

First submission

Please read the **Important notes** below, and the **Review guidance** on the next page.
When ready [submit online](#). The manuscript starts on page 3.

Important notes

Editor and deadline

Thomas Rattei / 12 Oct 2016

Files

7 Figure file(s)

6 Latex file(s)

1 Table file(s)

1 Other file(s)

Please visit the overview page to [download and review](#) the files not included in this review pdf.

Declarations

One or more DNA sequences were reported.




Please in full read before you begin

How to review






When ready [submit your review online](#). The review form is divided into 5 sections. Please consider these when composing your review:

- 1. BASIC REPORTING**
- 2. EXPERIMENTAL DESIGN**
- 3. VALIDITY OF THE FINDINGS**
4. General comments
5. Confidential notes to the editor



 You can also annotate this **pdf** and upload it as part of your review

To finish, enter your editorial recommendation (accept, revise or reject) and submit.







BASIC REPORTING

-  Clear, unambiguous, professional English language used throughout.
-  Intro & background to show context. Literature well referenced & relevant.
-  Structure conforms to [PeerJ standard](#), discipline norm, or improved for clarity.
-  Figures are relevant, high quality, well labelled & described.
-  Raw data supplied (See [PeerJ policy](#)).

VALIDITY OF THE FINDINGS

-  Impact and novelty not assessed. Negative/inconclusive results accepted. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.
-  Data is robust, statistically sound, & controlled.

EXPERIMENTAL DESIGN

-  Original primary research within [Scope of the journal](#).
-  Research question well defined, relevant & meaningful. It is stated how research fills an identified knowledge gap.
-  Rigorous investigation performed to a high technical & ethical standard.
-  Methods described with sufficient detail & information to replicate.
-  Conclusion well stated, linked to original research question & limited to supporting results.
-  Speculation is welcome, but should be identified as such.

The above is the editorial criteria summary. To view in full visit <https://peerj.com/about/editorial-criteria/>

Improving ancient DNA genome assembly

Alexander Seitz ^{Corresp.} ¹, Kay Nieselt ¹

¹ Center for Bioinformatics (ZBIT), Integrative Transcriptomics, Eberhard-Karls-Universität Tübingen, Tübingen, Germany

Corresponding Author: Alexander Seitz
Email address: alexander.seitz@uni-tuebingen.de

Most reconstruction methods for genomes of ancient origin that are used today require a closely related reference. In order to identify genomic rearrangements or the deletion of whole genes, de novo assembly has to be used. However, because of inherent problems with ancient DNA, its de novo assembly is highly complicated. In order to tackle the diversity in the length of the input reads, we propose a two-layer approach, where multiple assemblies are generated in the first layer, which are then combined in the second layer. We used this two-layer assembly to generate assemblies for an ancient sample and compared the results to current de novo assembly approaches. We are able to improve the assembly with respect to the length of the contigs and can resolve more repetitive regions.

1 Improving ancient DNA genome assembly

2 Alexander Seitz¹ and Kay Nieselt¹

3 ¹Center for Bioinformatics (ZBIT), Integrative Transcriptomics,
4 Eberhard-Karls-Universität Tübingen

5 ABSTRACT

6 Most reconstruction methods for genomes of ancient origin that are used today require a closely related reference. In order to identify genomic rearrangements or the deletion of whole genes, *de novo* assembly has to be used. However, because of inherent problems with ancient DNA, its *de novo* assembly is highly complicated. In order to tackle the diversity in the length of the input reads, we propose a two-layer approach, where multiple assemblies are generated in the first layer, which are then combined in the second layer. We used this two-layer assembly to generate assemblies for an ancient sample and compared the results to current *de novo* assembly approaches. We are able to improve the assembly with respect to the length of the contigs and can resolve more repetitive regions.

7 Keywords: ancient DNA, *de novo* assembly, genome reconstruction

8 INTRODUCTION

9 The introduction of next generation sequencing (NGS) made large scale sequencing projects feasible (Bentley et al., 2008). Their high throughput allows fast and cheap sequencing of arbitrary genomic material. It revolutionized modern sequencing projects and made the study of ancient genomes possible (Der Sarkissian et al., 2015). However, the resulting short reads pose several challenges for the reconstruction of the desired genome when compared to the longer Sanger reads (Li et al., 2010). For modern DNA samples, the problem of having only short reads can be mitigated by the sheer volume of sequenced bases and usage of long fragments with paired-end and mate-pair sequencing. The insert size is used to determine the distance between the forward and the reverse read, which are sequenced from both ends of the fragments. These distances can be important for *de novo* assembly as they are used for repeat resolution and scaffolding. However, samples from ancient DNA (aDNA) mostly contain only very short fragments between 44 and 172 bp (Sawyer et al., 2012). Paired-end sequencing of these short fragments therefore often results in overlapping forward and reverse reads (thus actually negative inner mate pair distances). This has two consequences: the usage of mate-pairs as well as sequencing technologies producing long reads is not beneficial. Additionally, post-mortem damage of aDNA, most importantly the deamination of cytosine to uracil, can result in erroneous base incorporation (Rasmussen et al., 2010). Using reference based approaches, these errors can be detected, as they always occur at the end of the fragments. This is not possible using *de novo* assembly approaches and these errors can lead to mistakes in the assembly. Deeper sequencing does not yield better results as the amount of endogenous DNA contained in aDNA samples is often very low (Sawyer et al., 2012).

10 In order to achieve a higher content of endogenous DNA, samples are often subject to enrichment using capture methods (Avila-Arcos et al., 2011). The principle of these capture methods relies on selection by hybridization (Maricic et al., 2010). Regions of interest are fixed to probes prior to sequencing. These probes can be immobilized on glass slides, called array capture (Hodges et al., 2007), or recovered by affinity using magnetic beads, referred to as in-solution capture (Gnirke et al., 2009). Using these capture methods, only DNA fragments that can bind to the probes are used for amplification, which increases the amount of the desired DNA. However, as these methods only amplify sequences that are contained in the probes, regions that were present in ancient samples and lost over time cannot be amplified and thus not be identified. Nevertheless, most of the current aDNA projects use these capture methods.

