1  # Comparison of target enrichment strategies for ancient

2  # pathogen DNA

3  Anja Furtwängler[1,2*], Judith Neukamm[1,3,4], Lisa Böhme[5], Ella Reiter[1], Melanie Vollstedt[5],

4  Natasha Arora[6], Pushpendra Singh[7], Stewart T. Cole[8], Sascha Knauf[9,10], Sébastien Calvignac-

5  Spencer[11], Ben Krause-Kyora[5,12], Johannes Krause[1,2,12], Verena J. Schuenemann[1,2,3#,*],

6  Alexander Herbig[1,12#,*]

7

8  [1]Institute for Archaeological Sciences, Archaeo- and Palaeogenetics, University of

9  Tübingen, Germany

10  [2]Senckenberg Centre for Human Evolution and Palaeoenvironment, University of Tübingen,

11  Germany

12  [3]Institute of Evolutionary Medicine, University of Zurich, Switzerland.

13  [4]Institute for Bioinformatics and Medical Informatics, University of Tübingen, Germany

14  [5]Institute of Clinical Molecular Biology, Kiel University, Germany

15  [6]Zurich Institute of Forensic Medicine, University of Zurich, Switzerland

16  [7]Indian Council of Medical Research-National Institute of Research in Tribal Health,

17  Jabalpur, MP, India

18  [8]Institute Pasteur, Paris, France

19  [9]Deutsches Primatenzentrum GmbH, Leibniz-Institute for Primate Research, Goettingen,

20  Germany

21  [10]Department for Animal Sciences, Georg-August-University, Goettingen, Germany

22  [11]Robert Koch Institut, Berlin, Germany

23  [12]Max Planck Institute for the Science of Human History, Jena, Germany

24    #These authors jointly supervised this study.

25    * Corresponding authors: anja.furtwaengler@uni-tuebingen.de, verena.schuenemann@iem.uzh.ch,
26    herbig@shh.mpg.de

27

# Abstract

29    In ancient DNA research, the degraded nature of the samples generally results in poor yields

30    of highly fragmented DNA, and targeted DNA enrichment is thus required to maximize

31    research outcomes. The three commonly used methods – (1) array-based hybridization capture

32    and in-solution capture using either (2) RNA or (3) DNA baits – have different characteristics

33    that may influence the capture efficiency, specificity, and reproducibility. Here, we compared

34    their performance in enriching pathogen DNA of *Mycobacterium leprae* and *Treponema*

35    *pallidum* of 11 ancient and 19 modern samples. We find that in-solution approaches are the

36    most effective method in ancient and modern samples of both pathogens, and RNA baits usually

37    perform better than DNA baits.

# Method summary

38

39    We compared three targeted DNA enrichment strategies used in ancient DNA research for

40    the specific enrichment of pathogen DNA regarding their efficiency, specificity, and

41    reproducibility for ancient and modern *Mycobacterium leprae* and *Treponema pallidum*

42    samples. Array-based capture and in-solution capture with RNA and DNA baits were all tested

43    in three independent replicates.

# Main Text

44

45    The field of ancient DNA (aDNA), which studies DNA retrieved from paleontological and

46    archaeological material, was revolutionized by the invention of high-throughput sequencing

47    (HTS). In combination with HTS, the development of targeted DNA enrichment protocols has

48    made a crucial contribution in advancing aDNA research during the last decade.

49      As DNA decays over time, aDNA is usually only present in trace amounts of highly

50      fragmented sequences (**1, 2, 3**). Detecting endogenous pathogen aDNA from archaeological

51      material is additionally compounded by the larger amount of background DNA from the

52      environment including soil microorganisms. Furthermore, the background of host DNA in

53      ancient remains is an additional obstacle in order to obtain ancient pathogen DNA. Shotgun

54      sequencing of libraries from aDNA extracts to sufficient genomic coverage is, therefore, cost-

55      intensive (**4**). To circumvent this problem, specific regions of interest such as bacterial

56      chromosomes, mammalian mitochondrial genomes, or regions with single-nucleotide-

57      polymorphisms (SNP) are often target-enriched before sequencing (**4**). Aside from its

58      application in aDNA sequencing, targeted DNA enrichment is also useful to retrieve pathogen

59      DNA from clinical samples, particularly for infectious agents that are found in low quantities

60      in the host organism and which are difficult to culture, as is the case for *Mycobacterium leprae*

61      and *Treponema pallidum.* Removal of background DNA prior to sequencing increases the yield

62      of pathogen DNA, and thus allows valuable information for epidemiologists investigating

63      outbreaks to be obtained.

