNA (aDNA) studies, especially in increasing isolated DNA yields and quality, opened the possibility of analysis of ancient host microbiome. However, this analysis could lead to numerous pitfalls, including spurious identification of pathogens based on fragmentary data or environmental contamination, leading to incorrect epidemiological conclusions. Within the Mycobacterium genus, MTBC (Mycobacterium tuberculosis complex) members responsible for tuberculosis share up to ~99% genomic sequence identity, while other more distantly related MOTT (Mycobacteria other than tuberculosis) can be causative agents for pulmonary diseases or soil dwellers. Therefore, reliable determination of species complex is highly relevant for interpretation of sequencing results.<br><br>Results<br><br>Here we present a novel bioinformatical approach, used for screening of ancient tuberculosis in sequencing data, derived from 28 individuals (dated 4400 - 4000 BC and 3100 - 2900 BC) from Central Poland. We demonstrate that cost effective next generation screening sequencing data (c.a 20M reads per sample) could yield enough information to provide statistically supported identification of probable ancient disease cases.<br><br>Conclusions<br><br>Application of appropriate bioinformatic tools, including an unbiased selection of genomic alignment targets for species specificity, makes it possible to extract valid data from full-sample sequencing results (without subjective targeted enrichment procedures). This approach broadens the potential scope of paleoepidemiology both to older, suboptimally preserved samples and to pathogens with difficult intrageneric taxonomy. |
|---|---|

| Corresponding Author: | Dominik Strapagiel<br>University of Lodz<br>Lodz, POLAND |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | University of Lodz |
| Corresponding Author's Secondary Institution: | |
| First Author: | Paulina Borówka |
| First Author Secondary Information: | |
| Order of Authors: | Paulina Borówka |

| | Lukasz Pułaski |
| | Błażej Marciniak |
| | Beata Borowska-Strugińska |
| | Jarosław Dziadek |
| | Elżbieta Żądzińska |
| | Wiesław Lorkiewicz |
| | Dominik Strapagiel |
| | Łukasz Pułaski |
| **Order of Authors Secondary Information:** | |
| **Response to Reviewers:** | Lodz, 7th of April, 2019<br>Dear GigaScience Editor,<br>We appreciate the insightful and detailed comments of the reviewers. We included all the minor language-related corrections in the resubmitted manuscript. Following is a point-by-point reply to the major points raised by the reviewers.<br>Reviewer 1<br>1. Our statistical approach, as we state directly in the manuscript, is aimed at finding positive outliers in a pool of samples with unknown presence of mycobacterial sequences. When applied to a dataset like Kay et al., which consists exclusively of individuals with previously confirmed ancient mycobacterial infection, its outcome is therefore necessarily limited to identifying the outliers with highest microbial load (in the case of tuberculosis, potentially those individuals who died during the active phase of the disease) - these outliers being by definition always a minority of analysed samples. This explanation is provided in the text of the manuscript.<br>2. We have performed the MapDamage analysis suggested by the reviewer and it indeed yielded a positive result - we thank the reviewer for this suggestion as this strengthened our conclusions significantly. We have now included a new Supplementary Fig. 4, and we have reworded both the legend to Supplementary Fig. 3 and the sentences in the manuscript that refer to it.<br>3. Libraries were build using Meyer et al. (2010) protocol with modifications proposed by Gamba et al. (2014). We have performed mapDamage analysis in the way which fit to double stranded libraries. Information about different Meyer protocol and single stranded libraries was incorrectly added to previous version of the manuscript.<br>4. We have expanded Supplementary Tab. 2 to include the absolute numbers of reads - both total and aligning to each alignment target.<br>Reviewer 2<br>1. We have now expanded the analysis of state of the art in biochemical detection of ancient mycobacteria by citing recent articles mentioning improvements in cell wall component analysis.<br>2. Indeed, Pott's disease is usually regarded as pathognomonic signature of TB. When we said that many pathological conditions of the spine can mimic Pott's disease thought that they can be diagnosed mistakenly as tuberculosis, especially in practice with poorly preserved skeletons. The present text has been appropriately modified (both in the Introduction and Discussion) to clarify this statement. Moreover, according to the Reviewer's suggestion, the list of pathological conditions has been replaced with the Table S1, which include a short description of basic differences between these lesions and bone tuberculosis.<br><br>Dear Editor, regarding to your comment:<br>"As your revised manuscript focuses more on a method, it may be more suitable as a "Technical Note" rather than a "research article" "<br>We would like to proceed this manuscript as research paper. In our work we present original dataset which allow us to present a novel bioinformatical approach, used for screening of ancient tuberculosis in sequencing data, derived from 28 individuals (dated 4400 - 4000 BC and 3100 - 2900 BC) from Central Poland. That dataset was not previously published elsewhere.<br><br>Sincerely,<br>Dominik Strapagiel |

| Additional Information: | |
|---|---|
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript. | Yes |

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

**Screening methods for detection of ancient *Mycobacterium tuberculosis* complex fingerprints in NGS data derived from skeletal samples**

Paulina Borówka[1], Łukasz Pułaski[2,3], Błażej Marciniak[4,5], Beata Borowska-Strugińska[1], Jarosław Dziadek[3], Elżbieta Żądzińska[1], Wiesław Lorkiewicz[1,#], Dominik Strapagiel[4,5,#]

1 Department of Anthropology, Faculty of Biology and Environmental Protection, University of Lodz, 90-237 Łódź, Poland;

2 Department of Molecular Biophysics, Faculty of Biology and Environmental Protection,University of Lodz, 90-237 Łódź, Poland;

3 Institute of Medical Biology, Polish Academy of Sciences, 93-232 Łódź, Poland;

4 Biobank Lab, Faculty of Biology and Environmental Protection, Department of Molecular Biophysics, University of Lodz, 90-231 Łódź, Poland;

5 BBMRI.pl Consortium, 54-066 Wrocław, Poland.

Corresponding author:

Dominik Strapagiel, Biobank Lab, Department of Molecular Biophysics, Faculty of Biology and Environmental Protection, University of Lodz, Pilarskiego 14, 90-231 Łódź, Poland, +48426655702, dominik.strapagiel@biol.uni.lodz.pl

Wiesław Lorkiewicz, Department of Anthropology, Faculty of Biology and Environmental Protection, University of Lodz, 90-237 Łódź, Poland, +48426354456, wieslaw.lorkiewicz@biol.uni.lodz.pl

**Abstract**

1    **Background:** Recent advances in ancient DNA (aDNA) studies, especially in increasing isolated DNA yields and quality, opened the

2    possibility of analysis of ancient host microbiome. However, this analysis could lead to numerous pitfalls, including spurious identification of

3    pathogens based on fragmentary data or environmental contamination, leading to incorrect epidemiological conclusions. Within the

4    *Mycobacterium* genus, MTBC (*Mycobacterium tuberculosis complex*) members responsible for tuberculosis share up to ~99% genomic

5    sequence identity, while other more distantly related MOTT (*Mycobacteria* other than *tuberculosis*) can be causative agents for pulmonary

6    diseases or soil dwellers. Therefore, reliable determination of species complex is highly relevant for interpretation of sequencing results.

7    **Results:** Here we present a novel bioinformatical approach, used for screening of ancient tuberculosis in sequencing data, derived from

8    28 individuals (dated 4400 - 4000 BC and 3100 - 2900 BC) from Central Poland. We demonstrate that cost effective next generation screening

9    sequencing data (c.a 20M reads per sample) could yield enough information  to provide statistically supported identification of probable ancient

10   disease cases.

11   **Conclusions:** Application of appropriate bioinformatic tools, including an unbiased selection of genomic alignment targets for species

12   specificity, makes it possible to extract valid data from full-sample sequencing results (without subjective targeted enrichment procedures). This

13   approach broadens the potential scope of paleoepidemiology both to older, suboptimally preserved samples and to pathogens with difficult

14   intrageneric taxonomy.

15   **Keywords**

16   **ancient DNA, aTB, ancient tuberculosis, NGS**

**Background**

A rapid population growth initiated in the Neolithic period, connected with domestication of animals and increase of human sedentism, played a key role in pathogen transmission within the so-called first epidemiological transition[1-4]. The identification of infectious diseases and selection of unique fingerprints of their causative agents, especially those derived from skeletal elements, are still of the greatest interest for paleopathologists and anthropologists, which is evidenced by the range of available analysis methods. Members of the *Mycobacterium tuberculosis* complex (MTBC) are genetically very closely related and are causative agents for one of the oldest human infectious diseases – tuberculosis (TB). It is a disease that may leave lesions on patients' bones, enabling a diagnosis based on bone morphology [5]. The main problem of paleopathological diagnoses based solely on dry bones is that TB-related bone changes are often nonspecific. The most reliable skeletal indicator of TB are destructive lesions in thoracic and lumbar spine sections, which can lead to destruction and collapse of vertebral bodies, resulting in spinal kyphosis, or gibbus, known as Pott's disease [5-7]. However, there are several pathological conditions which could mimic TB in dry bone leading to erroneous diagnosis, especially when they affect the spine (Supplementary Tab. 1). Although their differential diagnosis from TB is well known in paleopathology it could be problematic to use it in analysis of often poorly preserved archaeological human remains[8, 9]. Diagnoses based on bone lesions in other region of the skeleton are even more tentative; these are primarily based on osteomyelitis of the joints (most commonly the hip and knee, but also ankle and elbow) and periosteal reactive lesions (mainly in the ribs or diaphysis of the long bones, including tubular bones of the hands and feet in children [6, 9]. Bone lesions from TB in nonspinal locations may be indistinguishable from those of other etiologies[5, 6] . Lastly, morphological studies of bones do not permit detection of many individuals affected with TB in past human populations: data from the pre-antibiotic era show that bone changes occur only in about 3–7% of individuals with active TB [9].

Since the 1990s, new possibilities to diagnose TB in archaeological specimens have arisen, offered by the detection and analysis of mycobacterial DNA and other biomolecules specific to MTBC at the molecular level [10-21]. A common complication in molecular studies for ancient MTBC detection is the presence of DNA and other metabolites from the whole microbiome of the individual whose remains are being analysed as well as from environmental bacteria that have colonised the skeleton *post-mortem* [22, 23]. These contaminants might include Mycobacteria other than *M. tuberculosis* (MOTT), some of which are prevalent in the environment, while others are associated with clinical cases of non-tuberculosis diseases [22, 24-26]. It should be emphasized that members of *Mycobacterium tuberculosis* complex themselves are characterized by a particular high sequence similarity [27, 28], which leads to often unsurmountable difficulties in distinguishing them on the molecular level.