11 Currently, there are two ways to reconstruct a genome from sequencing data. If there is a known, closely related genome, it can be used as a reference. Mapping programs like BWA (Li and Durbin, 2009) can then be used to align the reads against the reference genome. Single nucleotide variations (SNVs) or short indels between the DNA sequence of the sample and reference can be identified after all reads are aligned.

44 Because of the inherent characteristics of aDNA, specialized mapping pipelines for the recon-
45 struction of aDNA genomes, such as EAGER (Peltzer et al., 2016) and PALEOMIX (Schubert et al.,
46 2014), have recently been published. The mapping against a reference genome allows researchers to
47 easily eliminate non-endogenous DNA and identify erroneous base incorporations. These errors can
48 be identified after the mapping and used to verify that the sequenced fragments stem from ancient
49 specimen.

50 The reference-based mapping approaches cannot detect large insertions or other genomic archi-
51 tectural rearrangements. In addition, if the ancient species contained regions that are no longer
52 present in the modern reference, these cannot be identified via mapping against modern reference
53 genomes. In these cases a *de novo* assembly of the genome should be attempted. This is also true for
54 modern samples, if no closely related reference is available. If the ancient sample was sequenced
55 after amplification through capture arrays, genomic regions that are contained on the probe
56 also can't be identified. Using shotgun sequencing, sequences that stem from species that migrated
57 into the sample post-mortem are often more abundant (Knapp and Hofreiter, 2010). However, if
58 shotgun data is available an effort for assembly can be made to identify longer deletions or genomic
59 rearrangements. The introduction of NGS has led to new assembly programs that can handle short
60 reads such as SOAPdenovo2 (Luo et al., 2012), SPARK (Bankevich et al., 2012) and many more.

61 The assembly of modern NGS data is still a hard problem (Chao et al., 2015) and methods to
62 improve them are constantly developed. Among these is ALLPATHS-LG (Gnerre et al., 2011),
63 arguably the winner of the so-called Assemblathon (Earl et al., 2011). ALLPATHS-LG uses the
64 information provided by long fragments from paired-end and mate-pair sequencing to improve the
65 assembly, and has therefore been shown to be one of the best assembly programs that are available
66 today (Utturkar et al., 2014). However, because of the short fragments contained in aDNA samples,
67 this approach is not feasible for aDNA samples and other methods have to be employed.

68 *De Bruijn* graph assemblers highly rely on the length of the *k*-mer to generate the graph (Li
69 et al., 2012). The choice of an optimal value is already a hard problem for modern sequencing
70 projects (Durai and Schulz, 2016).

71 Because of the short fragments of aDNA samples, the sequencing adapter is often partially or fully
72 sequenced. After the adapter is removed, the length of the resulting read is then equal to the length of
73 the fragment. Furthermore, overlapping forward and reverse reads can be merged to generate longer
74 reads, which is usually done in aDNA studies to improve the sequence quality (Peltzer et al., 2016).
75 Thus the length distribution of reads from aDNA samples is often very skewed. This implies that the
76 choice of one single fixed *k* for the *k*-mer in *de Bruijn* graph-based assembly approaches is not ideal
77 in aDNA studies. Long *k*-mers miss all reads that are shorter than the value of *k* and shorter *k*-mers
78 cannot resolve repetitive regions.

79 We have developed a two-layer assembly approach where in the first layer, the contigs are
80 assembled from short reads using a *de Bruijn* graph approach with multiple *k*-mers. These contigs
81 are then used in the second layer in order to combine overlapping contigs contained in the different
82 assemblies resulting from the first layer. This is done using an overlap-based approach.

83 **Outline** This article is organized as follows. The next section contains the methods we used to
84 improve the *de novo* assembly for aDNA samples. In short, we used multiple assemblies with different
85 *k*-mers and then merge these assemblies into longer contigs. In the results section, we used our
86 two-layer assembly to improve the assembly of the sample Jorgen625 published by Schuenemann
87 et al. (2013). Finally, we conclude our findings and give an outlook.

88 METHODS

89 The general structure of our two-layer approach is as follows: In the first layer, the raw *fastq*
90 files are preprocessed, followed by a *de Bruijn* graph-based assembly using multiple *k*-mer sizes to
91 generate several different, yet similar assemblies. All produced contigs are quality filtered before
92 they are combined and used in the second layer. There, an overlap-based approach is used to identify
93 contigs in the different assemblies that represent the same genomic region. These can be merged
94 into longer contigs. Afterwards small contigs are removed from the result. The rest of this section
95 explains these steps in more detail.

96 We used the tool Clip & Merge (Peltzer et al., 2016) to remove the sequencing adapters. It was
97 also used to quality trim all bases in reads below a minimum phred score of 20. This threshold was left
98 at this default value as the low-quality ends of the reads are merged and thus the base call is confirmed
99 by two reads. The value was not changed for the unmerged reads in order to be able to compare the
100 experiments. In order to evaluate how different preprocessing affects the assembly, the reads were

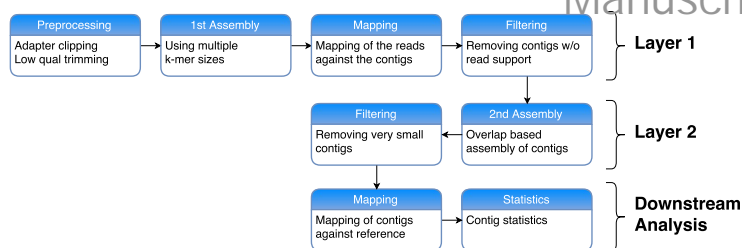


Figure 1. Workflow of our two-layer assembly approach. First the reads are preprocessed by removing sequenced adapters and clipping low-quality bases. After that, multiple *de novo* assemblies are generated using a *de Bruijn* graph approach with multiple values for k . The reads are then mapped back against each of these resulting contigs and the contigs with no read support are filtered out. In Layer 2, these filtered contigs are then combined and assembled again using an Overlap-Layout-Consensus approach. Very short contigs are removed. The resulting contigs are mapped against a reference genome and contig statistics are calculated in order to assess the quality of the assembly.