64      For the enrichment of entire bacterial and mammalian chromosomes, there are currently

65      three methods available, which are based on hybridization capture (**5**): DNA microarrays (here

66      represented by SureSelect from Agilent Technologies), in-solution capture with DNA baits

67      (represented by SureSelect from Agilent Technologies according to Fu and colleagues (**6**)) and

68      in-solution capture with RNA baits (here represented by myBaits® from Arbor Biosciences).

69      In the case of the DNA array-based method, up to a million artificial DNA baits are printed

70      on the surface of a glass slide (**7**). Additionally, there is the possibility to perform in-solution

71      capture with baits cleaved from the glass slides and used right away or immortalized in DNA

72      bait libraries (**6**). The second in-solution approach uses up to 100,000 artificial RNA baits. The

73 three approaches rely on the hybridization of target fragments to the complementary sequence

74 of the baits (immobilized or in-solution), which can be levered to wash background DNA away.

75 To date there has been to our knowledge, no statistical comparison of the performance of all

76 three methods: microarrays, in-solution capture with DNA baits, and in-solution capture with

77 RNA baits (**6**). So far only microarrays and the in-solution capture with DNA baits were

78 compared for *Salmonella enterica* and no replicates for statistical assessment were produced

79 (**8**).

80 Here, we present results from the enrichment of modern and ancient samples containing

81 pathogen DNA, using the three aforementioned approaches. All samples had previously tested

82 positive but had also shown low amounts of target DNA for *M. leprae* or *T. pallidum*

83 (Supplementary Table 1).

84 The different enrichment concepts tested were chosen to represent methods as they are

85 applied in ongoing research and therefore not only differ in the technology used (DNA *vs.* RNA

86 baits, immobilized *vs.* in-solution) but also in the design such as bait length and number of

87 unique baits, which might have an effect on the performance.

88 We used eight ancient samples positive for *M. leprae* and six modern libraries from leprosy

89 patients that were shown to contain *M. leprae* DNA (Supplementary Note 1). Genetic data from

90 the ancient and modern *M. leprae* samples were previously published in **9** and **10**. Samples with

91 less than 0.6 % endogenous bacterial DNA were selected.

92 Modern *T. pallidum* samples (n=13) were previously published in **12** and **13**. Three ancient

93 extracts of *T. pallidum* were used from **14**. The portion of endogenous DNA for the selected

94 *T. pallidum* samples was below 0,01 % for ancient and modern samples.

95 Starting from existing sequencing libraries all three methods were applied with three

96 independent replicates each (see Figure. 1 and Supplementary Note 1 for a detailed description

97 of the methods, the newly generated data is available at the Sequence Read Archive under the

4

98      BioProject PRJNA645054). Following the manufacturer's suggestion for libraries with low

99      yields of target DNA, we performed two successive rounds of hybridization for all methods. To

100     investigate the effectiveness of this procedure, we compared results from the first and second

101     rounds for the in-solution capture with RNA baits. We then evaluated differences in efficiency,

102     reproducibility, and specificity across the three approaches by calculating mean coverage,

103     standard deviation of the mean coverage, enrichment factor (calculated by dividing the % of

104     target DNA after enrichment by the % of target DNA in the shotgun data), and the % of the

105     genome covered 5-fold or more after normalizing the data of each bacterial species to the same

106     number of raw reads (Supplementary Tables 2, 3 & 5 and Supplementary Figures 1 & 2).

