

1 **Assembly by Reduced Complexity (ARC): a hybrid approach for targeted assembly of**
2 **homologous sequences.**
3
4

5 Samuel S. Hunter¹, Robert T. Lyon¹, Brice A. J. Sarver^{1,2}, Kayla Hardwick²,
6 Larry J. Forney^{1,2}, Matthew L. Settles^{1,2}
7
8

9 ¹Institute for Bioinformatics and Evolutionary Studies, University of Idaho, Moscow, ID 83844-3051,

10 ²Department of Biological Sciences, University of Idaho, Moscow, ID 83844-3051
11
12
13

14 **Keywords:**

15 Bioinformatics software, Subgenome assembly
16
17
18

19 **Correspondence:**

20 Matthew L. Settles
21 Institute for Bioinformatics and Evolutionary Studies
22 University of Idaho
23 Moscow, ID 83844-3051
24 Phone: (208) 885-6051
25 Fax: (208) 885-5003
26 Email: msettles@uidaho.edu
27
28

29 **Running Title:**

30 Assembly by Reduced Complexity (ARC)

31 **Abstract**

32 Analysis of High-throughput sequencing (HTS) data is a difficult problem, especially in the
33 context of non-model organisms where comparison of homologous sequences may be hindered by the
34 lack of a close reference genome. Current mapping-based methods rely on the availability of a highly
35 similar reference sequence, whereas *de novo* assemblies produce anonymous (unannotated) contigs that
36 are not easily compared across samples. Here, we present Assembly by Reduced Complexity (ARC) a
37 hybrid mapping and assembly approach for targeted assembly of homologous sequences. ARC is an
38 open-source project (<http://ibest.github.io/ARC/>) implemented in the Python language and consists of the
39 following stages: 1) align sequence reads to reference targets, 2) use alignment results to distribute reads
40 into target specific bins, 3) perform assemblies for each bin (target) to produce contigs, and 4) replace
41 previous reference targets with assembled contigs and iterate. We show that ARC is able to assemble high
42 quality, unbiased mitochondrial genomes seeded from 11 progressively divergent references, and is able
43 to assemble full mitochondrial genomes starting from short, poor quality ancient DNA reads. We also
44 show ARC compares favorably to *de novo* assembly of a large exome capture dataset for CPU and
45 memory requirements; assembling 7,627 individual targets across 55 samples, completing over 1.3
46 million assemblies in less than 78 hours, while using under 32 Gb of system memory. ARC breaks the
47 assembly problem down into many smaller problems, solving the anonymous contig and poor scaling
48 inherent in some *de novo* assembly methods and reference bias inherent in traditional read mapping.

49 INTRODUCTION

50 High-throughput sequencing (HTS) techniques have become a standard method for producing
51 genomic and transcriptomic information about an organism (Schbath et al. 2012). The Illumina, Roche,
52 and Life Sciences sequencing platforms produce millions of short sequences referred to “reads” that range
53 in length from 50 to 700 base pairs (bp) depending on chemistry and platform. In shotgun sequencing,
54 these short reads are typically produced at random, making them effectively meaningless without further
55 analysis. The primary challenge in the analysis of HTS data is to organize and summarize the massive
56 number of short reads into a form that provides insight into the underlying biology. Two analysis
57 strategies, *de novo* sequence assembly and sequence mapping have been widely adopted to achieve this
58 end.

59 The objective of *de novo* assembly is to piece together shorter read sequences to form longer
60 sequences known as contigs. Sequence assembly is a challenging problem that is made more difficult by
61 characteristics of the sequenced genome (e.g., repeated elements and heterozygosity) and by sequencing
62 technology characteristics (e.g., read length and sequencing errors). Additionally, assembly algorithms are
63 computationally intensive for all but the smallest datasets, thus limiting their application (Li et al. 2012).
64 Finally, *de novo* assembly of large datasets typically produces many short contigs that require additional
65 organization and analysis. Despite many advances and a large selection of assembly software packages,
66 fragmentation and misassembly remain common problems and improving the quality of *de novo* sequence
67 assemblies continues to be an area of active research (Bradnam et al. 2013).

