

Chapter II: Methods for Microbial Genome Reconstruction from Ancient DNA

Luis A. Arriola

October 2017

Working Title: Comparative assessment of methods for reconstruction of whole genomic sequences from ancient microbial organisms

ABSTRACT

Technological and analytical advances in the last decade revolutionized the study of ancient microorganisms by allowing researchers to reconstruct draft genomes of historical pathogenic organisms. Reference mapping, using *BWA aln* with parameters optimized for ancient DNA (aDNA), is the preferred used algorithm for microbial genome reconstruction. However, the amount of information that this approach can recover is highly dependent on genomic sequences of modern organisms and presents certain limitations and biases. New mapping algorithms and alternative approaches, such as *de novo* assembly or iterative mapping, have shown promising results in reconstructing microbial and mitochondrial chromosomes and could help in the recovery of genomic information from highly divergent microbial organisms.

Here we present a comparative assessment of methods for reconstruction of whole genomic sequences using simulated data mirroring conditions commonly found in ancient metagenomic datasets (contamination, damage, and divergence).

We use information from this comparative assessment to guide the reconstruction of *Methanobrevibacter oralis*, the oldest ancient microbial genome draft reported, providing % more information with % fewer errors.

Iterative mapping approaches provide the means to study the evolution of ancient microbial organisms using ancient DNA as they allow to recover highly variable genomic sequences inaccessible through commonly used reference mapping methods.

INTRODUCTION

In recent years significant technological and analytical advances have revolutionized the study of ancient microbial organisms. High-throughput sequencing technologies, such as Illumina NGS and enrichment capture, have improved the recovery of endogenous DNA sequences of microorganisms hundreds to thousands of years old (Bos *et al.*, 2011, 2016; Feldman *et al.*, 2016; Hofreiter *et al.*, 2015). However, few progress has been made in the approaches used to reconstruct and analyze whole genomic sequences of these organisms (Hofreiter *et al.*, 2015; Parks and Lambert, 2015; Schubert *et al.*, 2012). While single or multi-locus analyses offer a convenient approach to identify and classify organisms, they offer reduced power to resolve phylogenetic relationships, study the evolution and biology of an organism, and in some cases can lead to misleading results due to its low resolution (REF). On the other hand, genome-wide information recovered from historical sources provides important insights into the biology and evolution of ancient organisms (REF) as it allows researchers to identify organisms unambiguously, refine models, and test hypothesis created from current genetic or genomic data at a higher resolution (REF). Differences between genomes from the same species or closely related species are used to establish more detailed phylogenetic relations, calculate mutation rates, and identify patterns of adaptation and evolution over time at the individual and population level (Parks *et al.*, 2015). Moreover, dated ancient genomic information adds a temporal dimension to genomic variation and provides details of the geographic distribution of the organisms in the past (REF).

The reconstruction of a genomic sequence is a critical step in the study of an organism and its evolution since all downstream analyses, comparisons and inferences depend and will be affected by the quality and fidelity of it (Giese *et al.*, 2014; Leonardi *et al.*, 2016; Schubert *et al.*, 2012). Currently, two main approaches are used to reconstruct whole microbial genomes from high-throughput sequencing (HTS) data: ‘reference mapping’ and ‘*de novo* assembly’ [REF]. Mapping assembly approaches use an already available reference sequence as a guide to locate and align fragments to their most likely locations based on sequence similar-