107     For most ancient samples, the highest mean coverage (Figure 2A) is reached with the RNA

108     bait in-solution capture (eight out of eleven, more details can be found in Supplementary Note

109     2 & 3, and SSupplementary Tables 1 & 2). On average the RNA bait capture results in a 1.5

110     and 20.0 times higher mean coverage than the DNA bait or the array capture, respectively. As

111     illustrated in Figure. 2B, the highest enrichment factor is obtained in the RNA bait capture of

112     ancient *T. pallidum* DNA (all three samples) and *M. leprae* (four samples showed best results

113     for the RNA bait, three for the DNA bait, and one for the array), with values between 2-150x

114     higher, compared to the other two approaches. An in-solution approach seems, therefore, to be

115     advantageous for enriching ancient pathogen DNA.

116     A similar pattern can be observed in the data of the modern *M. leprae* and *T. pallidum*

117     samples (Figures. 2A and 2B) further highlighting the performance of the in-solution approach

118     in general and RNA baits in particular.

119     In-solution capture with DNA baits was used with robot-assistance in this study whereas the

120     in-solution capture with RNA baits was performed in two different labs. Unsurprisingly, the

121     DNA bait capture showed the smallest differences (2- to 50-fold lower) between the replicates

5

122 whereas the RNA bait capture showed the largest and the DNA array capture was intermediate.

123 Consistent conditions are therefore crucial for reproducibility.

124 Another important feature of targeted enrichment is specificity. We estimated the specificity of

125 the three tested methods by comparing the number of reads specific to either *M. leprae* or *T.*

126 *pallidum* in comparison to general mycobacterial or treponemal reads, respectively (Figure 2

127 C). Here, differences between the two pathogens can be observed. In the ancient and modern

128 *T. pallidum* samples, the RNA bait capture consistently shows the highest proportion (up to 1.5

129 times higher) of specific reads. The same trend was observed for the libraries prepared from

130 recent leprosy patient samples, i.e. modern samples of *M. leprae*. Only for ancient *M. leprae*

131 samples, the DNA bait capture is more specific. The highest percentages of specific reads are

132 not necessarily found in samples with high percentages of endogenous DNA in the shotgun data

133 before enrichment.

134 For ancient and modern samples, due to high efficiency, reproducibility and specificity in-

135 solution approaches are highly recommendable.

136 Two rounds of hybridization are routinely performed in aDNA research, which is expected to

137 improve enrichment but may also reduce data complexity in terms of portions of unique reads.

138 To formally investigate the effect of the second round of capture, we also sequenced the

139 libraries only enriched with one round of hybridization with the RNA baits and compared the

140 results to the second round of hybridization. The second round of hybridization resulted in an

141 increase in the enrichment factor for ancient and modern *M. leprae* samples (with an average

142 of 2x increase) as well as for *T. pallidum* samples (with an average of 17x increase),

143 demonstrating the utility of such a second round of hybridization capture (Supplementary Table

144 5). On the other hand, when comparing the library complexity (Figure. 2 D and Supplementary

145 Note 2 & 3, Supplementary Figure 3), we found a substantial loss of complexity after the second

146 round of hybridization in all modern and ancient samples. This loss was reflected in the higher

147   percentage of unique reads in all the reads mapped after the first round. Therefore, if the portion

148   of endogenous DNA in a sample is high in the beginning it may be worthwhile considering

149   whether a single round of capture combined with deeper sequencing is sufficient or even

150   advantageous.

151   The three protocols also differ in terms of cost and effort. The most cost-intensive is the array-

152   capture approach (~673 € per sample), which requires additional equipment that is not usually

153   necessary with the other approaches. The in-solution capture with DNA baits is, by contrast,

154   cheaper once the baits are cleaved from the glass slide (~56,23€ per sample), but the version

155   that can be used for the immortalization of the baits by transforming them into a library is not

156   freely available. The in-solution capture with RNA baits is more comparable to the DNA bait

157   capture than to the array with ~109 € per sample and it also needs the lowest number of

158   additional equipment and reagents (Supplementary Table 7).

159   After a detailed comparison of the three tested methods it can be concluded that for ancient

160   and modern pathogen samples, the RNA bait capture with two rounds of hybridization seems

161   to be the most suitable. The generally high performance of the in-solution approach (mainly the

162   one with RNA baits) for both bacterial species suggests that the findings are highly

163   representative and comparable performance is also expected for a variety of other

164   bacterial/microbial organisms.