68 Sequence mapping is often the first step carried out in resequencing projects where a good
69 reference sequence exists. The objective of mapping is to align short reads against a reference sequence,
70 thereby permitting direct sequence comparisons between a sample and the reference. This approach is
71 significantly faster than *de novo* sequence assembly and has proven to be very effective at identifying
72 sequence variants at a large scale (The 1000 Genomes Project Consortium et al. 2010). Unfortunately, this
73 approach is entirely dependent on a reference sequence that is similar to the organism being sequenced.

74 Differences between a sample and reference sequence (e.g., structural variations (SVs), novel sequences,
75 an incomplete or misassembled reference, or sequence divergence) can result in unmapped or poorly-
76 mapped reads, which may result in false variant calls (Li, 2011). In the context of RNA-Seq experiments,
77 unmapped reads result in counting errors, and can affect the identification of differentially expressed
78 genes (Pyrkosz et al. 2013). Resequencing projects are performed to identify differences between a
79 sample and an established reference; however, the regions that are most divergent can also be the most
80 difficult to map reads against. Because of this, mapping based approaches are inherently biased by the
81 reference and only provide reliable results when sequence divergence is below the threshold at which
82 reads can be mapped accurately.

83 The two approaches described above (mapping and *de novo* assembly) have been developed and
84 optimized for whole-genome analysis; however, another class of problems exists in which specific
85 regions of a genome or subsets of the sequenced DNA are analyzed. This type of analysis is appropriate
86 in many instances, including sequence capture, viral genome assembly from environmental samples,
87 RNA-Seq, mitochondrial or chloroplast genome assembly, metagenomics, and more. In cases like these, it
88 has been necessary to develop custom pipelines to carry out analyses. In order to assemble the mammoth
89 mitochondrial genome from whole genome shotgun data, Gilbert et al. (2007) first mapped the reads to a
90 reference mitochondrial sequence, filtered the mapped reads and then “assembled using scripts to run
91 existing assembly software”. Other tools that have been developed to address sub genome assembly
92 include: MITObim an extension of the MIRA assembler (Hahn et al. 2013; Chevreur et al. 1999). It
93 requires the user to first perform a mapping based assembly with MIRA, then to use the output of this
94 assembly to do iterative read recruitment and assembly; however, according to the documentation,
95 MITObim does not take advantage of paired-end reads for recruitment or extension. Further, it is not
96 optimized for multiple targets or multiple samples, due to many steps that are manually carried out. The
97 Mapping Iterative Assembler (MIA) uses an iterative mapping and consensus calling approach
98 (<https://github.com/udo-stenzel/mapping-iterative-assembler>). The algorithm is tuned for ancient DNA

99 and was reported by Hahn et al. (2013) to be very slow. It also appears to only function with a single
100 sample and reference; Other groups such as Malé et al. (2014), and Picardi and Pesole (2012) have also
101 developed strategies for assembling smaller subsets from larger datasets; however, none of these were
102 developed as a general purpose, highly parallelized homologous sequence assembler.

103 To address this problem we introduce a hybrid strategy, Assembly by Reduced Complexity
104 (ARC) that combines the strengths of mapping and *de novo* assembly approaches while minimizing their
105 weaknesses. This approach is designed for the myriad of situations in which the assembly of entire
106 genomes is not the primary objective, but instead the goal is the assembly of one or many discreet,
107 relatively small subgenomic targets. ARC is an iterative algorithm that uses an initial set of reference
108 sequences (subgenomic targets) to seed *de novo* assemblies. Reads are first mapped to reference
109 sequences, and then the mapped reads are pooled and assembled in parallel on a per-target basis to form
110 target-associated contigs. These assembled contigs then serve as reference sequences for the next iteration
111 (see Figure 1). This method breaks the assembly problem down into many smaller problems, using
112 iterative mapping and *de novo* assembly steps to address the poor scaling issue inherent to some *de novo*
113 assembly methods and the reference bias inherent to traditional read mapping. Finally, ARC produces
114 contigs that are annotated to the reference sequence from which they were initiated from, making across
115 sample comparisons possible with little additional processing.

116 **RESULTS**

117 Experiments were conducted to determine how well ARC performs across an array of
118 progressively more divergent references, assembly of short, poor quality reads produced from ancient
119 DNA samples, and to measure ARC's performance on a large dataset. ARC was tested using two datasets.
120 The first dataset is made up of Illumina sequence reads from two chipmunk (*Tamias* sp.) exome capture
121 experiments (Bi et al. 2012; Sarver et al. in prep). The second dataset consists of Roche 454 FLX
122 sequence reads from a whole-genome shotgun sequencing experiment using ancient DNA extracted from
123 a mammoth hair shaft sample (Gilbert et al. 2007). The workflow and results of these experiments are

124 presented below.