ity. ‘*De novo* assembly’ approaches, are reference-free reconstructions that match overlapping sequences to create longer contiguous DNA reads (contigs) (Hofreiter *et al.*, 2015). In palaeomicrobiology (i.e., the study of ancient microbial organisms), the most widely used approach for whole genome reconstruction is reference mapping assembly. Using sequences from modern microbial organisms as references, researchers have been able to reconstruct draft genomic sequences from a variety of ancient microbial pathogens at various points in the past, including *Yersinia pestis* (Bos *et al.*, 2011, 2016; Feldman *et al.*, 2016; Rasmussen *et al.*, 2015; Wagner *et al.*, 2014), *Mycobacterium leprae* (Mendum *et al.*, 2014; Schuenemann *et al.*, 2013), *Mycobacterium tuberculosis* (Bos *et al.*, 2014; Kay *et al.*, 2015), *Helicobacter pylori* (Maixner *et al.*, 2016), *Treponema denticola* (Maixner *et al.*, 2014), *Brucella melitensis* (Kay *et al.*, 2014), and *Methanobrevibacter oralis* subsp. *neandertalensis* (Weyrich *et al.*, 2017). The reference mapping assembly approach is advantageous because it works well with low abundance and short fragment lengths characteristic of aDNA (Dabney *et al.*, 2013; Schubert *et al.*, 2012), and it allows the recovery of endogenous sequences in a complex system (i.e., metagenomic shotgun sequence datasets) (Maixner *et al.*, 2016). However, due to the stringency of this approach, it recovers highly conserved regions of the genome, and this hinders the ability to find regions that have greatly changed through time between the target organism and the used reference (Hofreiter *et al.*, 2015). Moreover, the amount of sequence recovered is highly dependent on the availability, quality, and choice of the sequence used as a reference, and the divergence and the heterogeneity of the target organism. Mapping assembly algorithms are affected by DNA damage patterns which result in lower yields of coverage, as reads are mapped with lower confidence, and in some cases allow reads to artificially map to other genomic locations with similar composition or the introduction of non-endogenous reads (Schubert *et al.*, 2012). The most widely used algorithm in aDNA studies is the Burrows-Wheelers Aligner (BWA), which specializes in short fragment mapping (Li and Durbin, 2009; Orlando *et al.*, 2015). In 2012, Schubert *et al.* explored a series of parameters to improve the mapping of aDNA reads to modern reference genomes. In this survey, they showed that by deactivating the use of a seed region and increasing tolerance to higher edit distances, it was possible to increase the number of high-quality endogenous hits recovered, albeit in low proportions for Illumina data (Schubert *et al.*, 2012). Since then, these findings have been used by the aDNA community as the gold standard for mapping assembly. In contrast, the de

novo assembly approach has only been successful in reconstructing an ancient microbial genome under very special circumstances (Schuenemann *et al.*, 2013, [and *Yersinia pestis*]). As this method does not depend on a reference sequence, it allows (in theory) a less biased reconstruction. The main disadvantage of the method is that it is highly dependant on the quality of the dataset, as the low fragment length and the presence of non-endogenous sequences complicate a de novo reconstruction (Westbury *et al.*, 2017).

In addition to reference-based mapping and de novo assembly, other methods have been developed to assemble small genomes. These include reference guided assembly, which makes use of the advantages of de novo assembly to reconstruct contigs that could contain highly variable regions, and a known reference sequence as a scaffold to locate the most likely positions for these contigs (Rajaraman *et al.*, 2013). Another novel method is iterative mapping. This method uses an initial mapping or de novo step to create self-standing reference sequences that will be cyclically extended via a reference mapping approach until no more new reads can be added (REF). This method has been used successfully to reconstruct mitochondrial sequences, with the advantage of not having to rely on a full genomic reference sequence and be able to work on complex datasets (metagenomics) (Hahn *et al.*, 2013; Westbury *et al.*, 2017).

Reconstructing microbial genomes presents certain advantages, as they count with small haploid genomes, which allow for reduced mapping/assembly times compared to higher organisms; however, some disadvantages also need to be considered. Ancient bacterial genome reconstruction represents a challenge due to lower yields of coverage derived from post-mortem damage. Other complications include genetic heterogeneity, a higher genetic divergence from modern organisms, lateral gene transfer, lower presence in samples and contamination. Many different strategies have been used in palaeomicrobiology studies to reconstruct whole genome information (**Table 01**). Nevertheless, to our knowledge, no formal comparison of different genome reconstruction approaches has been conducted, nor has the impact of post-mortem aDNA damage on the reconstruction of genomes from microorganisms been investigated. Therefore it is essential to explore how different conditions found on aDNA, such as short read sizes, DNA damage, and genomic divergence, and contaminant reads affect the mappability of aDNA reads and the reconstruction of microbial genomes.