165   # References (max. 20 References)

166   1. Sawyer, Susanna; Krause, Johannes; Guschanski, Katerina; Savolainen, Vincent; Pääbo, Svante (2012): Temporal
167      patterns of nucleotide misincorporations and DNA fragmentation in ancient DNA. In: *PloS one* 7 (3), e34131. DOI:
168      10.1371/journal.pone.0034131.

169   2. Allentoft, Morten E.; Collins, Matthew; Harker, David; Haile, James; Oskam, Charlotte L.; Hale, Marie L. et al.
170      (2012): The half-life of DNA in bone: measuring decay kinetics in 158 dated fossils. In: *Proceedings. Biological*
171      *sciences* 279 (1748), S. 4724–4733. DOI: 10.1098/rspb.2012.1745.

172   3. Briggs, Adrian W.; Stenzel, Udo; Johnson, Philip L. F.; Green, Richard E.; Kelso, Janet; Prüfer, Kay et al. (2007):
173      Patterns of damage in genomic DNA sequences from a Neandertal. In: *Proceedings of the National Academy of*
174      *Sciences of the United States of America* 104 (37), S. 14616–14621. DOI: 10.1073/pnas.0704665104.

175   4. Krause, Johannes (2010): From Genes to Genomes: What is New in Ancient DNA? In: *Mitteilungen der*
176      *Gesellschaft für Urgeschichte* 19, S. 11–33.

177   5.  Spyrou, Maria A.; Bos, Kirsten I.; Herbig, Alexander; Krause, Johannes (2019): Ancient pathogen genomics as an
178       emerging tool for infectious disease research. In: *Nature reviews. Genetics* 20 (6), S. 323–340. DOI:
179       10.1038/s41576-019-0119-1.

180   6.  Fu, Qiaomei; Meyer, Matthias; Gao, Xing; Stenzel, Udo; Burbano, Hernán A.; Kelso, Janet; Pääbo, Svante (2013):
181       DNA analysis of an early modern human from Tianyuan Cave, China. In: *Proceedings of the National Academy of*
182       *Sciences of the United States of America* 110 (6), S. 2223–2227. DOI: 10.1073/pnas.1221359110.

183   7.  Vågene, Åshild J.; Herbig, Alexander; Campana, Michael G.; Robles García, Nelly M.; Warinner, Christina; Sabin,
184       Susanna et al. (2018): Salmonella enterica genomes from victims of a major sixteenth-century epidemic in Mexico.
185       In: *Nature ecology & evolution* 2 (3), S. 520–528. DOI: 10.1038/s41559-017-0446-6.

186   8.  Burbano, Hernán A.; Hodges, Emily; Green, Richard E.; Briggs, Adrian W.; Krause, Johannes; Meyer, Matthias et
187       al. (2010): Targeted investigation of the Neandertal genome by array-based sequence capture. In: *Science (New*
188       *York, N.Y.)* 328 (5979), S. 723–725. DOI: 10.1126/science.1188046.

189   9.  Schuenemann, Verena J.; Singh, Pushpendra; Mendum, Thomas A.; Krause-Kyora, Ben; Jäger, Günter; Bos,
190       Kirsten I. et al. (2013): Genome-wide comparison of medieval and modern Mycobacterium leprae. In: *Science*
191       *(New York, N.Y.)* 341 (6142), S. 179–183. DOI: 10.1126/science.1238286.

192   10. Schuenemann, Verena J.; Avanzi, Charlotte; Krause-Kyora, Ben; Seitz, Alexander; Herbig, Alexander; Inskip,
193       Sarah et al. (2018): Ancient genomes reveal a high diversity of Mycobacterium leprae in medieval Europe. In:
194       *PLoS pathogens* 14 (5), e1006997. DOI: 10.1371/journal.ppat.1006997.