101 treated using three different methods: First, the reads were only adapter clipped and trimmed. Reads
 102 that no longer have a partner were removed. These reads were then used in a paired-end assembly.
 103 Second, after the reads were adapter clipped and quality trimmed, all resulting forward and reverse
 104 reads were combined into one file, each read given a unique identifier so that they could be used in a
 105 single-end assembly. Third, after the adapter clipping the forward and reverse reads were merged into
 106 longer reads whenever possible. For the merging of the reads, we used the standard parameters of
 107 Clip & Merge defining a minimum overlap length of 10 bp with a maximum mismatch rate of 5%.
 108 The resulting reads were then quality trimmed as described above. Unique identifiers were assigned
 109 to forward and reverse reads that could not be merged and added to the resulting *fastq* file. These
 110 reads were then used in a single-end assembly. In all three sets, resulting reads that were shorter than
 111 25 bp were removed before the assembly.


112 After the preprocessing, the resulting reads are of different lengths. The reason for this are the
 113 different fragment lengths contained in the sample. This is why we propose assembly of aDNA using a
 114 two-layer approach. In the first layer, we use a k -mer based assembly program like SOAPdenovo2 (Luo
 115 et al., 2012), MEGAHIT (Li et al., 2014), or any other assembly program for which different values
 116 for k can be chosen.

117 *De Bruijn* based programs first generate all possible k -mers based on the input reads. Matching
 118 k -mers are used to generate the *de Bruijn* graph. This can lead to random overlaps of k -mers contained
 119 in different reads and therefore to read incoherent contigs (Myers, 2005). To filter out these contigs,
 120 the reads are mapped back against the resulting contigs. This can be done by using modern mapping
 121 programs like B-MEM (Li, 2013). Contigs that are not supported by any read are removed before
 122 the next step.

123 To combine the results of the different assemblies, each contig is given a unique identifier before
 124 they are combined into one file. This file is the input of the second layer assembly. Here, the assembly
 125 is based on string overlaps instead of k -mers, a concept originally introduced by Myers (2005). An
 126 assembly program that uses this approach is the String Graph Assembler (SGA) (Simpson and Durbin,
 127 2012). It efficiently calculates all overlaps of the input using suffix arrays (Manber and Myers, 1993).
 128 These overlaps are then used to generate an overlap graph and the final contigs are generated based
 129 on this graph. We used this method to merge the contigs from the different assemblies based on their
 130 overlap.



131 As SGA uses string-based overlaps and modern sequencing techniques are not error-free, it
 132 provides steps to correct for these errors. There is a preprocessing step that removes all bases that are
 133 not A,G,C or T. There is also a correction step that performs a k -mer based error correction and a
 134 filtering step that removes input reads with a low k -mer frequency. Because the input for SGA are
 135 already pre-assembled contigs, these errors are already averaged out and these steps are not used
 136 for the assembly of the second layer. However, the assemblies with the different k -mers produce
 137 similar contigs, which is why the duplicate removal step of SGA is performed. SGA can also use the
 138 Ferragina Manzini (FM) index (Ferragina and Manzini, 2000) to merge unambiguously overlapping
 139 sequences, which is used to further remove duplicate information. Afterwards the overlap graph
 140 is calculated and the new contigs are assembled. All these steps are performed using the standard
 141 parameters provided by SGA. Afterwards, contigs shorter than 1 000 bp are removed from the final

142 assembly. In order to evaluate our two-layer assembly method, the resulting contigs are then aligned
 143 with the reference genome of interest. We use again BWA-MEM for this step. Finally various
 144 statistics for the assembly are computed.


145 An overview of this methodology can be seen in Figure 1 

146 RESULTS

147 To evaluate our two-layer assembly, we applied it to
 148 a published ancient sample containing DNA from *My-*
 149 *cobacterium leprae*. We used the sample Jorgen625 pub-
 150 lished by Schuenemann et al. (2013). The bones from
 151 which the DNA was extracted, are approximately 700
 152 years old. Two different sequencing libraries are avail-
 153 able for this sample. In order to get an overview of
 154 the two libraries, we used the EAGER pipeline (Peltzer
 155 et al., 2016) to map the two libraries against the refer-
 156 ence genome of *Mycobacterium leprae TN*. One of the
 157 two libraries contained relatively long fragments with a
 158 mean fragment length of 173.5 bp and achieved an aver-
 159 age coverage on the reference genome of 102.6X. The
 160 other library was sequenced on an Illumina MiSeq with
 161 a read length of 151 bp. It was produced from shorter
 162 fragments with a mean fragment length of 88.1 bp and
 163 a mean coverage of 49.3X. With its shorter fragments
 164 and lower achieved coverage, the second library better re-
 165 flects typical sequencing libraries generated from aDNA
 166 samples (Sawyer et al., 2012), so we focused our experi-
 167 ments on this library.

168  The distribution of the different read lengths after the different preprocessing steps were performed
 169 is shown in Figure 2. There are many reads that were clipped, trimmed or merged and thus not of
 170 equal length. 

171 Each of the three input read files (generated from the three different preprocessing methods) were
 172 then subject to our two-layer assembly approach. We used both SOAPdenovo2 (Version 2.04) and
 173 MEGAHIT (v1.0.4-beta-3-g027c6b6) in the first layer of the assembly. In order to cover a broad
 174 range of k -mers representing both short and long reads contained in the input, we used ten different
 175 k -mer sizes (37, 47, 57, ..., 127). After removing contigs with no read support, the contigs were then
 176 reassembled with SGA. To identify contigs that belong to the genome of *Mycobacterium leprae*, the
 177 results were mapped against the reference sequence of *Mycobacterium leprae TN*. Contigs that could
 178 be mapped against the reference were extracted and used to compare the assemblies generated in the
 179 different layers.