Table 01. Whole-genome reconstruction

RESULTS

In this study, we compared three different approaches for reconstructing microbial genomes from ancient DNA and explored how they are affected by characteristics commonly found in ancient metagenomic datasets. We used the genomic reference sequence of *Streptococcus mutans* UA159 and simulated insertions, deletions, and substitutions on it at five increasing levels to generate five target genomes with variable divergence levels. We included each target genome in three different contexts of contamination (No-Contamination (**0C100E**); Inter-species Contamination (**90C10E**); Intra-genera Contamination (**80Smu20E**)), and simulated reads with ancient damage profiles at three levels of deamination (0.1, 0.3, 0.5), for a total of 45 read pools. Then, we used *reference mapping*, *reference-guided de novo assembly*, and *iterative mapping* approaches to reconstruct the original target genomes from each pool.

In the reference mapping approach, we used the original genomic sequence of *S. mutans* UA159 as a reference and mapped reads with *MIRA* (using default parameters specific for Illumina reads) and the programs *Bowtie2* and *BWA* (using three different sets of parameters, including parameters optimized for aDNA). In the reference assisted assembly, we generated contigs *de novo* using *Velvet* (v.1.2.10) through *VelvetOptimiser* (v2.2.5) and mapped the resulting contigs to the reference with the previously mentioned mappers and set of parameters. Finally, in the iterative mapping approach, we used *MITObim* with mapping and *de novo* initial steps, and different values for allowed mismatches. The resulting bam alignments were filtered by mapping quality ($MQ \geq 30$), and consensus files were called using *samtools* (v.1.3.1)/*bcftools* (v.1.3.1). For *MIRA* and *MITObim*, we used the final ungapped *FASTA* consensus sequence generated by these programs. The final consensus sequences of each approach were divided into “consensus contigs” wherever more than five consecutive uncalled bases (N’s) were found. Finally, we aligned the “consensus contigs” of each approach to their original target genomes and evaluated the resulting reconstructions using the Qualitative Assessment Tool (QUAST).

Genome fraction recovered (Alignment coverage)

We define *genome fraction recovered* as the percentage of the target genome covered by the alignment of *consensus contigs*. For all reconstruction methods, the *genome fraction recovered* decreased with the increase in divergence between the reference sequence used and target genome. *Iterative mapping* approaches consistently recovered the largest amount of original genomic sequence, with the ones including an initial mapping step recovering more than 60% of the target genome, even at the highest level of divergence. Not surprisingly, the presence of non-endogenous reads affected the *reference guided assembly* approaches the most, in many cases halving the amount of genome fraction recovered by these tools. *Reference mapping* and *iterative mapping* approaches were only slightly affected by contaminants reads, with *iterative mapping* approaches being the most affected by non-endogenous reads from closely related species reducing the genome fraction recovered by these approaches more than 5%.