125 **Assembly by Reduced Complexity Workflow**

126 The iterative mapping and assembly principle (Figure 1) and workflow (Figure 2) behind ARC
127 consists of several steps: 1) align sequenced reads to reference targets, 2) use alignment results to
128 distribute reads into target specific bins, 3) perform assemblies for each bin (target) to produce contigs,
129 and 4) replace initial reference targets with assembled contigs and iterate the process until stopping
130 criteria have been met. During the read alignment step (1), either the sequence aligner BLAT, or Bowtie 2,
131 is used to identify reads that are similar to the current reference targets. The assembly step (3) is
132 performed using either the Roche GS De Novo Assembler (aka “Newbler”) or SPAdes assemblers.

133 ARC accepts a plain text configuration file, a FASTA formatted file with reference target
134 sequences, and either FASTA or FASTQ formatted files containing reads for each sample. An output
135 folder is generated for each sample that contains the final set of contigs, the reads recruited on the final
136 iteration, and ARC statistics.

137 ARC is open source software implemented in the Python programming language with source
138 code available for download from GitHub (<http://ibest.github.io/ARC/>). Prerequisite software packages
139 include: Python 2.7.x, Biopython (Cock et al., 2009), BLAT (Kent, 2002) or Bowtie 2 (Langmead and
140 Salzberg, 2012) and Newbler (Margulies et al., 2005) or SPAdes (Bankevich et al., 2012). These software
141 packages are all free and easy to obtain, and may already be available on systems previously used for
142 HTS analysis. ARC can be installed on most Linux servers, but will also work on many desktops or
143 laptops, provided the required prerequisites are installed. The installation size is only 3Mb, and system
144 administrator access is not required, making it easy to download and use. Configuration is done via a
145 plain text file that can be distributed to make replication of results simple.

146 **ARC performs well with divergent references**

147 A divergent reference sequence can result in unmapped and misaligned reads (Li, 2012). To test
148 how robust ARC is to reference sequence divergence, we assembled mitochondrial genomes using reads

149 from an exome sequence capture experiment performed on 55 chipmunk specimens representing seven
150 different species within the *Tamias* genus (*T. canipes*, *T. cinereicollis*, *T. dorsalis*, *T. quadrivittatus*, *T.*
151 *rufus*, *T. umbrinus*, and *T. striatus*) (Bi et al. 2002; Sarver et al. In prep.). We ran ARC using a set of 11
152 mitochondrial references spanning Mammalia, including Eastern long fingered bat, Cape hare, Edible
153 dormouse, Gray-footed chipmunk, Guinea pig, House mouse, Human, Platypus, Red squirrel, Ring-tailed
154 lemur, and Tasmanian devil. Sequence divergence of mitochondrial references with respect to *Tamias*
155 *cinereicollis* ranged (in percent identity) from 71.2% (Platypus) to 94.9% (Gray-footed chipmunk).
156 Generally speaking, the more divergent the reference sequence, the more ARC iterations were needed in
157 order to complete the assembly process (see Figure 3), while still producing the same resulting
158 mitochondrial genome sequence.

159 Supplemental Table 1 reports ARC results for final number of reads recruited and used for
160 assembly (as well as the common count of reads across all 11 reference targets), contig size (total sum of
161 bases across all contigs produced), contig count, ARC iterations needed before stopping criteria were met,
162 and final ARC status (completed or killed) across the 11 reference target sequences for each of the 55
163 samples. Results show that the choice of reference sequence did not qualitatively impact the final result,
164 ultimately producing, in most cases, the same final mitochondrial genome sequence. In general each of
165 the 11 reference target species recruited the same number of reads, produced the same number of
166 contigs, and resulted in the same length product, with the primary difference being the number of ARC
167 iterations conducted before stopping conditions were met. The relationship between target and read
168 recruitment is further illustrated in Supplemental Figure 1, which shows that the most similar target, the
169 Gray-footed chipmunk (*T. canipes*), typically recruits almost the full set of reads in the first iteration and
170 finishes by the third iteration. At the other extreme, platypus recruited a significantly smaller proportion
171 of reads on the first iteration, but continues to recruit more reads at each iteration until it acquires the full
172 (and often, the same) set of reads. Quality of the original read dataset attributed more to determining the
173 success of assembly than choice of reference sequence.