Detection of cell wall components such as mycolic, mycocerosic and mycolipenic acids [13, 15, 18, 19] with matrix-assisted laser desorption/ionization tandem time of flight (MALDI-TOF) which present profiles specific for MTBC is considered a reliable method to identify ancient causative agents in human archaeological samples. Initial attempts to use mass spectrometry to detect cell wall lipids were shown to be erroneous in some cases [15, 29, 30]. In more recent studies, the combination of cell wall lipid analysis with genetic markers showed significant improvement in discriminative ability for ancient mycobacteria [31, 32]. Polymerase chain reaction, followed by gel electrophoresis, is still a popular method for detection of MTBC ancient DNA in human samples such as bones and teeth [32-34], mummified soft tissues [35, 36], or calcified pleura [10]. Known cases of tuberculosis disease proven on the basis of ancient DNA derived from human material are old as 9000 BC [37], through Iron Age [38] and up to modern times [39]. However, PCR-based methods have not been without controversy due to the possibility of cross-contamination as well as limitations of selection of proper sequences. While repetitive insertion sequences, e.g. IS6110 and IS1081, are widely used and sometimes considered as a biomarker specific to MTBC bacteria [34], the current consensus recommends strong caution in their

53   use due to their presence in MOTT bacteria. Those commonly used markers have even been found to occur in soil mycobacteria [40-45], and

54   even weak homology can cause false-positive PCR results for unrelated microbes [40, 46].

55   Recently, next generation sequencing (NGS) methods were introduced for detection of causative agents of ancient diseases [47, 48],

56   including MTBC, with or without pre-enrichment of MTBC aDNA [49-52]. The increasing quantity of data generated by NGS and efficiency of

57   non-Sanger-based sequencing platforms requires a new approach in processing tools: suitable bioinformatic pipelines are required for reliable

58   DNA analysis of ancient causative agents. Similar to PCR, where the use of only short conserved regions considered as specific for MTBC may

59   lead to false positive results, improper analysis of NGS data can misinterpret sequences from modern known or unknown environmental

60   *Mycobacteria* which are present in ancient human skeletons [26]. New analytical tools for more unequivocal answers to questions of

61   identification and differentiation of *ante-mortem* causative from *post-mortem* non-causative microbial agents are urgently needed. Application of

62   specifically designed *in silico* (bioinformatical approach) verification methods for improved downstream processing of molecular fingerprint

63   data from ancient samples is necessary for drawing conclusions on clinical prevalence and epidemiology of pathogenic mycobacteria in history.

64   Here we present an improved strategy for specific identification of bacteria from the *M. tuberculosis* complex in ancient non-enriched NGS data.

65   The main purpose of this study was to design an unbiased genomic marker alignment query composed of sequences belonging strictly to MTBC

66   members. Therefore, we present a workflow including appropriate bioinformatic alignment algorithms and statistical tools that allowed the

67   identification of tuberculosis causative agents, using fragment length variation to balance selectivity (species specificity) with sensitivity of

68   detection.

69   **Sample Description**

70   Ancient bone samples come from skeletal remains of 28 individuals representing two Neolithic populations from the Kujawy region in

71   Central Poland: the Middle Neolithic Brześć Kujawski Group of the Lengyel culture (BKG), dated to ca. 4400-4000 BC (26 individuals) and the

72   Late Neolithic Globular Amphora culture (GAC), dated to ca. 3100-2900 BC (2 individuals), previously described in [18, 53] (Supplementary

73   Tab. 2). The skeletons come from two archaeological sites, BK 3 and BK 4, which represent relics of a settlement and cemetery of the BKG

74   culture with some secondary objects within them, like the GAC grave. Both sites overlap each other, thus soil conditions and diagenetic agents

75   were similar for all skeletal remains analyzed. Bone material was taken mainly from vertebral bodies of individuals with well-preserved

76   skeletons. One of two individuals belonging to the GAC revealed bone lesions consistent with Pott's disease. BKG samples provided more

77   ambiguous evidence of skeletal lesions. One individual showed destructive lesions of the thoracic and lumbar vertebrae with central collapse of

78   the vertebral bodies which may indicate tuberculous spondylodiscitis. Three other individuals of this population revealed only relatively mild

79   and nonspecific inflammatory bone changes in the postcranial skeleton which were located on the internal surface of the ribs, tibia and femur

80   shafts, as well as foot bones.

81   **Analyses**

82   **Reference target construction (alignment target)**

83   As our main reference sequence, we used the most commonly applied modern laboratory strain of *M. tuberculosis* (MTB), H37Rv, for

84   which the whole genomic sequence is available. In order to select a subset of this reference sequence as an alignment target providing enhanced

85   specificity for tuberculosis-causing agents (MTBC members), we first derived a set of all protein-coding sequences (CDS) from the H37Rv

86   genome using the RAST tool [54]. These 4,360 sequences were screened using the BLAST tool (Megablast) at the National Library of Medicine

87   sequentially against 12 available genomic sequences of selected MOTT: *M. kansasii, M. avium subsp. paratuberculosis, M. ulcerans, M.*

88   *smegmatis, M. fortuitum, M. haemophilum, M. marinum, M. simiae, M. asiaticum, M. xenopi, M. phlei, M. abscessus.* Any detected similarities

89   (gapless alignments >10 bp) between a H37Rv CDS and any MOTT genomic sequences resulted in the exclusion of this CDS from the result

90    dataset, which was therefore restricted to sequences fully specific for MTBC, having no homologs in any MOTT genome. The resulting set of

91    sequences was subsequently called the BoRówka et al. alignment target and consisted of 1,534 coding sequences with total sequence length of

92    0.814 Mbp. Since no sequences from other MTBC species were used at this stage, and it is known that they exhibit up to 99.9% nucleotide

93    sequence similarity [55], the constructed alignment target cannot be considered specific only for *M. tuberculosis*, but rather for the whole

94    MTBC; this is justified in epidemiological studies on ancient samples by the need to include all clinically equivalent causative agents for the

95    same disease entity: tuberculosis. For comparison purposes, we prepared and used two literature-derived, knowledge-based H37Rv sequence

96    subsets as alternative alignment targets: the c. 0.046 Mbp sequence used for capture enrichment in Bouwman et al. (2012) [52] for sequencing

97    mycobacterial samples from a 19th century skeleton, subsequently called the Bouwman et al. alignment target, and the two genes (*katG* and

98    *mpt40*, total length 0.004 Mbp) listed as MTBC-specific among the capture enrichment probes used by Bos et al. (2014) [50] for sequencing

99    mycobacterial samples from 11th-13th century Peruvian skeletons, subsequently called the Bos et al. alignment target. All the reference

100   sequences were prepared for alignment by indexing with the suffix array - induced sorting algorithm, implemented in the BWA software

101   package (BWA).

102         Since the construction of the Borówka et al. alignment target was based on elimination of sequences similar to other mycobacterial

103   species, we reasoned that the performance of an alignment target is directly linked to the number of similarities between the MTB genome and

104   other potentially interfering mycobacterial species (both ancient and environmental) present in the ancient host-derived sample. In order to

105   quantify this, we subjected the publicly available genome sequences of *Mycobacterium* species to an *in-silico* procedure to generate collections

106   of short sequences broadly analogous to authentic NGS reads. Including reads below a certain length in similarity analysis of ancient microbial

107   DNA leads to non-specific matches (for both evolutionary and statistical reasons); this threshold is usually arbitrarily set to around 30 bp, but a

108   broader analysis might make it easier to construct a reliable algorithm for detection of specific ancient pathogens. Therefore, in our further

109   analysis both of reference and authentic ancient NGS sequences we extracted groups (bins) of non-human sequences over several length

110   thresholds: ≥20bp, ≥25bp, ≥30bp and ≥35bp, to enable a thorough analysis of specificity gain upon increase in minimal sequence length. For

111   reference *Mycbacterium* genomes, k-mers of specified length (corresponding to the lower limit of read length for NGS bins: 20, 25, 30 or 35)

112   were filtered against the human genome assembly hg19, and the resulting "short read" collections were aligned to the full MTB reference

113   genome or its selected subsets (Borówka et al., Bouwman et al. and Bos et al. alignment targets). Table 1 shows the respective number of

114   genomic k-mers from MTB complex and MOTT species which match the MTBC alignment targets as well as the total lengths of assayed

115   genomes for comparison. Since the various subsets of the MTB genome differ in length and thus the probability of random match increases with

116   target length, we standardised the obtained data by presenting it as a percentage of k-mers from a given mycobaterial genome that match the

117   alignment target, divided by the ratio of target length to the full MTB genome length (genomic coverage of the target). These values, which are

118   an inverse measure of alignment target specificity (they increase if more "reads" from a species which is not MTB or MTBC can be mistaken for

119   MTBC), are shown in Table 1. As a reference, the MTB genome itself was also subjected to this procedure - obviously, the match percentage

120   values are almost 100% here. Several conclusions can be drawn from these data: firstly, it is obvious that selecting longer reads (in this case

121   longer k-mers) for comparison increases specificity, with reads 30 bp long or longer optimal for specific identification of the MTB complex,

122   reflecting a common consensus in the field. However, it is important to note that shorter reads still add important information to the analysis, as

123   the rate of specificity increase (decrease in matching read percentage with increase in read length) varies between species (i.e. some species have

124   longer stretches of highly similar sequence to MTB). For example, while *M. smegmatis* has a very high match percentage to the Borowka et al.

125   alignment target at low read length, this is rapidly lost at longer (more genuine) read lengths; the opposite is true e.g. for *M. marinum*. It is a

126   derivation of the evolutionary history of the genus, but in this case also a practical caveat for further interpretation of sequence matches in actual

127   aDNA samples. Moreover, the specificity of various alignment targets varies, with the Borówka et al. target being consistently the most specific

128  (for longer k-mers) for distinguishing MOTT, while it is (by design) not well suited to distinguishing other members of the MTB complex from

129  MTB itself.

130      Since we intended to develop a highly specific screening test (based on low depth sequencing strategy) for verification of MTBC

131  infection in Neolithic samples with *a priori* relatively low degree of aDNA preservation, we decided on a statistical approach. Since any

132  preserved ancient mycobacterial DNA would be only a fraction of total aDNA, and it in turn would only be a fraction of total reads (the balance

133  being the modern environmental metagenome), a balance between sensitivity and specificity in verifying this very low number of reads must be

134  struck. In sedentary, communal populations MTBC infection tends to be epidemic in character, but in most individuals with latent infection the

135  microbial load (and thus the probability of DNA survival in ancient samples) is relatively low and constant. Any similarity analysis based on

136  sequence alignment will also invariably generate false positive alignment hits, thus, it would be impossible to construct a test with sufficient

137  statistical power to distinguish individuals genuinely free of ancient MTBC and those with average/modest latent infection. Therefore, we

138  concentrated on the detection of outlier individuals with high microbial load (which may be later selected for enrichment-based further genetic

139  analysis, such as phylogenetic studies or genome reconstruction), measured by the positive read ratio (the intrinsically very low ratio of reads

140  matching the MTBC alignment target to all eligible reads). Based on the epidemiology of MTBC infection, we assumed a quasi-normal

141  distribution of positive read ratios in a randomly selected sample of ancient individuals, with outliers as candidates for active tuberculosis and for

142  selection for more in-depth studies. Thus, our method was based on standardising read ratio values to normal distribution parameters (arithmetic

143  mean and standard deviation) and, as a further step in the detection algorithm for ancient tuberculosis (aTB), we applied a typical cutoff value of

144  1.5xSD to detect outliers.