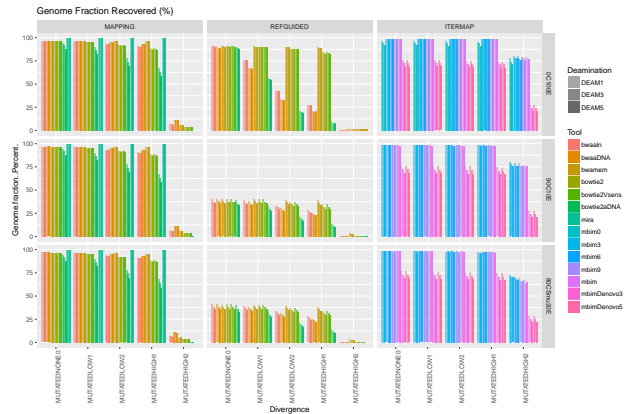
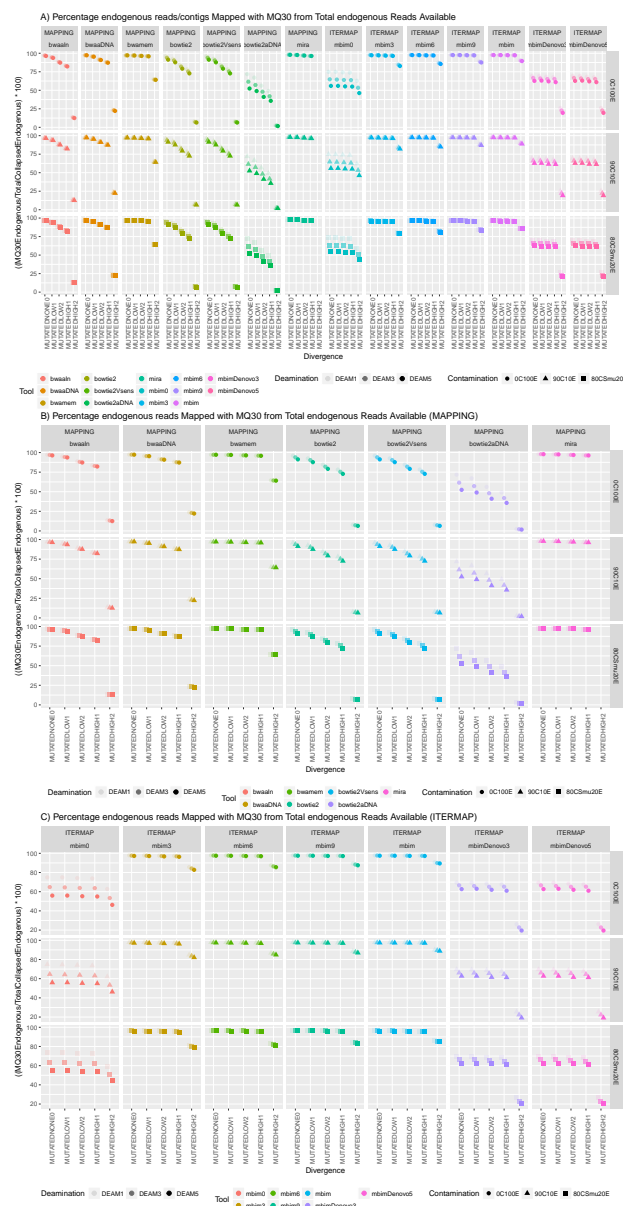


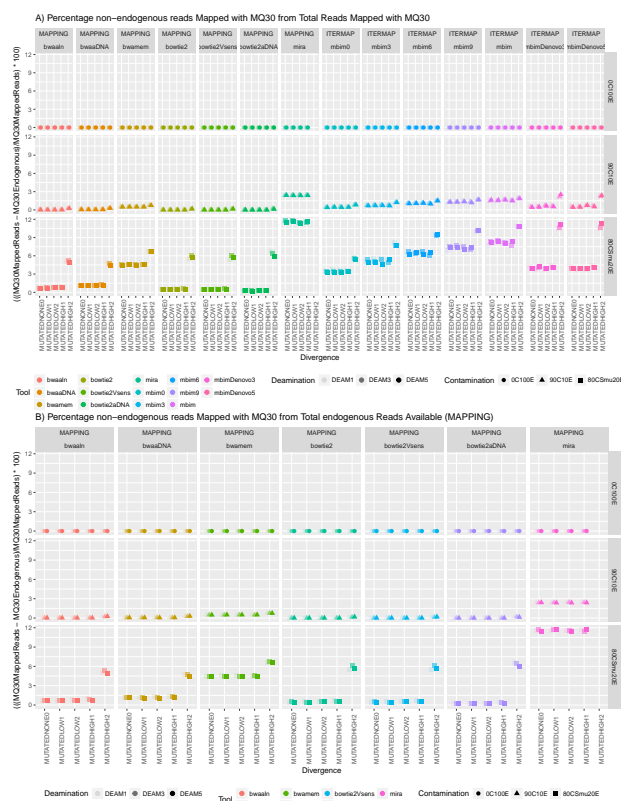
Figure #.- Percentage of target genome covered by the alignment of consensus contigs generated by each method.