174 We observed some variation in final contig lengths across the reference sequences; however, this
175 can be attributed to the linearization of the circular mitochondrial genome. As an example, ARC
176 assembled a single contig for sample S160 across all 11 references, with the length of contig differing by
177 29 bp between two groups of targets: Edible dormouse, Ring-tailed lemur, and Eastern long-fingered bat
178 targets produced an identical 16,642 bp contig, and all other references produced an identical 16,671 bp
179 contig. A combination of pairwise alignments and dot-plots (data not shown) indicate that these
180 differences are due to the way in which this circular sequence was linearized. The 16,642 bp contig has a
181 90 bp overlap between the beginning and end of the contig, while the 16,671 bp contig has a 119 bp
182 overlap, caused from group 2 recruiting one additional read relative to group 1. Therefore, even though
183 the assembled length differed slightly the resulting mitochondrial genomes were identical and equal in
184 length after trimming overlapping ends.

185 **ARC assembles large contigs from short, poor quality reads produced from ancient DNA**

186 Methods that permit investigators to extract DNA from samples that are as much as 50,000 years
187 old and prepare libraries for HTS have been developed (Gilbert et al. 2007, 2008; Knapp and Hofreiter
188 2010). The DNA from these ancient samples tends to be partially degraded resulting in shorter, poorer
189 quality reads (Knapp and Hofreiter, 2010). As described previously, ARC relies on an iterative process to
190 extend assemblies into gaps. Recruiting reads with partial, overhanging alignments at the edge of a contig
191 eventually fills these gaps. To test the effectiveness of ARC with short, single-end reads produced from
192 ancient samples, we used ARC to assemble the mammoth (*Mammuthus primigenius*) mitochondrial
193 genome using reads sequenced by Gilbert et al. (2007) from DNA collected from hair samples.

194 Sequenced reads were obtained for *Mammuthus primigenius* specimen M1 from the Sequence
195 Read Archive (SRA001810) and preprocessed as described in the Methods section. ARC was run using
196 three mitochondrial references target sequences: the published sequence from *Mammuthus primigenius*
197 specimen M13, Asian elephant (*Elephas maximus*) the closest extant relative of the mammoth (Gilbert et

198 al. 2008), and a more divergent reference, the house mouse (*Mus musculus*) (accessions: EU153445,
199 AJ428946, NC_005089 respectively).

200 We evaluated ARC results by alignment to the published *Mammuthus primigenius* M1 sequence
201 (EU153444), which is 16,458 bp in length. Results of this comparison are presented in Table 1. Percent
202 coverage (> 99%) and identity (> 98%) is high for both the mammoth and elephant references. The mouse
203 reference resulted in a slightly smaller assembly (total length 15,781 bp), however coverage (95.9%) and
204 identity (99.4%) were still high. Not surprisingly, the mouse reference required 78 ARC iterations to
205 build its final set of contigs, recruiting only 223 reads on the first iteration. Despite starting from such a
206 small number of initial reads, the final iteration recruited 4,507 reads, almost the same number as the
207 other reference sequences, but from a significantly more divergent reference sequence.

208 All contigs assembled by ARC could be aligned to the published reference sequence, however the
209 lengths of contigs assembled using the mammoth (16,620 bp) and elephant references (16,603 bp) were
210 both longer than the published sequence length (16,458 bp). To investigate whether this was due to a poor
211 quality assembly on the part of ARC, or an error in the published sequence, we aligned the ARC contigs
212 produced from the mammoth reference (*Mammuthus primigenius* M13) and the published *Mammuthus*
213 *primigenius* M1 sequence against the published Asian elephant sequence (Supplemental Figure 2). The
214 alignment showed a number of gaps existed in the ARC assembly as compared to the published contigs.
215 Each of these gaps was associated with a homopolymer (consecutive identical bases, e.g., AAA), a known
216 issue with Roche 454 pyrosequencing technology. More interesting was that the D-loop region of the
217 published *Mammuthus primigenius* M1 sequence contains 10 'N' characters followed by a 370bp gap
218 when aligned against the Asian elephant reference. ARC assembled 220 bp of this sequence, including
219 sequence that crosses the unknown, "N" bases in the published sequence. These assembled bases align
220 with high identity against the Asian elephant reference, suggesting that they represent an accurate
221 assembly of this locus and that the published M1 mitochondrial sequence is either missing sequence or is
222 misassembled in this region.