145      As a first stage of testing our screening approach on actual NGS data from ancient material, we used a control dataset based on published

146  NGS results of confirmed tuberculosis-infected individuals - 18th/19th-century mummified bodies from a crypt in Vác, Hungary, described by

147  Kay et al. (2015) [48]. The aim of <u>that study</u> was to reconstruct and analyse historical genome sequences of *M. tuberculosis*, which resulted in

148  sequencing results with high coverage. Since all these samples (26 bodies) were previously demonstrated by PCR to come from infected

149  individuals [56], application of our screening procedure did not aim at distinguishing "positive" from "negative" samples, but at validating the

150  selection of individuals with highest microbial load (especially since some of them were sampled from 1-3 different parts of the body), at the

151  same time enhancing specificity (vs. MOTT). We used the Kay et al. dataset for verification of specificity of all applied alignment targets:

152  Borówka et al., Bouwman et al., Bos et al. and the whole genome sequence of *M. tuberculosis* H37Rv, with our algorithm aimed at detection of

153  strongest aTB outliers. While application of the Borówka et al. target sequence (with 30 bp read length cutoff) detected four samples as outliers,

154  they turned out to belong only to two individuals (bodies 68 and 92) (Supplementary Tab. 3). This validated our approach as a suitable method

155  for selecting ancient samples with highest MTBC genetic material content, especially since, despite our alignment target consisting only of

156  sequences specific exclusively for MTBC, it turned out that those four samples were also those that showed the highest ratio of aligned reads to

157  the full *M. tuberculosis* reference sequence (and thus the highest number of reads used to reconstruct the ancient genome) in the original study

158  by Kay et al. (shown there in Supplementary Tab. 3). Moreover, only the two alignment targets prepared with both specificity and sensitivity in

159  mind (Borówka et al. and Bouwman et al.) led to identification of all three samples from body 68 as outliers.

160          Subsequently, we applied the full statistical approach (with all four NGS read length bins) and the four selected genomic

161  alignment targets: full reference *Mycobacterium tuberculosis* H37Rv genome (broadest possible target), two published targets consisting of

162  rationally selected genes (applied previously to enrichment-based sequencing: Bouwman et al. and Bos et al.) as well as the novel specificity-

163  tailored target (Borówka et al.), to the Neolithic samples from Brześć Kujawski. Table 2 presents the number of reads in each read length bin

164  used for alignment with targets and statistical analysis, while Supplementary Tables 4-7 show the alignment results as numbers and ratios of

165  matching reads. Fig. 1 presents the results of statistical analysis as outlying standardised ratio values in different read length bins. Overall, the

166  expected population structure of majority of individuals with few positive reads and outlier individuals with an exceptional number of positive

reads is confirmed. However, it is immediately obvious that the composition of outlier individuals depends strongly not only on the genomic alignment target, but also on minimum length of reads used for the alignment. There are individuals who remain positive (with a high relative ratio of reads aligning to the respective target) for all four length bins (e.g. 4_BK4 for the *Mycobacterium tuberculosis* H37Rv target), i.e. the share of putative MTBC-derived sequences remains constant despite the decrease in number of analysed sequences and increase in sequence complexity. There are individuals who, despite being outliers for the bins including shorter reads, lose this status for the more restrictive bins (e.g. 55_BK4 for the Borówka et al. target), i.e. the majority of their MTBC-like sequences were of low complexity. Contrastingly, in some individuals the share of MTBC-like sequences increases above the cut off value only for bins with longer reads (e.g. 31_BK4 for the Borówka et al. target), i.e. most specifically aligned fragments are relatively long. It is again apparent that since most of this change concerns reads between 20 and 29 bp in length, the optimal threshold for read aligning to a genomic target for specificity towards MTBC is ≥30 bp. Thus, the three individuals which exceed the threshold of 1.5xSD for the MTBC-specific Borówka et al. target (17_BK4, 29_BK4 and 31_BK4) are considered with high probability to be ancient cases of MTBC infection and merit selection for further in-depth studies by a more cost-intensive approach.

Since the cut off-based detection algorithm, while robust for the presented dataset, may be less suitable for other, less homogenous groups of ancient individuals, we also set out to construct an objective, parametric testing-based outlier detection algorithm. Since the main objective of our overall study is specificity of MTBC detection, we applied this algorithm to the original Borówka et al. genomic alignment target. Based on the observation that positive read ratio tends to depend monotonically on read length bin – either consistently increasing or decreasing for outlier individuals – we decided to calculate a monotonicity parameter. We first standardised positive read ratios as percentage of average positive read ratio (without assumptions towards normal distribution, Supplementary Fig. 1) and then calculated ratios of these values for adjoining read length bins (≥25bp/≥20bp, ≥30bp/≥25bp and ≥35bp/≥30bp). The arithmetic mean of these values (Supplementary Tab. 8) depended on monotonicity of the studied relationship and had a normal distribution among individuals in our study. For outlier detection, we applied a one-tailed critical z value test on both tails on the sample. We consider the positive outliers (individuals with consistently increasing share of positively aligned reads with increasing read length) to be potential individuals with high MTBC loads, suitable for further analysis both by virtue of good mycobacterial genomic material preservation and high certainty of this material belonging to ancient MTBC. On the other hand, negative outliers may either be individuals with ancient MOTT infection (we suggest this as highly probable for 4_BK4) or samples with high proportion of short, non-specific alignments, probably due to environmental contamination (most probably 55_BK4) - to distinguish these two groups, a comparison with the more Mycobacterium-generic whole-genome alignment target is necessary (see below). This approach, while retaining the strong specificity of the cut off approach, gains increased sensitivity due to inclusion of individuals with high background of environmental sequences (low initial positive alignments in the short-read bin) which nevertheless retain specific long positively aligned sequences upon read length restriction, e.g. 21_BK4.

An immediately obvious result of our analysis was that the comparison of alignment targets constructed with different assumptions leads to surprisingly large differences in assignation of individuals. Aligning aDNA sequences versus the whole MTB genome results in identification of two strong outliers (4_BK4 and 32_BK4). The same two individuals are identified, albeit with a smaller divergence, by using the enrichment bait sequence set uses by Bouwman et al. as alignment target. Since this subset of genomic sequences was originally selected for enrichment of lineage-distinguishing polymorphisms rather than for MTB complex specificity, this result is expected and confirms the efficiency of the outlier detection method and ≥30bp as optimal read length. On the other hand, our Borówka et al. genomic subset selected on the basis of MTB complex specificity led to identification of three different individuals as outliers (17_BK4, 29_BK4 and 31_BK4), while 4_BK4 and 32_BK4 had positive read values close to average. This is even more conspicuous when positive ratio values for the two different alignment targets (whole genome and specific subset Borówka et al.) are plotted against each other (Fig. 2). In our opinion this points to the broadly recognized risk of mistakenly identifying ancient infections caused by MOTT as tuberculosis based on the extensive similarity between the respective mycobacterial genomes. While restricting the alignment target leads to loss of sensitivity due to unavoidable decrease of absolute number of

aligned reads, which is a significant problem for ancient DNA, it is offset by the increase in specificity of detection. This distinction is crucial for epidemiological hypotheses where elimination of false positives is of paramount importance. We further show this by aligning our reads to the purportedly MTBC-specific target sequences selected by Bos et al. (sequences of only two *M. tuberculosis* specific genes), where increase of specificity leads to detection of the 29_BK4 individual, but the extreme loss of sensitivity linked to minuscule absolute number of reads (the highest number of positive reads in the ≥30bp bin is 13 – see Supplementary Tab. 7) leads to high experimental noise and low reliability of assignment of individuals, and it is not recommended.

Since for two individuals which were strongly enriched in mycobacterial sequences (4_BK4 and 32_BK4) we posit the existence of an ancient MOTT infection (as they do not score highly in comparison with the specific Borówka et al. alignment target), we decided to verify if this assumption is supported by aligning the optimal read bin (≥30bp) to full genomes of other mycobacterial species as targets. Indeed, as seen in Supplementary Fig. 2, those two individuals are also strong outliers in read ratio values after aligning to the *M. marinum* genome - moreover, when plotted against read ratio values for the MTB genome, it is apparent that they show higher similarity to *M. marinum*, since they are located on the *M. marinum* side of the read ratio regression line. This finding validates our workflow in that it corroborates the usefulness of read length binning while further demonstrating the advantages of read aligning to targets selected for species discrimination (like the Borówka et al. target) which allow for immediate flagging of suspicious samples with spuriously high absolute similarity to the MTB genome. We have also attempted to verify the possibility of distinguishing samples with predominantly ancient mycobacterial sequences from samples with recent environmental MOTT contamination by performing mapDamage analysis. MapDamage analysis shows that the low absolute number of reads that map to all *M. tuberculosis* alignment targets (including the full MTB genome) in the case of our samples prevents us from drawing meaningful conclusions in this regard (even for the samples with highest read numbers - 4_BK4, 32_BK4, 17_BK4, 29_BK4, 31_BK4). For general confirmation of ancient status of analysed reads, MapDamage analysis was performed for human sequences (aligning to the human genome build 37) and is presented in Supplementary Fig. 3 for all 6 individuals with potential MOTT and MTBC infections. Since the samples with potential MOTT infection (4_BK4, 32_BK4) included a substantial number of reads that aligned to the *M. marinum* genome (Supplementary Tab. 9), we were also able to perform MapDamage analysis for these reads (Supplementary Fig. 4), confirming the ancient character of mycobacterial sequences.

**Discussion**

The evolutionary and ecological complexity of mycobacteria, including the existence of a group of closely related pathogens known as the *Mycobacterium tuberculosis* complex, consists of a large number of more distantly related human and animal pathogens causing diseases other than tuberculosis, and an abundance of free-living (including soil- and water-borne) mycobacterial species in the environment. These all contribute to the difficulty in the unequivocal determination of ancient tuberculosis on the basis of MTBC aDNA. Present-day paleoepidemiology uses tools of classical biological anthropology as well as modern clinical diagnostics at the molecular level. Morphological diagnosis of tuberculosis is based on certain bone changes, especially those described as Pott's disease. This approach is not optimal from the point of view of sensitivity, since bone lesions are present only in 2% of all cases of tuberculosis infection and 10-20% of cases of extrapulmonary tuberculosis [41, 57]. The specificity of this tool is also relatively low: even in the case of Pott's disease, which is regarded by paleopathologists as the pathognomonic skeletal signature of TB, there are several lesions that may be difficult to differentiate from TB in archaeological skeletal remains. . In spite of that limitations, osteological analysis is often the main starting point of a study and cannot be disregarded. However, in our study the occurrence of bone lesions that could be linked in any way with tuberculosis did not correlate with the results of our genetic analyses. There are two possible explanations for this fact. First, the bone changes were not caused by tuberculosis, which is in accordance with a lack of pathognomonic characteristics of the disease on the skeleton alone, as was clarified before; it applies primarily to the graves 12_BK4, 18_BK4, 47_BK4, and 73_BK4. It may also be that the preservation of MTBC aDNA was too poor to pass the sensitivity/specificity threshold of the method proposed here.