180 Table 1 shows statistical results of the contigs that could be mapped against the reference genome
 181 of *Mycobacterium leprae TN*. The results that were generated in the second layer are shown as well as
 182 the assembly that generated the longest contig in the first layer using the respective assembly program.
 183 Additionally, results from SGA directly on the *fastq* files as well as results from programs that can
 184 use multiple k -mers in their assembly are shown. It can be seen that when using SOAPdenovo2 in the
 185 first layer, the longest contig, the N50 and the mean contig length could be improved by using SGA to
 186 merge the different assemblies in the second layer. Here, the overall best assembly was derived with
 187 the preprocessing method using the combined trimmed and clipped reads for a single-end assembly
 188 in the first layer. SOAPdenovo2 can also generate its graph using multiple k -mers. The result of
 189 this method is better than using only one k -mer but not as good as our two-layer approach. Using
 190 MEGAHIT, the merging in the second layer with SGA also improved the assemblies generated in the
 191 first layer. MEGAHIT also provides the possibility to generate an assembly using multiple k -mers.
 192 As with SOAPdenovo2, they improve the assembly compared to using only one k -mer but the result
 193 is worse than our two-layer methodology. Another assembly program that can use multiple k -mers to
 194 generate a result is the “interactive *de Bruijn* graph *de novo* assembler” (IDBA) (Peng et al., 2010).
 195 Its res  are very good but not as good as the second layer assembly with SOAPdenovo2 in the first
 196 layer.

197 The length distribution of the resulting contigs is shown in Figure 3. After the second layer
 198 assembly, the number of contigs at the upper end of the length distribution has increased, compared
 199 to the first layer. With MEGAHIT, this is also true, even though it is not as pronounced as in the

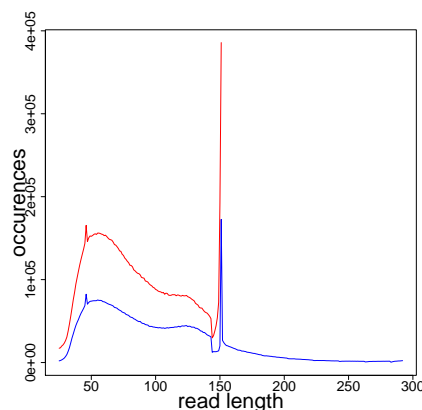



Figure 2. Read length distribution for the different preprocessed *fastq* files. Blue: merged reads, red: RAW reads. 

Table 1. Results using our two-layer assembly with SOAPdenovo2 and MEGAHIT as primary assemblies compared with the standard assemblies of SGA, SOAPdenovo2, MEGAHIT and IDBA on the short fragment library. The results show only values for contigs that could be mapped against the genome of *Mycobacterium leprae*. Here only the best assemblies (based on the longest mapped contig) for the different preprocessing methods and k -mers are shown. “SOAP” alone represents the results using the parameter (-m) resulting in an assembly using multiple different k -mers for the generation of their underlying graph structure. “MEGAHIT” and “IDBA” alone also represent an assembly using multiple internal k -mers. “SOAP K57” and “MEGAHIT K77” represent the best assemblies in the first layer of our pipeline using the respective k -mers of 57 and 77. “SOAP SGA” and “MEGAHIT SGA” show the results of the second layer using SOAPdenovo2 and/or MEGAHIT in the first layer. The column “preprocessing” describes the preprocessing method that was used to generate the result. Values in bold represent the best value that could be achieved. All other statistical values can be found in the supplementary material.

	name	prepro- cessing	# contigs	N50	mean con- tig length	longest contig	# gaps
Layer 1	SOAP	single	249	21909	13210.3	99866	103
	MEGAHIT	merged	175	28410	16777.5	91499	106
	IDBA	paired	164	35419	20152.7	118220	118
	SGA	single	1157	2199	1997.3	8640	952
	SOAP K57	single	215	24962	14918.6	72345	120
	MEGAHIT K77	merged	253	21863	12765.4	87880	108
Lv	SOAP SGA	single	133	42136	25225.0	135656	88
	MEGAHIT SGA	merged	668	19758	12245.3	109259	80

200 assembly using SOAPdenovo2. Using MEGAHIT, the total number of contigs that could be mapped
 201 against the reference genome after the second layer assembly with SGA is significantly higher than
 202 in the individual assemblies of the first layer. There are several more shorter contigs, whereas using
 203 SOAPdenovo2 in the first layer leads to fewer shorter contigs and more longer contigs after the second
 204 layer. Using SGA directly on the preprocessed *fastq* files did not result in good assembly results.

205 Since one normally is interested in one genome of interest (here the genome of the leprosy causing
 206 bacterium), we computed the genome coverage after mapping all contigs of length at least 1000
 207 bases against *Mycobacterium leprae* *TN*. We used Qualimap2 (Okonechnikov et al., 2015) for the
 208 analysis of the mapping. The percentage of the genome that could be covered using only contigs