195   11. Knauf, Sascha; Gogarten, Jan F.; Schuenemann, Verena J.; Nys, Hélène M. de; Düx, Ariane; Strouhal, Michal et al.
196       (2018): Nonhuman primates across sub-Saharan Africa are infected with the yaws bacterium Treponema pallidum
197       subsp. pertenue. In: *Emerging microbes & infections* 7 (1), S. 157. DOI: 10.1038/s41426-018-0156-4.

198   12. Arora, Natasha; Schuenemann, Verena J.; Jäger, Günter; Peltzer, Alexander; Seitz, Alexander; Herbig, Alexander
199       et al.: Origin of modern syphilis and emergence of a pandemic Treponema pallidum cluster. In: *Nat Microbiol* 2
200       (1), S. 1–6. DOI: 10.1038/nmicrobiol.2016.245.

201   13. Schuenemann, Verena J.; Kumar Lankapalli, Aditya; Barquera, Rodrigo; Nelson, Elizabeth A.; Iraíz Hernández,
202       Diana; Acuña Alonzo, Víctor et al. (2018): Historic Treponema pallidum genomes from Colonial Mexico retrieved
203       from archaeological remains. In: *PLoS neglected tropical diseases* 12 (6), e0006447. DOI:
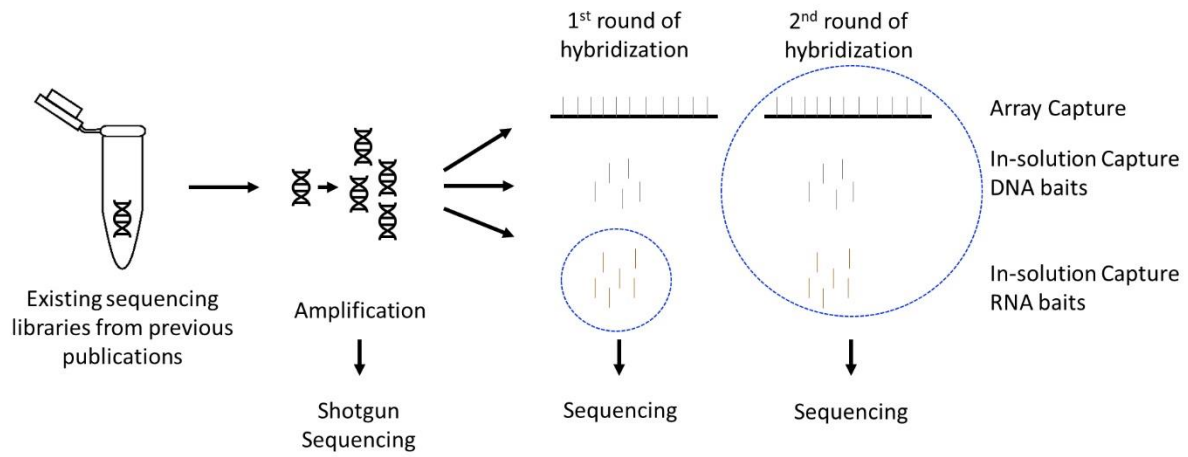204       10.1371/journal.pntd.0006447.