244      Among molecular techniques which are used for diagnosis of ancient tuberculosis cases, both biochemical methods based on mass

245      spectrometry and PCR amplification of marker sequences have been successfully used in literature, e.g. for preliminary description of the

246      Hungarian mummies used subsequently to reconstruct aTB genomes [48, 56]. However, both these groups of methods suffer from a number of

247      drawbacks which make them less useful in an ancient epidemiological context than in a contemporary one: environmental contamination from

248      modern soil mycobacteria can overwhelm both traces of ancient MTBC mycolic acids and less specific PCR amplicons, while strong care must

249      be taken to prevent in-lab cross-contamination with genuine MTBC samples. Therefore, NGS has a number of advantages in diagnosis of ancient

250      tuberculosis, having the potential to be both highly sensitive and highly specific; but the balance between sensitivity and specificity depends on

251      the selection of reference genomic sequences and crucially on the method of alignment. A large quantity of generated data allows potentially to

252      detect ancient mycobacteria selectively, unequivocally and semi-quantitatively, while making possible additional analyses such as preservation

253      period-related DNA damage pattern detection (e.g. mapDamage [58, 59], phylogenetic analysis of genetic kinship [50] or even full genome

254      reconstruction [48]). Due to small absolute amounts of actual ancient pathogen DNA in most types of human body samples, a common approach

255      is to use pre-sequencing enrichment (usually using probe capture, e.g. [50]). Only in bodies preserved in exceptional, isolated conditions, such as

256      the Hungarian mummies from a 18th century crypt, was a non-enriched metagenomics approach used [48]. Use of enrichment techniques

257      strongly increases sensitivity, but comes with its own drawbacks (apart from increased cost), the most relevant of which is the need to pre-design

258      a set of sequences (probes or primers) that will define and limit the scope of subsequently obtained NGS data. A full metagenome approach is

259      often more relevant when dealing with a highly ancient sample like in the present study, when neither the infection prevalence nor the pathogen

260      identity are known to any precision and a preliminary NGS study is needed for formulation of specific hypotheses and pre-selection of

261      individuals for further analysis.

262      However, in the case of ancient MTBC (especially samples more than a thousand years old), specificity is a more important consideration

263      than sensitivity. While modern MTBC contamination in the laboratory is a risk factor, it would not mask ancient data in a semi-quantitative

264      study and would be obvious if DNA damage analysis were performed.  A more important consideration is the possible presence of ancient

265      MOTT which can be unpredictably genetically similar to MTBC. The sources of these MOTT can be either soil contamination (including dead

266      animals) which could have happened at any time since inhumation (preventing reliable elimination by DNA damage analysis), or actual ancient

267      MOTT which were pathogenic/infectious/commensal to ancient humans. Thus, the design of sequencing analysis workflow has to take into

268      account the necessity to filter out unknown related sequences that are not derived from MTBC - this was the main rationale behind the design of

269      our study. While contamination with mycobacterial sequences within the laboratory (amplicons, genuine Mycobacterium DNA) can be

270      prevented by correct workflow (separation of pre- and post-PCR areas etc.), equipment and strict procedures, contamination by environmental

271      DNA is inescapable and has to be taken into account in the case of archaeological bone samples preserved by inhumation. Since for ancient

272      samples direct contact of bones with the environment has lasted for a very long time (unlike more recent samples from vault inhumation),

273      mycobacterial DNA derived from environmental (soil) MOTT can have undergone accretion in bones throughout this period, with some of it

274      ancient enough to be indistinguishable in terms of location and state of preservation from DNA of infectious microbes buried with the body. All

275      MTBC are obligate pathogens and thus are an unlikely source of environmental contamination of ancient samples. Therefore, for preliminary

276      identification of potentially interesting samples in ancient inhumated bones, specificity in methods of detection of ancient infectious agents from

277      this group should be developed towards exclusion of MOTT, with distinction between members of MTBC as a secondary, much less important

278      goal. Since MTBC also share a very high proportion of coding sequences, achieving specificity for *M. tuberculosis* s.s. could occur only by

279      drastically limiting the size of the reference marker sequence, thus leading to very low sensitivity, especially for usually highly degraded aDNA.

280      Moreover, the division of MTBC into lineages is not entirely concordant with classical taxonomic division into species, so attempting an

281      artificial distinction between some lineage groups based on accumulated NGS data would not be recommended. Our approach is designed as a

282      relatively low-cost, first-pass classification of ancient samples based on whole-metagenome NGS data. When a highly specific method like the

283  one we propose is used to identify likely ancient MTBC infection, potential lineage determination or any other phylogenetic studies (in pre-

284  selected samples) should proceed by other methods developed specifically for this purpose, based on the presence of lineage-specific

285  polymorphisms (with the caveat that enrichment for specificity-related sequences before NGS will certainly lead to loss of the majority of

286  phylogenetically important loci, so a full metagenomic sequencing round with sufficient coverage is inevitable).

287  We postulate that a combination of read length-based genomic alignment analysis and a careful knowledge-based selection of the

288  alignment target makes it possible to achieve relatively high specificity of aTB detection against all potential false positive sources. Therefore, a

289  robust tool for specifically identifying NGS-derived sequences that belong to ancient MTBC with high confidence is a priority task in molecular

290  paleoanthropology. Even more relevant to paleoanthropological studies, confusion between MOTT and MTBC can lead to spurious

291  identification of ancient individuals as tuberculosis sufferers or carriers, invalidating conclusions relevant to paleoepidemiology. We

292  demonstrate that read length selection is not only highly relevant (as has been shown before and by us, only reads above ca. 30 bp can be used

293  with high confidence), but when a statistics-based approach to multiple length thresholds is used, it can yield a substantial increase in specificity

294  of MTBC identification. At the same time, selection of pre-filtered alignment target, with combined knowledge-based (selection of transcribed

295  sequences) and automated (exclusion of sequences aligning with MOTT genomes) delineation of MTBC-specific sequences (which we call the

296  Borówka et al. target), makes it possible to perform in-depth specificity analysis by comparing the alignments of *in silico* fragmented

297  mycobacterial genomes (mimicking actual NGS data). Combining the novel alignment target and the read length binning approach, we were able

298  to select with high confidence three ancient individuals with probable ancient MTBC infection and two further individuals with highly probably

299  ancient mycobacteriosis caused by MOTT (which would be misidentified as tuberculosis if another alignment target or to short reads were taken

300  into account). Of course the limitations of our data make these identifications preliminary and another round of directed (e.g. enrichment-based)

301  sequencing would be required both for positive identification of the infectious agent and for potential phylogenetical analysis of its spatial and/or

302  temporal kinship. However, in our case read length analysis allowed us to suggest *M. marinum* as the potential ancient infectious agent based on

303  statistical analysis; obviously, positive confirmation of this diagnosis would require tools that are currently unavailable such as proven *M.*

304  *marinum*-specific enrichment probes as well as a much better sequence coverage that could be achieved in a preliminary study (Supplementary

305  Fig. 2). Still, this possible pathogen identification is not at odds with the archeological context as the inhumation site is next to a lake (Smetowo)

306  and within a geographical region rich in post-glacial lakes (Kujawy), so some individuals could have had routine professional contact with fish.

307  Our combined procedures used robust tools but cannot be treated as definite proof. Our samples are relatively old (in comparison to most other

308  ancient tuberculosis cases studied by molecular means before) and thus the absolute read numbers from an unbiased NGS approach is low. We

309  demonstrate that this disadvantage makes it relatively difficult to perform DNA damage analysis (except for samples with a very high absolute

310  number of reads). However, we provide a consistent proof of concept for a tool which allows relatively cheap and unbiased selection of samples

311  (e.g. individuals) for further analysis, e.g. by enrichment capture NGS. Thus, we suggest that it is possible to use global NGS results from

312  ancient samples as an economical pre-screening tool for more complex methods, while applying bioinformatic tools to maximise the number of

313  reliable conclusions that can be drawn from a limited dataset.

314  **Methods**

315  **Ancient DNA extractions**

316  A dedicated ancient DNA sample preparation facility at the University of Lodz was used, taking standard precautions to avoid any

317  contamination. All disposable materials, buffers, water, clean room surfaces and bone material, were UV-irradiated for at least 30 minutes before

318  any subsequent steps were taken. The fragments of bone material were isolated using Dremel disks, (USA), surface-cleaned, UV-irradiated for

319  7.5 minutes on each side, and ground into a fine powder, further used for DNA extraction procedures following the protocol of Dabney et al.

320  with modifications [60-62]. Ancient DNA was successfully isolated from all bone samples (See Supplementary Fig. 3). Illumina libraries were

321  prepared in separate facility, according to Meyer et al. protocol [63]  with modifications proposed by Gamba et al. [60] without UDG treatment

322  of the samples. All libraries were subjected to the screening next-generation sequencing on the Illumina Nextseq 500 platform (100bp single-end

323  sequencing), yielding between 2.2 and 33.9 million reads per individual (median number of reads after incomplete and truncated read trimming –

324  16.9 million reads per individual, Tab. 2). This dataset contains ancient human sequences from the deceased individuals, ancient microbial

325  sequences from parasites, pathogens, commensals or symbionts of the deceased individuals, as well as genomic sequences from environmental

326  organisms (mainly microbes, but also potentially higher Eukaryotes), to which the skeletal remains were exposed *post-mortem*.


327  **Bioinformatical procedures**

328  Raw NGS reads were subjected to standard quality processing such as trimming and adapter sequence removal (-q 30 --phred33 --

329  illumina --length 20), using the Trim Galore! software package [64]. Since the predominant expected type of sequence in skeletal samples is

330  ancient human genomic DNA and its presence would unnecessarily complicate our analysis, the read datasets were subsequently subjected to

331  filtering by alignment to the standard (hg19) human genome reference sequence. This alignment was performed using the BWA_aln algorithm (-

332  n 0.04, -l 1000), with duplicate removal, using the AGAT software tool - ocwrapper3mt.py script [65]. Any read which aligned without gaps

333  within the default mismatch rate (dependent on sequence length, e.g. 2 mismatches per 17 bp) was eliminated from the sample dataset.

334  Subsequently, separate sub-datasets (bins) of reads were generated on the basis of (trimmed) read length: minimal read length threshold -≥20bp,

335  ≥25bp, ≥30bp and ≥35bp. These datasets were used for alignment to reference targets. These procedures were applied also to the Kay et al

336  dataset, used for the Borówka et al. method verification.