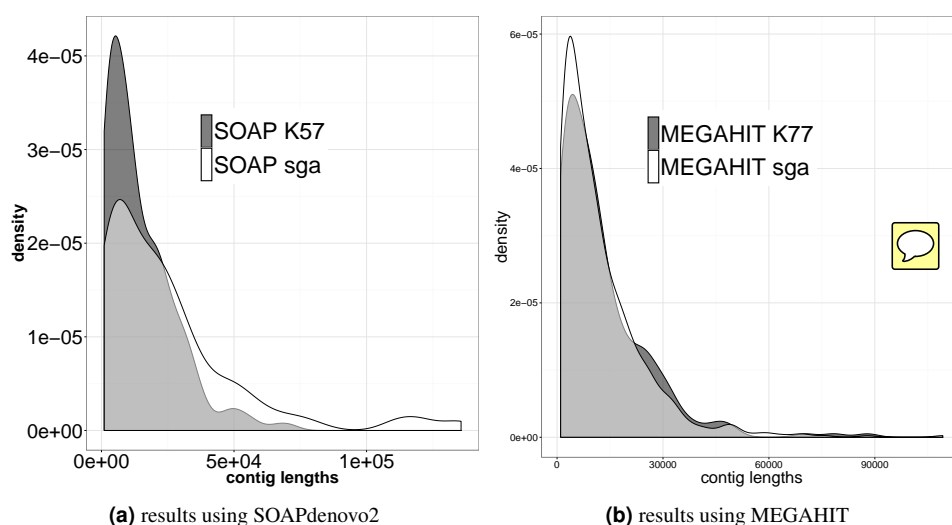


Figure 3. Distribution of the length of the contigs generated by the different assemblies. The results generated by the second layer assembly with SGA is shown in white. The results of first layer assembly is shown in dark grey. The light grey part represents values that belong to both methods. In 3a, the results using SOAPdenovo2 in the first layer are described. The results using MEGAHIT in this layer are shown in 3b.

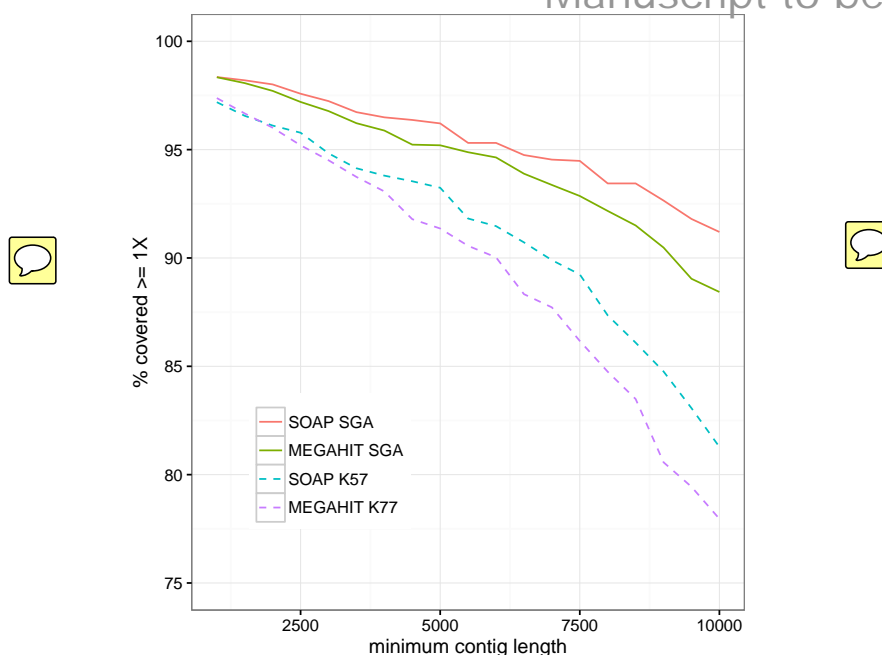


Figure 4. The percentage that could be covered with contigs longer than the minimum contig length.

longer than 1 000, 1 500, . . . , 10 000 bp is shown in Figure 4. It can be seen that the percentage of the genome that could be covered using different cutoffs for the minimum length of the contigs is always higher after the second layer assembly using SGA than using only the results generated in the first layer assemblies. This becomes more and more pronounced with increasing filter threshold for the minimum contig lengths. When using only contigs longer than 1 000 bp, the results are almost the same. Using only contigs longer than 10 000 bp, around 90% of the genome can be covered using the second layer assembly with SGA, whereas at most 80% of the genome is covered by contigs from assemblies generated in the first layer.

The percentage of the genome that was covered at least twice is around 1% for the assemblies generated in the first layer with SOAPdenovo2 and MEGAHIT. This value has increased after the second layer assembly where the contigs were assembled again with SGA, showing that not all overlapping contigs could be identified and merged by SGA.

In order to be able to merge more contigs, we performed a new experiment that also uses the internal error correction of SGA that were described in the previous section. The resulting assembly contained contigs of length $\geq 400,000$ bp that could be mapped against the reference genome. However, when analyzing these contigs, only subsequences of at most 500 bp actually mapped to the genome. The beginning and the end of these contigs were soft-clipped by BWA-MEM and did not map anywhere else on the reference genome. When analyzing the contigs from the assemblies generated without this internal error correction of SGA, the whole contig (with some small insertions and deletions) could be mapped against the reference genome.

The mapping of the contigs generated by the first layer assemblies of SOAPdenovo2 and MEGAHIT against the reference genome resulted in approximately 115 gaps. This value is reduced to around 84 gaps for the contigs generated by the second layer assembly with SGA (see Table 1). These gaps, together with annotated repeat regions of *Mycobacterium leprae*, are shown in Figure 5. It can be seen that the gaps in the mapping of the contigs mainly coincide with annotated repeat regions in the reference genome, as already shown by Schuenemann et al. (2013). Altogether, the percentage resolved regions has dropped from maximally 74.5% (using only SOAPdenovo2) down to 43.5% using our two-layer approach.

Up until now we showed that we were able to generate long, high quality contigs that can be mapped against the reference of *Mycobacterium leprae* TN. In order to show that the assembled contigs actually belong to the species of *Mycobacterium leprae* and not to other *Mycobacteria*, we took the ten longest contigs from each assembly and used BLASTN (Altschul et al., 1990) available on the NCBI webserver to align the contigs with all the genomes available from the genus *Mycobacterium*. The hits that generated the highest score for all of these contigs always belonged to a strain of *Mycobacterium leprae*.