Accuracy of reconstruction

To evaluate the accuracy of the reconstruction, we compared the average number of mismatches and indels per 100 000 aligned bases that each approach introduced. *Reference mapping* approaches had the lowest values of average mismatches at lower values of divergence, but this values raised with the increase in divergence. *MIRA* mapper was the most affected by divergence increases (**Sup. Figure #**). Interestingly, in the *reference mapping approaches* the presence of

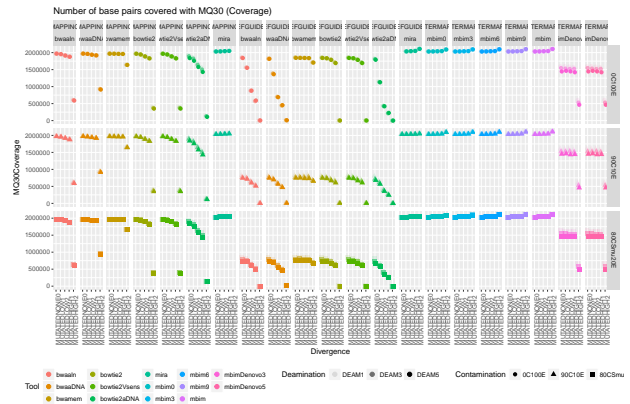
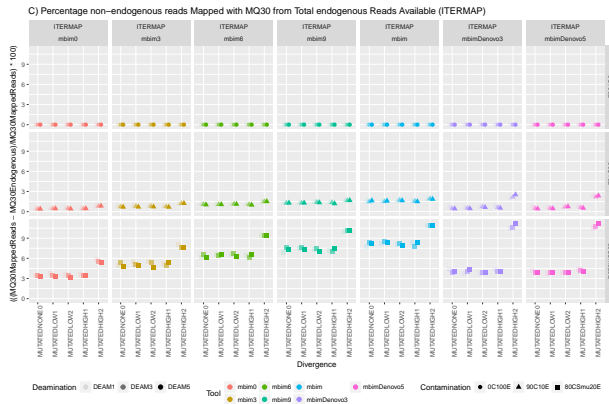


related species being the most affected across all tools and methods. The iterative approach is the most affected, allowing the introduction of more than 2% of non-endogenous reads in the ‘metagenomic’ dataset, and more than 10% when highly related species are present, and the divergence between the reference and the target genome is the highest. From the mapping approach, *MIRA* is the most affected mapper, followed by *BWA mem*. The approaches using *Bowtie2* performed the best and only included around 6% non-endogenous reads when highly similar reads were present, and the divergence was the highest. Again, the impact of deamination in the reads was not significant.



Mapping Accuracy

As a measurement of the accuracy of the Reference Mapping and Iterative Mapping approaches, we calculated the proportion of non-endogenous reads mapped with high-confidence ($\text{MQ} \geq 30$) from the total amount of mapped reads with mapping quality 30. This value gives us an idea of how much ‘contamination’ could be added from other genomes, even after filtering by mapping quality. As expected, we observed that the type of contamination has a substantial impact on the number of non-endogenous reads recovered, with the dataset composed of closely

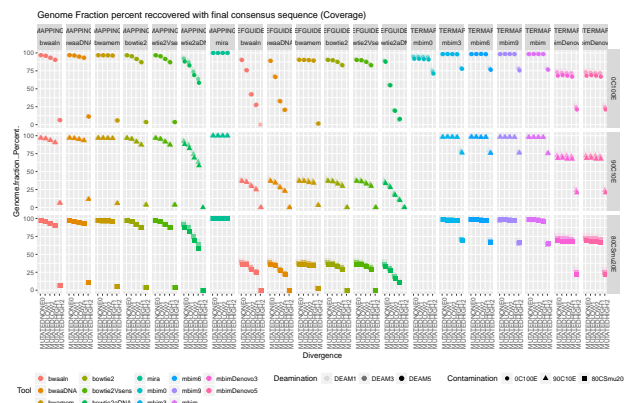


Warning: Removed 11 rows containing missing values (ge

Effects of Divergence

Effects of Contamination

Effects of Damage profiles



Warning: Removed 11 rows containing missing values (ge

COVERAGE, GENOME FRACTION PERCENT, TOTAL ALIGNED LENGTH

Number of base pairs covered with mapping quality $MQ \geq 30$ across three methods and the different tools used on each.