337  Estimation of terminal base deamination damage pattern was done by using mapDamage2.0 analysis with specifying a length (-l) of

338  75 bp (Supplementary Fig.3 and Supplementary Fig. 4).


340  **Query sequence preparation**

341  Selected 18 reference *Mycobacterial* genomes, including 5 of *M. tuberculosis* complex (underlined): *M. abscessus, <u>M. africanum</u>, M.*

342  *asiaticum, M. avium, <u>M. bovis</u>, <u>M. caprae</u>, M. fortuitum, M. haemophilum, M. kansasii, M. leprae, M. marinum, <u>M. microti</u>, M. phlei, M. simiae,*

343  *M. smegmati, <u>M. tuberculosis</u>, M. ulcerans, M. xenopi* were used. Nucleotide sequences of each organism have been subjected to fragmentation

344  with FA_TOOL script (small_tool.py) [66] respectively for 20 bp, 25 bp, 30 bp and 35 bp-long fragments and allocated in same manner to

345  length bins. Further, fragmented genomes were used for specificity testing of each constructed target which allowed to overcome the problem of

346  very short and non-specific fragments with threshold estimation.


**Verification of specificity and sensitivity of NGS screening method**

348  Due to the lack of available NGS data of positive M. tuberculosis cases, we tested in-silico methods by using the Kay et al. (2015) dataset

349  (PRJEB7454), derived from Hungarian mummies tissue microbiome sequencing. SRA files for each sample were identified and downloaded,

350  further fastq files passed through trimming with deprivation of the adapter sequences [65]. Raw sequencing files were conducted to human

351  genome reference sequence (hg19) filtration in spite the fact that host DNA material could be dominant in the sample. Alignment was performed

352  to the tested targets M. tuberculosis H37Rv, Borówka et al., Bos et al., and Bouwman et al. using the AGAT software tool [65]. Statistics for

353  each individual are presented in Supplementary Table 3. Summarized results of aTB cases from Brześć Kujawski are included in Supplementary

354  Tables 4-7.

**Statistical processing and parametric testing-based outlier detection algorithm**

Collected unmapped sequences from the original dataset, as well as from the Kay et al. dataset, were aligned to constructed marker sequences: *M. tuberculosis H37Rv*, Borówka et al. (Supplementary Table 10), Bos et al., and Bouwman et al. with application of experimentally determined minimal read length threshold ≥17 bp, ≥20bp, ≥25bp, ≥30bp and ≥35bp for detection of potential ancient MTBC cases. For detection of outlier individuals with high microbial load/positive read ratio, we standardised read ratio values to normal distribution parameters (arithmetic mean and standard deviation) and, as a further step in the aTB detection algorithm, applied a typical cut off value of 1.5xSD to detect outliers, postulating these to be candidates for active tuberculosis.

Based on the observation that positive read ratio tends to depend monotonically on read length bin – either consistently increasing or decreasing for outlier individuals – we decided to calculate a monotonicity parameter. We first standardised positive read ratios as percentage of average positive read ratio and then calculated ratios of these values for adjoining read length bins (≥25bp/≥20bp, ≥30bp/≥25bp and ≥35bp/≥30bp). For outlier detection, we applied a one-tailed critical z value test on both tails of the sample. We consider the positive outliers (individuals with consistently increasing share of positively aligned reads with increasing read length) to be confirmed ancient tuberculosis sufferers (See Supplementary tables 3-7).

## Availability of supporting data and materials

The datasets supporting the conclusions of this article are available under the NCBI repository project "Identification of ancient tuberculosis in human archaeological remains" (acc. num. PRJNA422903) including Biosamples and related Sequence Read Archive (SRA). Other supporting data are available via the Gigascience database, GigaDB [67].

## Additional files

**Borówka_et_al_Supplemetary_Tables.xls**

**Borówka_et_al_Supplementary_Tables_legends.doc**

**Borówka_et_al_Supplementary_Figures.pdf**

## Declarations

**Abbreviations**

**aDNA – Ancient DNA**

**aTB – Ancient tuberculosis**

**NGS – Next Generation Sequencing**

**MTBC – *Mycobacterium Tuberculosis* Complex**

**MOTT – Mycobacteria other than tuberculosis**

**SRA - Sequence Read Archive**

**Ethics approval and consent to participate**

Not applicable.

**Author's contributions**

P.B. and D.S. conceived the study, were responsible for extraction of aDNA, preparation of NGS libraries and Next Generation Sequencing of

samples. P.B, D.S and Ł.P analyzed the data, discussed the results, and wrote the manuscript. Ł.P. participated in the statistical analysis and

figure preparation. B.M wrote and ran AGAT primary analysis. B.B-S. precipitated in sample selection and preparation for laboratory phase.

J.D., WL analyzed the samples for pathological changes, participated in the study design, analyzed and discussed the data, and participated in

drafting the manuscript. E.Ż. participated in the study design, analyzed and discussed the data, and participated in drafting the manuscript. D.S.

coordinated studies and was responsible for the final version of the manuscript; all authors read and approved the final manuscript.

**References**

1.  Barrett, R., et al., *Emerging and re-emerging infectious diseases: the third epidemiologic transition.* Annual review of anthropology, 1998. **27**(1): p. 247-271.
2.  Armelagos, G.J. and M.N. Cohen, *Paleopathology at the Origins of Agriculture.* 1984: Academic Press Orlando (FL).
3.  Armelagos, G.J., A.H. Goodman, and K.H. Jacobs, *The origins of agriculture: Population growth during a period of declining health.* Population & Environment, 1991. **13**(1): p. 9-22.
4.  Weiss, R.A. and A.J. McMichael, *Social and environmental risk factors in the emergence of infectious diseases.* Nature medicine, 2004. **10**(12s): p. S70.
5.  Ortner, D.J. and W. Putschar, *Identification of Pathological Conditions in Human Skeletal Remains.* Smithsonian Contributions to Anthropology, 1985. **28**.
6.  Aufderheide, A.C., Rodriguez-Martin, Conrado and O. Langsjoen, *The Cambridge encyclopedia of human paleopathology.* Vol. 478. 1998: Cambridge University Press Cambridge.
7.  Roberts, C., Buikstra J. E.,, *The bioarcheology of tuberculosis: a global perspective on re-emerging disease.* University Press of Florida, Gainesville, FL, 2003.
8.  Holloway, K.L., et al., *Skeletal lesions in human tuberculosis may sometimes heal: an aid to palaeopathological diagnoses.* PLoS One, 2013. **8**(4): p. e62798.
9.  Steinbock, R.T., *Paleopathological diagnosis and interpretation: bone diseases in ancient human populations.* 1976: Charles C Thomas Pub Limited.
10. Donoghue, H., et al., *Mycobacterium tuberculosis complex DNA in calcified pleura from remains 1400 years old.* Letters in Applied Microbiology, 1998. **27**(5): p. 265-269.
11. Donoghue, H.D., *Palaeomicrobiology of tuberculosis*, in *Paleomicrobiology.* 2008, Springer. p. 75-97.
12. Donoghue, H.D., *Human tuberculosisis-an ancient disease, as elucidated by ancient microbial biomolecules.* Microbes and infection, 2009. **11**(14): p. 1156-1162.
13. Redman, J.E., et al., *Mycocerosic acid biomarkers for the diagnosis of tuberculosis in the Coimbra Skeletal Collection.* Tuberculosis (Edinb), 2009. **89**(4): p. 267-77.
14. Mark, L., et al., *High-throughput mass spectrometric analysis of 1400-year-old mycolic acids as biomarkers for ancient tuberculosis infection.* Journal of Archaeological Science, 2010. **37**(2): p. 302-305.
15. Minnikin, D.E., et al., *The interplay of DNA and lipid biomarkers in the detection of tuberculosis and leprosy in mummies and other skeletal remains.* 2011, Verlag Dr. Friedrich Pfeil.
16. Tran, T., et al., *Beyond ancient microbial DNA: nonnucleotidic biomolecules for paleomicrobiology.* Biotechniques, 2011. **50**(6): p. 370-380.
17. Masson, M., et al., *Osteological and biomolecular evidence of a 7000-year-old case of hypertrophic pulmonary osteopathy secondary to tuberculosis from neolithic hungary.* PLoS One, 2013. **8**(10): p. e78252.
18. Borowska-Strugińska, B., et al., *Mycolic acids as markers of osseous tuberculosis in the Neolithic skeleton from Kujawy region (central Poland).* AnthropologicAl review, 2014. **77**(2): p. 137-149.
19. Gernaey, A.M., et al., *Mycolic acids and ancient DNA confirm an osteological diagnosis of tuberculosis.* Tuberculosis (Edinb), 2001.