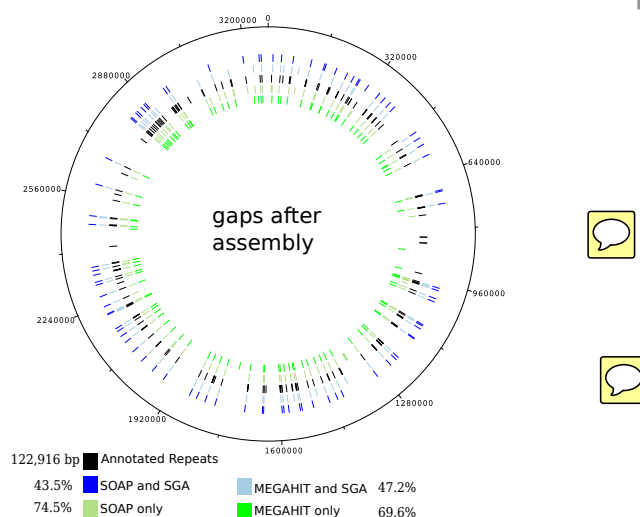


Figure 5. Gaps in the mapping of the contigs against the reference genome of *Mycobacterium leprae* TN together with annotated repeat regions in the reference genome. The outer ring represents the gaps that occur after the mapping of the contigs that were generated by the second layer assembly with SGA after a first layer assembly with SOAPdenovo2. The second outer ring shows the same but for a first layer assembly using MEGAHIT. The middle ring represents the annotated repeat regions of the reference genome. The second inner and innermost ring represent the gaps after using the best individual SOAPdenovo2 and MEGAHIT assemblies, respectively. The percentages represent the relative number of unresolved bases in annotated repeat regions (in total 122,916 bp).

244 While previous analyses confirmed the specificity of mapped contigs, there were several long
 245 contigs that could not be mapped against the reference of *Mycobacterium leprae* TN. This is not
 246 surprising, because DNA from ancient bones is often mixed with other DNA and thus a metagenomic
 247 sample. For this experiment, the longest contig that could not be mapped against the reference
 248 and aligned it against the whole nr/nt database with BLASTN. The best hits achieved only a query
 249 coverage of approximately 13%. These regions on the query are not consecutive and map to different
 250 genes. The most promising gene that can be identified is the heat shock protein 70, which is a highly
 251 conserved gene among several bacteria (Bukau and Horwich, 1998). The same is true for the very
 252 long contigs generated using the correction steps of SGA or the iterative graph construction approach
 253 of MEGAHIT. There is not one species in the database where more than 15% of these queries could
 254 be aligned to.

255 In order to see how this two-layer assembly handles sequencing libraries of lower mean coverage
 256 of the desired species, we performed several experiments of different samples that showed a mean
 257 coverage between 2X and 7X after being mapped with the EAGER pipeline against the respective
 258 reference genome. Here we could not achieve any meaningful results.

259 Furthermore we evaluated the scalability of our pipeline through subsampling. We used the library
 260 from the Jorgen625 sample with the longer fragments as it contained more than twice as many reads
 261 ($2 \times 15,101,591$ instead of $2 \times 6,751,711$ reads). We evaluated the whole pipeline using 1, 2, 5, 10
 262 and all 15.1 million reads. The calculations were performed on a server with 500GB memory and 32
 263 CPUs of type Intel® XEON® E5-416 v2 with 2.30 GHz using four threads wherever parallelization
 264 was possible. The results shown in Figure 6 show that the runtime scales linearly with the number of
 265 input reads. The time it would take to assemble a human genome using our two-layer approach can be
 266 estimated using this linear model. The ancient human LBK/Stuttgart sample published by Lazaridis
 267 et al. (2014) was sequenced using eight lanes, each containing between 200 and 230 million reads.
 268 The assembly of one such lane would take approximately one week and the assembly of all 1.74
 269 billion reads almost two months.

270 DISCUSSION AND CONCLUSIONS

271 With ancient genome assembly one faces a number of challenges. The underlying dataset stems from
 272 a metagenomic sample with short fragments. When performing a paired-end sequencing experiment,
 273 this results in mostly overlapping forward and reverse reads. Because of the highly different read
 274 lengths after the necessary preprocessing steps, including adapter removal and quality trimming,

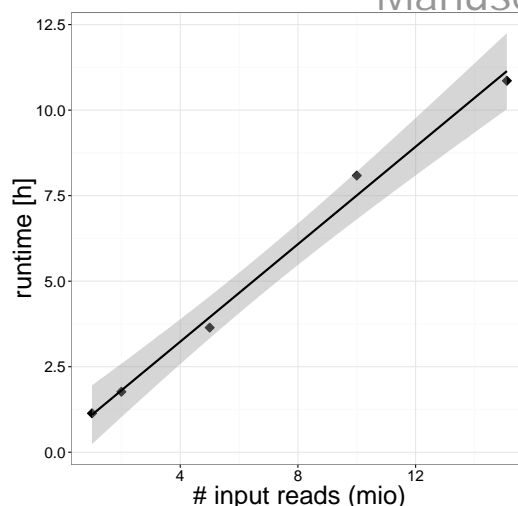


Figure 6. Runtime scaling of the two-layer assembly approach. The black dots show the runtime in minutes using 1, 2, 5, 10 and all 15.1 million input reads. The black line shows the fitted linear regression and the grey area represents the 95% confidence region.