- Coverage is calculated by counting every position with depth of coverage equal or greater than 1.
- Genome Fraction percent is the percentage of aligned bases in the target genome covered by aligning the resulting “contigs” of the consensus sequence generated for each alignment.
- Total aligned length is the total number of aligned bases in the assembly

DEPTH OF COVERAGE

Depth of coverage is calculated with samtools, and is the average of the depth of coverage across the full length of the genome.

Warning: Removed 12 rows containing missing values (ge



To establish the baseline efficiency of the methods in a contamination-free scenario, we used the five target genomes with different divergence levels and simulated on each three levels of cytosine deamination (i.e. 0.1, 0.3, 0.5), for a total of 15 subsets. Each subset contained only endogenous pair-end sequences of the target genomes at a depth of coverage of 5X.

- Dataset 2: Reconstruction of genomes with contamination (Metagenomic | interspecific)

To simulate background microbial contamination from a metagenomic context and to evaluate the effect of non-endogenous reads on the reconstruction of bacterial genomes, we generated simulated Illumina pair-end reads with aDNA damage profiles of each target genome plus 47 bacterial genomes (28 of those found in the Human Oral Microbiome) from 11 different Phyla (Table 03). Each subset contained simulated reads of the target genomes at a 5X depth of coverage equivalent to a 10% of the total reads, and 90% non-endogenous reads.

- Dataset 3: Reconstruction of genomes with contamination of similar species

We generated a dataset with 13 bacterial genomes from the genus *Streptococcus* (Table 04) to evaluate the ability of these approaches to reconstruct the original target genomes when intraspecific contamination is present. Each dataset contained simulated sequencing data with aDNA damage profiles and a composition of 80% non-endogenous reads and 20% endogenous reads, equivalent to a 5X depth of coverage of the target genomes.

Finally, we used AdapterRemoval (v2.2.0, <https://github.com/MikkelSchubert/adapterremoval>) to collapse reads and remove adapters, allowing a minimal length of 25, minimal quality of 4, and minimum alignment length of 11.

Genome Reconstruction Approaches

We tested three different approaches for whole genome reconstruction: reference mapping assembly; reference guided de novo assembly; and iterative mapping. For the mapping assembly approach we used the programs BWA and Bowtie2, each with different parameters combinations used in published palaeomicrobiology studies (Table 01 & 02), and the original genomic sequence of *Streptococcus mutans* UA159 as the reference sequence. For the reference assisted de novo as-

sembly approach we used VelvetOptimiser with kmer sizes 21-45 to create contigs de novo, and used the same algorithms of the reference mapping approach to map the resulting contigs to the reference sequence of *S. mutans* UA159. Finally for the iterative mapping approach we used MIRA 4 to create an initial mapping assembly to the reference sequence and MITObim (v1.9) to iteratively map reads, and MITObim (v1.9) with denovo assembly parameters.

- we changed the number of mismatches allowed by MITObim to check the stringency...

Evaluation and Comparison Stats

To evaluate the performance of the different approaches for whole genome assembly we used three main parameters: Percentage of original sequence recovered, Fidelity of the recovered sequence to the original genome, and the ratio of non-endogenous reads vs. endogenous reads.

For the reference mapping and reference guided de novo assembly, we used samtools to filter reads with mapping quality ≤ 30 and to calculate basic mapping statistics (Coverage, Depth of Coverage, Number of sequences mapped, etc.). We removed duplicates using picard tools and generated a consensus sequence fasta file using bcftools. We split the consensus sequence in “contigs” whenever we found more than 5 “N”. We used these “contigs” to evaluate the assembly using the Quality ASsessment Tool for genome assemblies (QUAST) and the original target genomes.