446      **81**(4): p. 259-65.
447    20.    Boros-Major, A., et al., *New perspectives in biomolecular paleopathology of ancient tuberculosis: a proteomic approach.* Journal of
448        Archaeological Science, 2011. **38**(1): p. 197-201.
449    21.    Harkins, K.M., et al., *Screening ancient tuberculosis with qPCR: challenges and opportunities.* Philos Trans R Soc Lond B Biol Sci,
450        2015. **370**(1660): p. 20130622.
451    22.    Campana, M.G., et al., *False positives complicate ancient pathogen identifications using high-throughput shotgun sequencing.* BMC
452        research notes, 2014. **7**(1): p. 111.
453    23.    Andam, C.P., et al., *Microbial Genomics of Ancient Plagues and Outbreaks.* Trends Microbiol, 2016. **24**(12): p. 978-990.
454    24.    Bouwman, A.S. and T.A. Brown, *The limits of biomolecular palaeopathology: ancient DNA cannot be used to study venereal syphilis.*
455        Journal of Archaeological Science, 2005. **32**(5): p. 703-713.
456    25.    Tsangaras, K. and A.D. Greenwood, *Museums and disease: using tissue archive and museum samples to study pathogens.* Ann Anat,
457        2012. **194**(1): p. 58-73.
458    26.    Müller, R., C.A. Roberts, and T.A. Brown, *Complications in the study of ancient tuberculosis: Presence of environmental bacteria in
459        human archaeological remains.* Journal of Archaeological Science, 2016. **68**: p. 5-11.
460    27.    Frothingham, R., H.G. Hills, and K.H. Wilson, *Extensive DNA sequence conservation throughout the Mycobacterium tuberculosis
461        complex.* Journal of clinical microbiology, 1994. **32**(7): p. 1639-1643.
462    28.    Brites, D. and S. Gagneux, *Co-evolution of Mycobacterium tuberculosis and Homo sapiens.* Immunol Rev, 2015. **264**(1): p. 6-24.
463    29.    Minnikin, D.E., et al., *Molecular biomarkers for ancient tuberculosis*, in *Understanding tuberculosis-deciphering the secret life of the
464        bacilli.* 2012, InTech.
465    30.    Minnikin, D.E., et al., *Essentials in the use of mycolic acid biomarkers for tuberculosis detection: response to â€śHigh-throughput mass
466        spectrometric analysis of 1400-year-old mycolic acids as biomarkers for ancient tuberculosis infectionâ€' by.* Journal of Archaeological
467        Science, 2010. **37**(10): p. 2407-2412.
468    31.    Donoghue, H.D., *Insights gained from ancient biomolecules into past and present tuberculosis—a personal perspective.* International
469        Journal of Infectious Diseases, 2017. **56**: p. 176-180.
470    32.    Donoghue, H.D., et al., *Ancient DNA analysis - An established technique in charting the evolution of tuberculosis and leprosy.*
471        Tuberculosis (Edinb), 2015. **95 Suppl 1**: p. S140-4.
472    33.    Spigelman, M. and E. Lemma, *The use of the polymerase chain reaction (PCR) to detect Mycobacterium tuberculosis in ancient
473        skeletons.* International Journal of Osteoarchaeology, 1993. **3**(2): p. 137-143.
474    34.    Müller, R., C.A. Roberts, and T.A. Brown, *Complications in the study of ancient tuberculosis: non-specificity of IS6110 PCRs.* STAR:
475        Science & Technology of Archaeological Research, 2015. **1**(1): p. 1-8.
476    35.    Pääbo, S., *Molecular cloning of ancient Egyptian mummy DNA.* Nature, 1985. **314**(6012): p. 644-645.
477    36.    Salo, W.L., et al., *Identification of Mycobacterium tuberculosis DNA in a pre-Columbian Peruvian mummy.* Proceedings of the National
478        Academy of Sciences, 1994. **91**(6): p. 2091-2094.
479    37.    Hershkovitz, I., et al., *Detection and molecular characterization of 9,000-year-old Mycobacterium tuberculosis from a Neolithic
480        settlement in the Eastern Mediterranean.* PLoS One, 2008. **3**(10): p. e3426.
481    38.    Mays, S. and G.M. Taylor, *A first prehistoric case of tuberculosis from Britain.* International Journal of Osteoarchaeology, 2003. **13**(4): p.
482        189-196.
483    39.    Zink, A.R., W. Grabner, and A.G. Nerlich, *Molecular identification of human tuberculosis in recent and historic bone tissue samples: The
484        role of molecular techniques for the study of historic tuberculosis.* Am J Phys Anthropol, 2005. **126**(1): p. 32-47.
485    40.    Dziadek, J. and A. Sajduda, *Specificity of insertion sequence-based PCR assays for Mycobacterium tuberculosis complex.* The
486        International Journal of Tuberculosis and Lung Disease, 2001. **5**(6): p. 569-574.
487    41.    Teo, H.E. and W.C. Peh, *Skeletal tuberculosis in children.* Pediatric radiology, 2004. **34**(11): p. 853-860.
488    42.    Kent, L., et al., *Demonstration of homology between IS6110 of Mycobacterium tuberculosis and DNAs of other Mycobacterium spp.?*
489        Journal of clinical microbiology, 1995. **33**(9): p. 2290-2293.
490    43.    McHugh, T., L. Newport, and S. Gillespie, *IS6110 homologs are present in multiple copies in mycobacteria other than tuberculosis-
491        causing mycobacteria.* Journal of clinical microbiology, 1997. **35**(7): p. 1769-1771.
492    44.    Picardeau, M., et al., *Genotypic characterization of five subspecies of Mycobacterium kansasii.* J Clin Microbiol, 1997. **35**(1): p. 25-32.
493    45.    Picardeau, M., et al., *Mycobacterium xenopi IS1395, a novel insertion sequence expanding the IS256 family.* Microbiology, 1996.
494        **142**(9): p. 2453-2461.
495    46.    Savelkoul, P.H., et al., *Detection of Mycobacterium tuberculosis complex with real time PCR: comparison of different primer-probe sets
496        based on the IS6110 element.* Journal of microbiological methods, 2006. **66**(1): p. 177-180.
497    47.    Rasmussen, S., et al., *Early divergent strains of Yersinia pestis in Eurasia 5,000 years ago.* Cell, 2015. **163**(3): p. 571-82.
498    48.    Kay, G.L., et al., *Eighteenth-century genomes show that mixed infections were common at time of peak tuberculosis in Europe.* Nature
499        communications, 2015. **6**.
500    49.    Chan, J.Z., et al., *Metagenomic analysis of tuberculosis in a mummy.* N Engl J Med, 2013. **369**(3): p. 289-90.
501    50.    Bos, K.I., et al., *Pre-Columbian mycobacterial genomes reveal seals as a source of New World human tuberculosis.* Nature, 2014.
502        **514**(7523): p. 494-7.
503    51.    Bos, K.I., et al., *Parallel detection of ancient pathogens via array-based DNA capture.* Philos Trans R Soc Lond B Biol Sci, 2015.
504        **370**(1660): p. 20130375.
505    52.    Bouwman, A.S., et al., *Genotype of a historic strain of Mycobacterium tuberculosis.* Proceedings of the National Academy of Sciences,
506        2012. **109**(45): p. 18511-18516.
507    53.    Lorkiewicz, W., et al., *Between the Baltic and Danubian worlds: the genetic affinities of a middle neolithic population from central
508        Poland.* PLoS One, 2015. **10**(2): p. e0118316.
509    54.    Aziz, R.K., et al., *The RAST Server: rapid annotations using subsystems technology.* BMC Genomics, 2008. **9**(1): p. 75.
510    55.    Djelouadji, Z., D. Raoult, and M. Drancourt, *Palaeogenomics of Mycobacterium tuberculosis: epidemic bursts with a degrading genome.*
511        Lancet Infect Dis, 2011. **11**(8): p. 641-50.
512    56.    Fletcher, H.A., et al., *Widespread occurrence of Mycobacterium tuberculosis DNA from 18thâ€"19th century Hungarians.* American
513        Journal of Physical Anthropology: The Official Publication of the American Association of Physical Anthropologists, 2003. **120**(2): p.
514        144-152.
515    57.    Peto, H.M., et al., *Epidemiology of extrapulmonary tuberculosis in the United States, 1993-2006.* Clinical Infectious Diseases, 2009.
516        **49**(9): p. 1350-1357.
517    58.    Ginolhac, A., et al., *mapDamage: testing for damage patterns in ancient DNA sequences.* Bioinformatics, 2011. **27**(15): p. 2153-5.
518    59.    Jonsson, H., et al., *mapDamage2.0: fast approximate Bayesian estimates of ancient DNA damage parameters.* Bioinformatics, 2013.
519        **29**(13): p. 1682-4.
520    60.    Gamba, C., et al., *Genome flux and stasis in a five millennium transect of European prehistory.* Nature communications, 2014. **5**: p.
521        5257.
522    61.    Pinhasi, R., et al., *Optimal ancient DNA yields from the inner ear part of the human petrous bone.* PLoS One, 2015. **10**(6): p. e0129102.
523    62.    Fernandes, D., et al., *A genomic Neolithic time transect of hunter-farmer admixture in central Poland.* Scientific reports, 2018. **8**(1): p.
524        14879.
525    63.    Meyer, M. and M. Kircher, *Illumina sequencing library preparation for highly multiplexed target capture and sequencing.* Cold Spring
526        Harbor Protocols, 2010. **2010**(6): p. pdb. prot5448.
527    64.    Krueger, F., *Trim Galore!: A wrapper tool around Cutadapt and FastQC to consistently apply quality and adapter trimming to FastQ files.*
528        2015.
529    65.    Marciniak, B., P. Borówka, and D. Strapagiel, *AGAT tool - Ancient Genomes Analysis Tool.* 2016.
530    66.    Marciniak, B., P. Borówka, and D. Strapagiel, *FA_TOOL-simple command line tool for fasta file editing.* 2016.

531    67.    Borówka P, Pułaski L, Marciniak B, Borowska-Strugińska B, Dziadek J, Żądzińska E et al. Supporting data for "Screening methods for
532             detection of ancient Mycobacterium tuberculosis complex fingerprints in NGS data derived from skeletal samples" GigaScience
533             Database 2019. http://dx.doi.org/10.5524/100598
534
535

536

538    Table 1. Number of genomic k-mers from MTBC and MOTT members after initial hg19 clearing step matching selected targets, with k-mer length distinction (≥20bp, ≥25bp, ≥30bp, ≥35bp). with estimation of percentage of
539    k-mers from a given mycobaterial genomes matching the *M. tuberculosis* genome for query length ≥30 and ≥35.
540