205

# Author contributions

V.J.S., A.H. and J.K. conceived of the study. B.K. and S.C-S. provided RNA baits and sequencing libraries. N.A., P.S., S.T.C., S.K. provided sequencing libraries. A.F., L.B., E.R., M.V. performed the laboratory work. A.F. and J.N. performed the data analysis. A.F. and A.H. conducted the statistical analysis. A.F. designed the figures. A.F., V.J.S, and A.H. wrote the manuscript with input from all authors. All authors reviewed the manuscript.
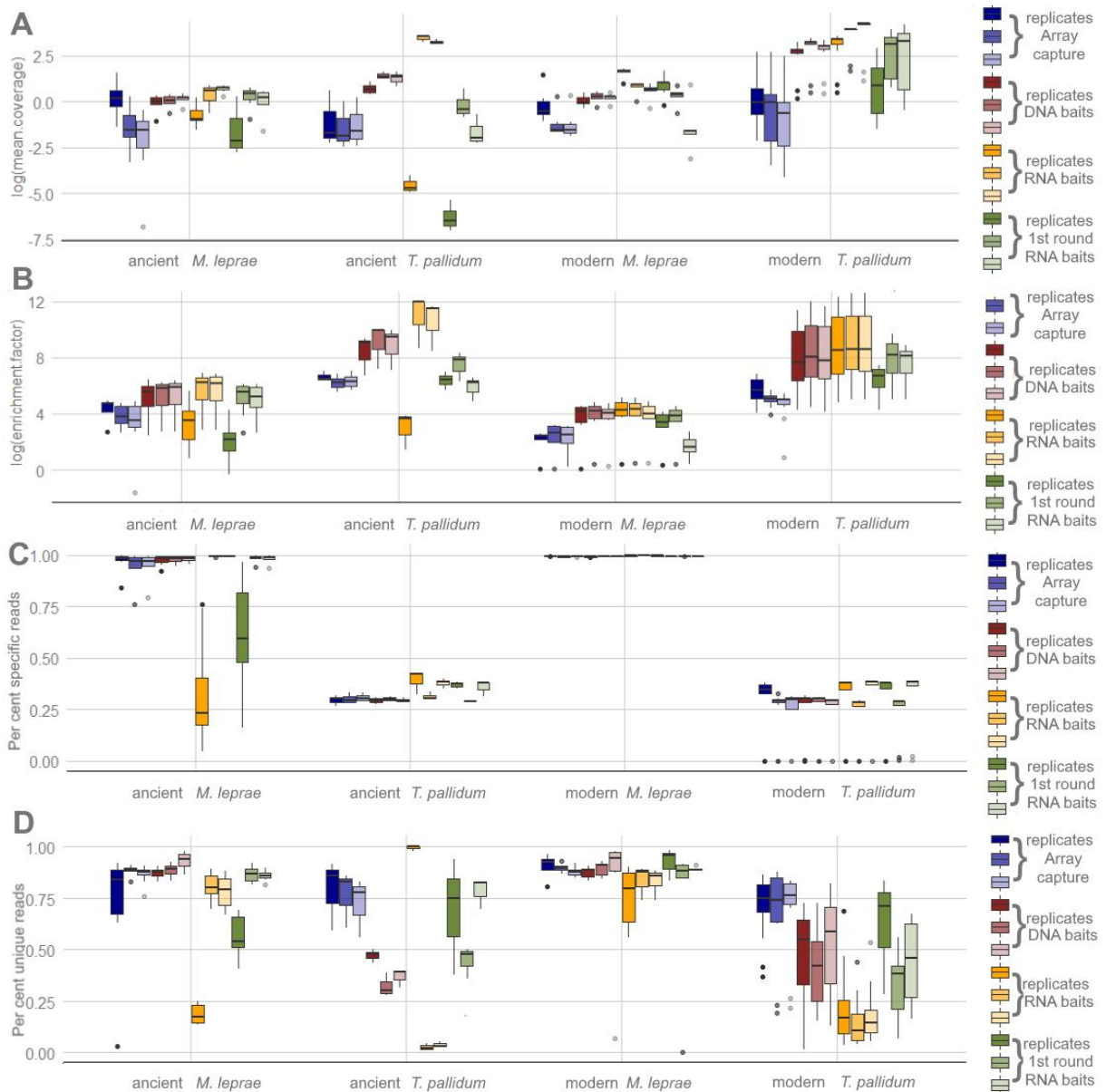
# Acknowledgments

8

240

241      **Figure 1. Schematic representation of the workflow.** For all samples, the three different

242      enrichment protocols were tested in three independent replicates. Blue circles indicate the

243      libraries that were sequenced at each particular step.

**Figure 2. Differences between the three tested protocols in ancient and modern *M. leprae* and *T. pallidum* samples.** A) Log-transformed values of the mean coverage. B) log-transformed values of the enrichment factor calculated by dividing the percentage of endogenous DNA by the percentage of endogenous DNA after shotgun sequencing. C) The proportion of specific reads corresponding to *M. leprae* and *T. pallidum* compared to other mycobacterial and treponemal reads, respectively. D) Percentage of unique reads calculated by the number of unique reads divided by the total number of sequences mapped to represent library complexity in *M. leprae* and *T. pallidum* samples.