275 typical *de Bruijn* approaches using a fixed k -mer size cannot sufficiently assemble the sample. On
 276 the other hand overlap-based approaches alone are also inferior. Our two-layer approach combining
 277 various assemblies using different k -mer sizes followed by a second assembly based on string overlaps
 278 is able to fuse the contigs generated in the first layer into longer contigs and reduce the redundancy.
 279 Additionally, we could show that longer, high quality contigs are generated after the second layer
 280 assembly. In particular, at least for our example of a *Mycobacterium leprae* genome, these longer
 281 contigs are able to close more gaps, mainly spanning repetitive regions. The different values for k that
 282 are used in the first layer assembly lead to similar contigs that can be combined in the second layer
 283 assembly. The percentage of the genome that is covered more than once is increased after the second
 284 layer assembly. This proves that SGA is not able to identify and merge all overlapping contigs. One
 285 reason for this could be the underlying metagenomic sample. Multiple species in the sample share
 286 similar but not identical sequences. As SGA is not designed to assemble metagenomic samples, these
 287 differences cannot be distinguished from different sequences of the same genome containing small
 288 errors. One possibility to solve this could be to optimize the parameters that SGA provides, as the
 289 current parameters for SGA cannot merge all relevant contigs. This probably has to be adapted for
 290 each sample. However, we showed that when using the steps to account for sequencing errors, the
 291 resulting contigs became worse, when considering the specificity of the contigs (of the organism of
 292 interest). We believe that it could be a problem of multiple *Mycobacteria* in the sample that share
 293 similar sequences which then combined to sequences that are built-up out of fragments of different
 294 species in the sample. The contigs that are generated without these steps are of high quality and
 295 map almost perfectly against the reference sequence that is known to be highly similar to the desired
 296 genome (Mendum et al., 2014). When assembling metagenomic and especially aDNA samples, the
 297 results always have to be regarded critically in order to avoid mistakes. Another possibility could
 298 be sequencing errors in the sample, leading to distinct contigs using different k -mers. However,
 299 these errors should be averaged out by the different assemblies (Schatz et al., 2010). Erroneous base
 300 incorporations can be out as the sample was treated with *Uracil-DNA Glycosylase* (UDG),
 301 removing these errors.

302 An important step is the preprocessing of the raw reads. We compared the performance using
 303 all reads as single reads, as paired reads or as merged reads. However, at least from our study, we
 304 can conclude that the results highly depends on the first layer assembler and probably also on the
 305 dataset itself. What is interesting is the fact that SOAPdenovo2 produces better results when using all
 306 input reads in a single-end assembly than in a paired-end assembly. One possible explanation is that
 307 the information between the pairs does not contain additional information as almost all paired-end
 308 reads overlap and can be merged. It is possible that the program then disregards some overlaps in
 309 order to fulfill the paired-end condition. Overlaps that were disregarded this way could be used in
 310 the single-end assembly leading to a better assembly. Additionally, reads that did not have a partner
 311 were removed before the paired-end assembly. These reads are of course available in the single-end
 312 assembly. It could be that they contained some relevant information.

The mapping of the assembled contigs against the reference show that in our case, all gaps align with annotated repeat regions. Using our two-layer assembly approach, more of these regions could be resolved, but many still remain. In sequencing projects of modern DNA, repetitive regions are resolved using other sequencing technologies such as PacBio. It can produce much longer sequences that span these regions. However, these technologies are not applicable to aDNA as most of the fragments contained in the sample are even shorter than the sequences that can be produced using the Illumina platforms.

In general, it can be concluded that assembly of aDNA is highly dependent on the amount of endogenous DNA in the sample. We are able to improve results generated by current assembly programs. However, the information gain generated by the second layer assembly is dependent on the quality of the first layer assemblies. Thus if the first layer assemblies are of low quality, the second layer assembly cannot significantly improve them.

The runtime scales linearly with the number of input reads, which is no problem for small bacterial datasets. However, big projects like the assembly of human specimen does not seem to be feasible. Nevertheless it has to be kept in mind that the current pipeline currently consists of bash scripts that have not been optimized for parallelization. Using more threads on optimized code might make this approach feasible even for large genomes.

We have shown that the concept of our two-layer approach can improve the assembly of aDNA samples. The results in this study were generated using several scripts. In order to facilitate other researchers to use our two-layer approach, we are currently developing an automated pipeline containing all the steps described above. In the meantime, we provide a shell script that can perform the two-layer assembly with SOAPdenovo2 in the first layer up to the removal of small contigs after the second layer assembly. This script as well as the supplementary material can be downloaded from <https://lambda.informatik.uni-tuebingen.de/gitlab/seitz/MADAM>.

ACKNOWLEDGMENTS

We want to thank Linus Backert and Vladimir Piven for their critical assessment and discussions on different ideas, as well as implementation specific details. Additionally we want to thank Alexander Peltzer, André Hennig, Michael Römer and Niklas Heinsohn for valuable insights and discussions on different ideas, implementation specific questions and for their comments on this paper.

REFERENCES

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–10.
- Avila-Arcos, M. C., Cappellini, E., Romero-Navarro, J. A., Wales, N., Moreno-Mayar, J. V., et al. (2011). Application and comparison of large-scale solution-based DNA capture-enrichment methods on ancient DNA. *Sci. Rep.*, 1:74.
- Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. a., Dvorkin, M., et al. (2012). SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology*, 19(5):455–477.
- Bentley, D. R., Balasubramanian, S., Swerdlow, H. P., Smith, G. P., Milton, J., et al. (2008). Accurate whole human genome sequencing using reversible terminator chemistry. - Supplement. *Nature*, 456(7218):53–9.
- Bukau, B. and Horwich, A. L. (1998). The Hsp70 and Hsp60 chaperone machines. *Cell*, 92(3):351–366.
- Chao, R., Yuan, Y., and Zhao, H. (2015). Recent advances in DNA assembly technologies. *FEMS Yeast Research*, 15(1):1–9.
- Der Sarkissian, C., Allentoft, M. E., Ávila-Arcos, M. C., Barnett, R., Campos, P. F., et al. (2015). Ancient genomics. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, 370(1660):20130387.
- Durai, D. A. and Schulz, M. H. (2016). Informed kmer selection for de novo transcriptome assembly. *Bioinformatics*, page btw217.
- Earl, D., Bradnam, K., St. John, J., Darling, A., Lin, D., et al. (2011). Assemblathon 1: A competitive assessment of de novo short read assembly methods. *Genome Research*, 21(12):2224–2241.
- Ferragina, P. and Manzini, G. (2000). Indexing Compressed Text. *Journal of the ACM*, 52(4):552–581.
- Gnerre, S., Maccallum, I., Przybylski, D., Ribeiro, F. J., Burton, J. N., et al. (2011). High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proceedings of the National Academy of Sciences of the United States of America*, 108(4):1513–8.