| | k-mer length | | Query length ≥20 | | | | | Query length ≥25 | | | | | Query length ≥30 | | | | | Query length ≥35 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Alignment target | Genome length (bp) | % of sequences mapped to M. tuberculosis genome | Full genome | Borowka et al. | Bos et al. | Bouwman et al. | % of sequences mapped to M. tuberculosis genome | Full genome | Borowka et al. | Bos et al. | Bouwman et al. | % of sequences mapped to M. tuberculosis genome | Full genome | Borowka et al. | Bos et al. | Bouwman et al. | % of sequences mapped to M. tuberculosis genome | Full genome | Borowka et al. | Bos et al. | Bouwman et al. |
| **Species group** | | | | | | | | | | | | | | | | | | | | | | |
| | *M. leprae* | 3268203 | 3.19% | 140922 | 19736 | 101 | 1683 | 4.88% | 215257 | 4349 | 85 | 2240 | 2.61% | 115138 | 1430 | 26 | 1201 | 1.45% | 63860 | 543 | 6 | 715 |
| **MOTT** | *M. abscessus* | 5067172 | 5.26% | 232228 | 46530 | 158 | 2915 | 2.87% | 126769 | 2816 | 103 | 1890 | 1.39% | 61160 | 283 | 46 | 1065 | 0.75% | 33175 | 62 | 14 | 644 |
| | *M. smegmatis* | 6988209 | 11.19% | 493570 | 107339 | 543 | 6917 | 5.68% | 250537 | 7793 | 340 | 2919 | 2.88% | 127219 | 1187 | 162 | 1610 | 1.64% | 72286 | 262 | 65 | 944 |
| | *M. fortuitum* | 6254616 | 8.48% | 374030 | 77785 | 291 | 5208 | 5.22% | 230382 | 5940 | 131 | 2774 | 2.69% | 118483 | 958 | 40 | 1534 | 1.53% | 67463 | 236 | 16 | 916 |
| | *M. phlei* | 5349645 | 8.98% | 396255 | 88582 | 391 | 5909 | 6.57% | 289788 | 9912 | 157 | 3597 | 3.45% | 152331 | 1665 | 97 | 1951 | 2.03% | 89593 | 377 | 56 | 1176 |
| | *M. simiae* | 5938797 | 9.33% | 411677 | 80142 | 339 | 5414 | 9.51% | 419641 | 12734 | 197 | 4578 | 5.35% | 235800 | 3904 | 93 | 2702 | 3.15% | 139050 | 1450 | 33 | 1575 |
| | *M. asiaticum* | 5910436 | 9.00% | 396854 | 76829 | 413 | 5597 | 10.69% | 471493 | 19780 | 392 | 5022 | 0.00% | 265638 | 5366 | 186 | 2806 | 3.54% | 156188 | 1531 | 71 | 1706 |
| | *M. xenopi* | 4434836 | 7.14% | 314850 | 60482 | 262 | 4336 | 8.17% | 360395 | 11534 | 207 | 4105 | 0.00% | 200395 | 3233 | 120 | 2126 | 2.62% | 115687 | 1060 | 68 | 1235 |
| | *M. marinum* | 6660144 | 9.48% | 418304 | 82499 | 466 | 5715 | 14.08% | 621166 | 52438 | 707 | 6366 | 7.88% | 347459 | 16301 | 266 | 3465 | 4.49% | 198076 | 4208 | 88 | 2046 |
| | *M. ulcerans* | 5805761 | 8.26% | 364492 | 71682 | 339 | 4800 | 12.26% | 540893 | 36626 | 448 | 5543 | 6.94% | 306075 | 10994 | 160 | 3094 | 4.04% | 178217 | 3088 | 61 | 1886 |
| | *M. kansasii* | 6402301 | 10.51% | 463445 | 89051 | 472 | 6353 | 15.93% | 702577 | 39990 | 596 | 7181 | 9.54% | 420814 | 13458 | 278 | 4032 | 5.82% | 256893 | 4132 | 129 | 2373 |
| | *M. avium* | 4829781 | 8.07% | 356159 | 71620 | 322 | 5128 | 12.08% | 532953 | 16610 | 194 | 5331 | 7.31% | 322606 | 4752 | 110 | 3271 | 4.58% | 202232 | 1475 | 65 | 2095 |
| | *M. haemophilum* | 4235765 | 7.08% | 312375 | 52214 | 274 | 4137 | 13.05% | 575862 | 22641 | 540 | 6284 | 7.98% | 352034 | 8023 | 374 | 3703 | 4.94% | 217744 | 2893 | 254 | 2322 |
| **MTBC** | *M. caprae* | 4288871 | 17.53% | 773238 | 181627 | 598 | 9935 | 94.85% | 4184378 | 734742 | 2306 | 37814 | 96.27% | 4245996 | 725608 | 2253 | 35935 | 96.21% | 4244109 | 713211 | 2214 | 34394 |
| | *M. microti* | 4370115 | 17.81% | 785606 | 188016 | 825 | 10498 | 96.71% | 4266542 | 772527 | 3989 | 40576 | 98.17% | 4330722 | 771950 | 3873 | 38507 | 98.12% | 4328596 | 758572 | 3785 | 36841 |
| | *M. africanum* | 4389314 | 17.87% | 788161 | 186939 | 850 | 10494 | 97.15% | 4285645 | 764554 | 4038 | 40740 | 98.63% | 4350937 | 764150 | 3893 | 38685 | 98.60% | 4349503 | 751103 | 3796 | 37018 |
| | *M. bovis* | 4345492 | 17.72% | 781857 | 184148 | 592 | 10161 | 96.31% | 4248729 | 750458 | 2304 | 39042 | 97.79% | 4313964 | 749050 | 2252 | 36990 | 97.76% | 4312566 | 735993 | 2213 | 35367 |
| | *M. tuberculosis* | 4411532 | 18.07% | 797099 | 192022 | 833 | 10844 | 98.41% | 4341179 | 791071 | 3947 | 42253 | 99.97% | 4410355 | 792717 | 3851 | 40180 | 100.00% | 4411458 | 779771 | 3777 | 38435 |

541

542    Table 2. Number of reads (per individual) used for alignment and statistical processing.

543

| Sample ID | Raw reads | Trimmed reads | Average read length | Non-human reads | Non-human reads | Non-human reads | Non-human reads |
|---|---|---|---|---|---|---|---|

| | | | | (>20) | (>25) | (>30) | (>35) |
|---|---|---|---|---|---|---|---|
| 1_BK4 | 17507911 | 17038725 | 57.6 | 16977024 | 16902603 | 16378765 | 15191086 |
| 4_BK4 | 18816573 | 18215498 | 51.7 | 18095660 | 17960494 | 17086604 | 15246279 |
| 6_BK4 | 16322105 | 15815995 | 55.0 | 15551094 | 15427193 | 14682610 | 13220243 |
| 7_BK4 | 2231650 | 2160395 | 59.7 | 2102955 | 2095297 | 2047913 | 1936435 |
| 9_BK4 | 14974057 | 14503433 | 53.5 | 14240738 | 14085752 | 13149549 | 11600503 |
| 11A_BK4 | 16432267 | 16000777 | 58.0 | 15766313 | 15695767 | 15172161 | 14034604 |
| 11B_BK4 | 18522995 | 18078222 | 55.7 | 725913 | 718941 | 674747 | 597601 |
| 12_BK4 | 23116936 | 22273434 | 55.6 | 21272850 | 21151065 | 20156692 | 18073071 |
| 14_BK4 | 17849685 | 17383629 | 58.8 | 17310864 | 17235014 | 16752835 | 15595926 |
| 15_BK4 | 16062102 | 15607381 | 58.2 | 15539859 | 15460941 | 14915585 | 13881414 |
| 17_BK4 | 14980797 | 14496468 | 58.1 | 14426404 | 14372805 | 14078235 | 13247545 |
| 18_BK4 | 24217412 | 23575201 | 59.1 | 23370869 | 23281268 | 22704123 | 21306454 |
| 21_BK4 | 11890953 | 11500254 | 60.1 | 11271958 | 11237968 | 11021676 | 10439448 |
| 22_BK4 | 17996717 | 17498339 | 58.8 | 17417850 | 17365274 | 17013067 | 16007094 |
| 25_BK4 | 17560698 | 16997518 | 57.7 | 16888515 | 16816770 | 16375850 | 15237575 |
| 29_BK4 | 8994172 | 8724285 | 58.1 | 8683928 | 8642230 | 8393680 | 7800006 |
| 31_BK4 | 20427813 | 19941632 | 58.4 | 19741741 | 19684774 | 19309226 | 18187574 |
| 32_BK4 | 35100769 | 33926405 | 54.9 | 33754943 | 33623260 | 32780233 | 30194531 |
| 33_BK4 | 24501712 | 23719299 | 58.3 | 21669095 | 21595959 | 21031538 | 19569420 |
| 34_BK4 | 16453473 | 16047224 | 57.3 | 14901123 | 14842998 | 14421402 | 13376818 |
| 47_BK4 | 18736966 | 18155651 | 55.6 | 17998648 | 17903991 | 17174561 | 15478180 |
| 55_BK4 | 17435264 | 16904284 | 48.0 | 16768595 | 16530082 | 14886541 | 12170210 |
| 65_BK3 | 17465925 | 16921732 | 50.6 | 16735483 | 16587671 | 15466034 | 13185810 |
| 71_BK4 | 17919758 | 17434181 | 50.4 | 17086979 | 17017135 | 16549441 | 15441174 |
| 72_BK4 | 16355009 | 15952974 | 57.9 | 15874302 | 15812384 | 15444022 | 14541576 |
| 73_BK4 | 17050731 | 16578547 | 57.8 | 16270896 | 16212738 | 15778509 | 14632126 |
| 77_BK4 | 14044420 | 13478859 | 56.0 | 13390126 | 13322735 | 12866845 | 11763625 |
| 78_BK4 | 17004599 | 16352717 | 60.1 | 16250859 | 16164585 | 15758397 | 15027226 |