- 369 Gnirke, A., Melnikov, A., Maguire, J., Rogov, P., Leproust, E. M., et al. (2009). Solution Hybrid
370 Selection with Ultra-long Oligonucleotides for Massively Parallel Targeted Sequencing. *27(2)*:182–
371 189.
- 372 Hodges, E., Xuan, Z., Baliya, V., Kramer, M., Molla, M. N., et al. (2007). Genome-wide in situ exon
373 capture for selective resequencing. *Nature genetics*, *39(12)*:1522–7.
- 374 Knapp, M. and Hofreiter, M. (2010). Next generation sequencing of ancient DNA: Requirements,
375 strategies and perspectives. *Genes*, *1(2)*:227–243.
- 376 Lazaridis, I., Patterson, N., Mittnik, A., Renaud, G., Mallick, S., et al. (2014). Ancient human genomes
377 suggest three ancestral populations for present-day Europeans. *Nature*, *513(7518)*:409–413.
- 378 Li, D., Liu, C. M., Luo, R., Sadakane, K., and Lam, T. W. (2014). MEGAHIT: An ultra-fast
379 single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph.
380 *Bioinformatics*, *31(10)*:1674–1676.
- 381 Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
382 *arXiv preprint arXiv*, 00(00):3.
- 383 Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform.
384 *Bioinformatics*, *25(14)*:1754–1760.
- 385 Li, R., Zhu, H., Ruan, J., Qian, W., Fang, X., et al. (2010). De novo assembly of human genomes
386 with massively parallel short read sequencing. *Genome Research*, *20(2)*:265–272.
- 387 Li, Z., Chen, Y., Mu, D., Yuan, J., Shi, Y., et al. (2012). Comparison of the two major classes of
388 assembly algorithms: Overlap-layout-consensus and de-bruijn-graph. *Briefings in Functional*
389 *Genomics*, *11(1)*:25–37.
- 390 Luo, R., Liu, B., Xie, Y., Li, Z., Huang, W., et al. (2012). SOAPdenovo2: an empirically improved
391 memory-efficient short-read de novo assembler. *GigaScience*, *1(1)*:18.
- 392 Manber, U. and Myers, G. (1993). Suffix Arrays: A New Method for On-Line String Searches. *SIAM*
393 *Journal on Computing*, *22(5)*:935–948.
- 394 Maricic, T., Whitten, M., and Pääbo, S. (2010). Multiplexed DNA sequence capture of mitochondrial
395 genomes using PCR products. *PLoS ONE*, *5(11)*:9–13.
- 396 Mendum, T. a., Schuenemann, V. J., Roffey, S., Taylor, G. M., Wu, H., et al. (2014). Mycobacterium
397 leprae genomes from a British medieval leprosy hospital: towards understanding an ancient
398 epidemic. *BMC genomics*, *15(1)*:270.
- 399 Myers, E. W. (2005). The fragment assembly string graph. *Bioinformatics*, *21(suppl 2)*:ii79–ii85.
- 400 Okonechnikov, K., Conesa, A., and García-Alcalde, F. (2015). Qualimap 2: Advanced multi-sample
401 quality control for high-throughput sequencing data. *Bioinformatics*, *32(2)*:292–294.
- 402 Peltzer, A., Jäger, G., Herbig, A., Seitz, A., Kniep, C., Krause, J., and Nieselt, K. (2016). EAGER:
403 Efficient Ancient Genome Reconstruction. *Genome Biology*, *17(1)*:60.
- 404 Peng, Y., Leung, H. C. M., Yiu, S. M., and Chin, F. Y. L. (2010). IDBA - A practical iterative De Bruijn
405 graph De Novo assembler. In *Annual International Conference on Research in Computational*
406 *Molecular Biology*, volume 6044, pages 426–440. Springer.
- 407 Rasmussen, M., Li, Y., Lindgreen, S., Pedersen, J. S., Albrechtsen, A., et al. (2010). Ancient human
408 genome sequence of an extinct Palaeo-Eskimo. *Nature*, *463(7282)*:757–762.
- 409 Sawyer, S., Krause, J., Guschanski, K., Savolainen, V., and Pääbo, S. (2012). Temporal patterns of
410 nucleotide misincorporations and DNA fragmentation in ancient DNA. *PLoS ONE*, *7(3)*.
- 411 Schatz, M. C., Delcher, A. L., and Salzberg, S. L. (2010). Assembly of Large Genomes using Cloud
412 Computing. *Genome research*, *20(9)*:1165–1173.
- 413 Schubert, M., Ermini, L., Der Sarkissian, C., Jónsson, H., Ginolhac, A., et al. (2014). Characterization
414 of ancient and modern genomes by SNP detection and phylogenomic and metagenomic analysis
415 using PALEOMIX. *Nature Protocols*, *9(5)*:1056–1082.
- 416 Schuenemann, V. J., Singh, P., Mendum, T. A., Krause-Kyora, B., et al. (2013). Genome-Wide
417 Comparison of Medieval and Modern Mycobacterium leprae. *Science*, *179(July)*:179–184.
- 418 Simpson, J. T. and Durbin, R. (2012). Efficient de novo assembly of large genomes using compressed
419 data structures. *Genome Research*, *22(3)*:549–556.
- 420 Utturkar, S. M., Klingeman, D. M., Land, M. L., Schadt, C. W., Doktycz, M. J., Pelletier, D. A., and
421 Brown, S. D. (2014). Evaluation and validation of de novo and hybrid assembly techniques to
422 derive high-quality genome sequences. *Bioinformatics*, *30(19)*:2709–2716.