544

545

546

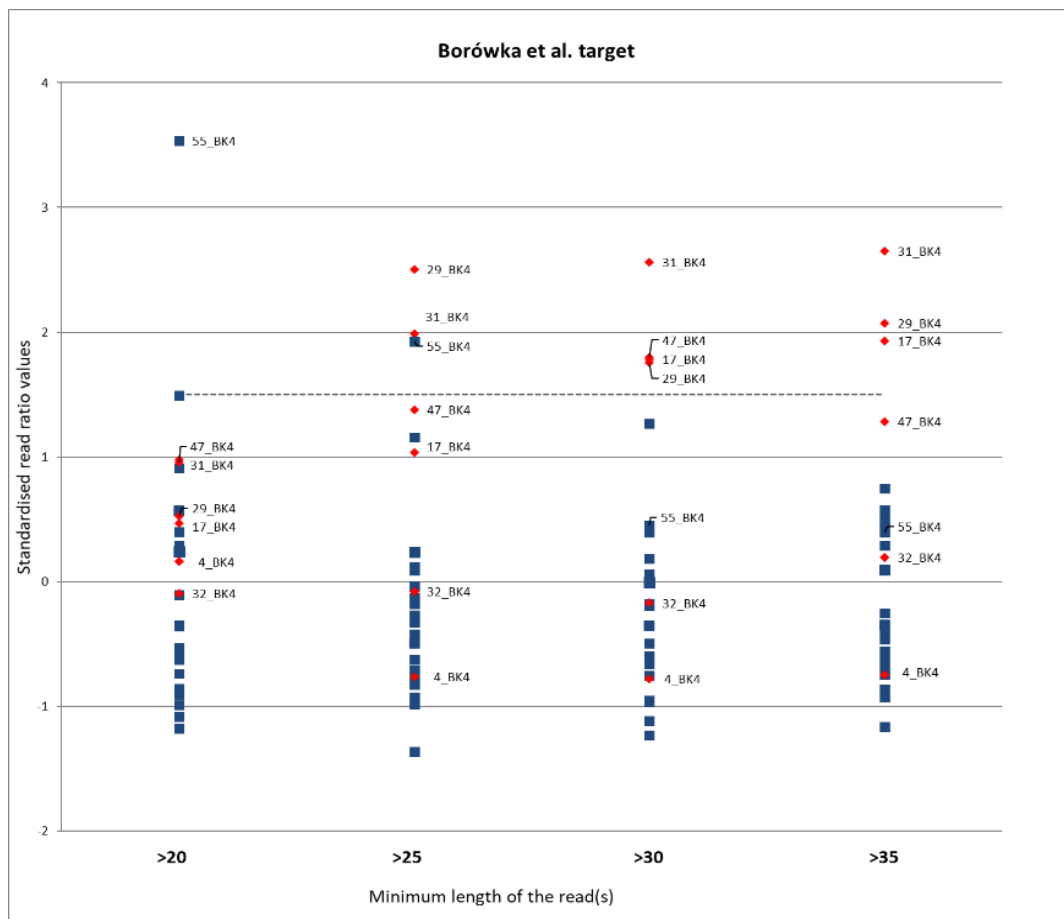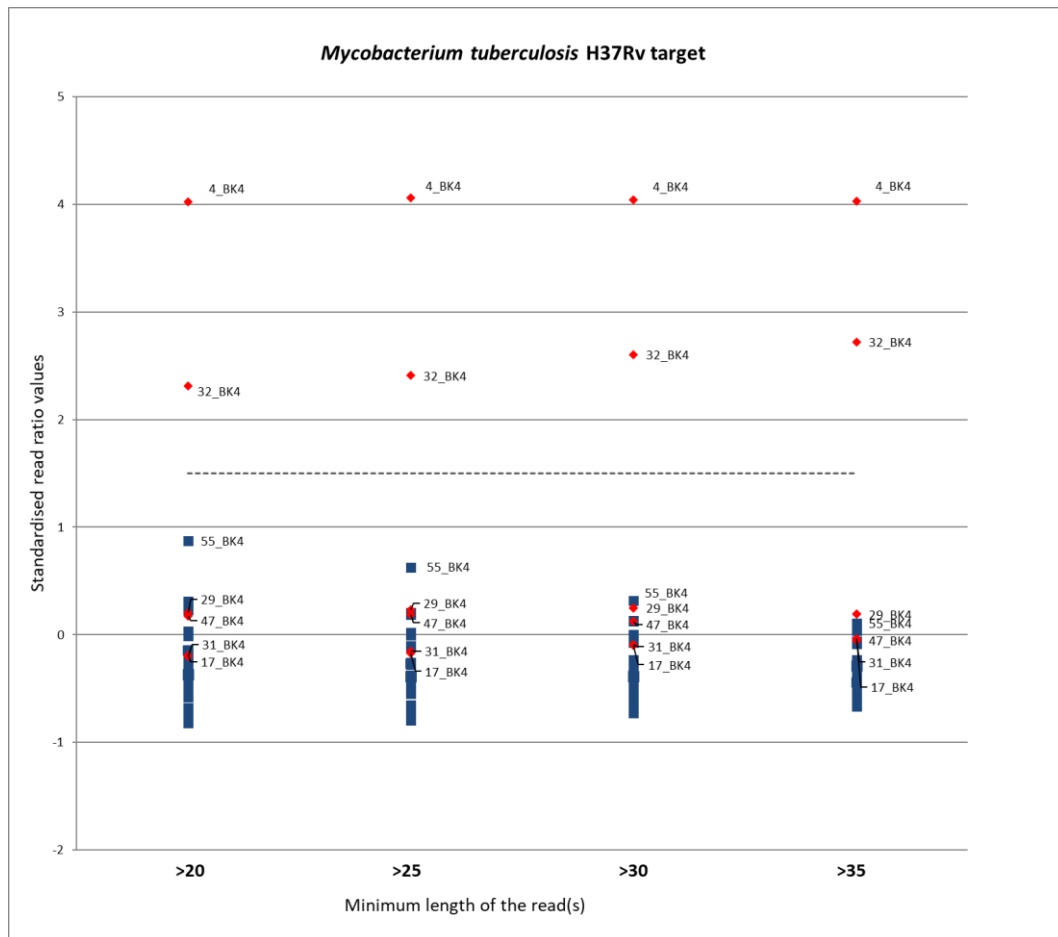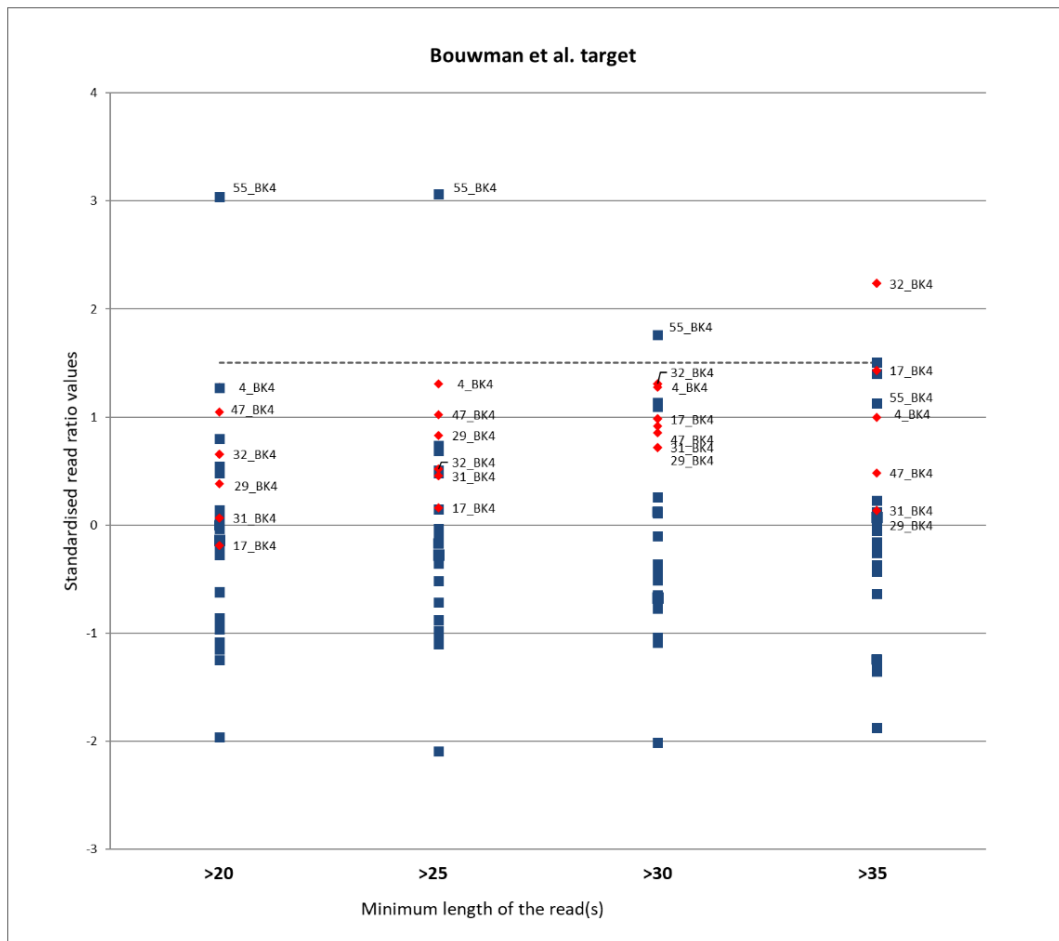547

548        **Fig 1. Changes in standardised ratio values in different read length bins and targets (red diamonds - outliers in** *Mycobacterium*

549            *tuberculosis* **H37Rv and Borówka et al. targets in bin of reads equal or longer than 30).**

550

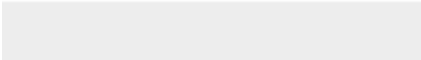551        **Fig 2. Comparison of alignment targets constructed with different assumptions (red diamonds indicate outliers in** *Mycobacterium*

552            *tuberculosis* **H37Rv and Borówka et al. targets in bin of reads equal or longer than 35).**

553

Fig 1

*Mycobacterium tuberculosis* H37Rv target



Borówka et al. target

**Fig 1.** Changes in standardized ratio values in different read length bins and targets (red diamonds – outliers in *Mycobacterium tuberculosis* H37Rv and Borówka et al. targets in bin of reads equal or longer than 30). Dotted line present cutoff values based on 1.5×SD.
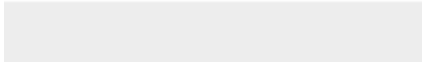
Fig 2

Click here to access/download;Figure;Fig 2.docx ±

**Fig 2.** Comparison of alignment targets constructed with different assumptions (red diamonds indicate outliers in

*Mycobacterium tuberculosis* H37Rv and Borówka et al. targets in bin of reads equal or longer than 35).
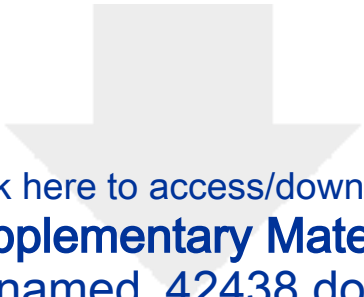
Supplementary figures

Click here to access/download
**Supplementary Material**
renamed_c76da.docx

Supplementary tables

Click here to access/download
**Supplementary Material**
renamed_9c49f.xlsx

Supplementary Tables Legend

Click here to access/download
**Supplementary Material**
renamed_42438.docx

Lodz, 7th of April, 2019

Dear GigaScience Editor,

We appreciate the insightful and detailed comments of the reviewers. We included all the minor language-related corrections in the resubmitted manuscript. Following is a point-by-point reply to the major points raised by the reviewers.

Reviewer 1

1. Our statistical approach, as we state directly in the manuscript, is aimed at finding positive outliers in a pool of samples with unknown presence of mycobacterial sequences. When applied to a dataset like Kay et al., which consists exclusively of individuals with previously confirmed ancient mycobacterial infection, its outcome is therefore necessarily limited to identifying the outliers with highest microbial load (in the case of tuberculosis, potentially those individuals who died during the active phase of the disease) - these outliers being by definition always a minority of analysed samples. This explanation is provided in the text of the manuscript.

2. We have performed the MapDamage analysis suggested by the reviewer and it indeed yielded a positive result - we thank the reviewer for this suggestion as this strengthened our conclusions significantly. We have now included a new Supplementary Fig. 4, and we have reworded both the legend to Supplementary Fig. 3 and the sentences in the manuscript that refer to it.

3. Libraries were build using Meyer et al. (2010) protocol with modifications proposed by Gamba et al. (2014). We have performed mapDamage analysis in the way which fit to double stranded libraries. Information about different Meyer protocol and single stranded libraries was incorrectly added to previous version of the manuscript.

4. We have expanded Supplementary Tab. 2 to include the absolute numbers of reads - both total and aligning to each alignment target.

Reviewer 2

1. We have now expanded the analysis of state of the art in biochemical detection of ancient mycobacteria by citing recent articles mentioning improvements in cell wall component analysis.

2. Indeed, Pott's disease is usually regarded as pathognomonic signature of TB. When we said that many pathological conditions of the spine can mimic Pott's disease thought that they can be diagnosed mistakenly as tuberculosis, especially in practice with poorly preserved skeletons. The present text has been appropriately modified (both in the Introduction and Discussion) to clarify this statement. Moreover, according to the Reviewer's suggestion, the list of pathological conditions has been replaced with the Table S1, which include a short description of basic differences between these lesions and bone tuberculosis.

Dear Editor, regarding to your comment:

"As your revised manuscript focuses more on a method, it may be more suitable as a "Technical Note" rather than a "research article" "

We would like to proceed this manuscript as research paper. In our work we present original dataset which allow us to present a novel bioinformatical approach, used for screening of ancient tuberculosis in sequencing data, derived from 28 individuals (dated 4400 - 4000 BC and 3100 - 2900 BC) from Central Poland. That dataset was not previously published elsewhere.


Sincerely,

Dominik Strapagiel