223 **ARC computational requirements for large datasets**

224 To be useful for modern genomic experiments ARC must be able to process large datasets with
225 multiple samples and potentially thousands of targets. We benchmarked ARC's performance with the
226 previously described chipmunk exome capture dataset that contains reads from 55 specimens and exonic
227 sequence captured from 7,627 genes as well as the full mitochondrial genome. After stringent read
228 cleaning to remove adapters, PCR duplicates, and overlapping of paired-end reads with short inserts, this
229 dataset contains 21.9 Gbp in 194,597,935 reads. For comparison purposes, we also carried out de novo
230 assemblies of three libraries using the Roche Newbler v2.6 assembler (Table 2).

231 ARC required 77 hours 45 minutes to process all 55 samples and 7,627 genes, carrying out a total
232 of 1.3 million assemblies and using a maximum of 31.19 GB of memory. On average this equates to 1
233 hour 25 minutes per sample. By comparison, individual whole dataset assemblies for a representative
234 three samples were variable, requiring between 6.71 GB and 17.54 GB of memory, with running times of
235 between 31 minutes and 13 hours 27 minutes to complete using Roche Newbler. Although time and
236 memory requirements are smaller for assembly of an individual sample, the total time required to
237 assemble 55 samples in serial would have been be much greater than the time required by ARC to process
238 all samples on a single machine. Likewise, the total memory usage needed to assemble all 55 samples
239 concurrently on a single machine would exceed the memory usage required by ARC on the same machine.
240 Further, assembly algorithms produce anonymous (unannotated) contigs, requiring significant additional
241 processing and analysis before homologous sequences are identified and can be compared between
242 samples. In contrast ARC contigs are annotated to the target from which they were initiated, facilitating
243 across sample comparisons.

244 Since ARC breaks a large assembly problem into smaller, more manageable pieces, we postulated
245 that memory requirements would scale as a function of the number of CPUs used to perform ARC
246 assemblies rather than as a function of the total number of reads as is normally the case with sequence
247 assembly (Li 2012). To test this, we performed nine ARC runs using between 10 and 50 CPU cores with

248 the 55-specimen chipmunk dataset. We used a random subset of 200 targets instead of the full 7,627
249 targets so that the experiment could be completed in a reasonable amount of time. During each assembly
250 we recorded maximum memory usage. The results indicate a linear increase in memory usage as the
251 number of cores increases (see Figure 4). A linear model was fit to this data resulting in an estimated
252 slope of 0.07 GB per CPU core ($P < .005$, $R^2 = 0.96$) for this dataset. It is important to note that even
253 though this dataset contains 21.9 Gbp of reads, analysis using a small number of CPU cores and a reduced
254 dataset required less than 3 GB of RAM total, making it possible to use ARC on any size dataset with any
255 modern desktop computer.

256 **DISCUSSION**

257 In this paper we introduce ARC, a software package that facilitates targeted assembly of HTS
258 data. This software is designed for use in situations where assembly of one or several discreet and
259 relatively small targets is needed and (potentially divergent) homologous reference sequences are
260 available for seeding these assemblies. ARC fills the gap between fast, mapping based strategies which
261 can fail to map, or misalign reads at divergent loci, and de novo assembly strategies which can be slow,
262 resource intensive, and require significant additional processing after assembly is complete. ARC was
263 evaluated in three ways: 1) we determined whether ARC results were biased by divergence of the
264 reference; 2) we tested the effectiveness of ARC to produce assemblies using short, low quality reads
265 produced from ancient DNA; and 3) we characterized performance on a large HTS dataset with 55
266 samples and thousands of subgenomic targets.

267 Assemblies using a divergent set of references with chipmunk specimens show that ARC does not
268 require a close reference to produce high quality final contigs. Supplemental Figure 1 illustrate that on the
269 initial iteration, ARC is able to map only a tiny fraction of the mitochondrial reads to all but the most
270 closely related gray-footed chipmunk reference, yet is able to recover, in most cases, a full set of reads
271 and complete mitochondrial genomes by iteration 50. This small set of reads represents the total number
272 of reads that would have been aligned using a traditional mapping strategy and further illustrates how

273 sensitive read mapping is to high levels of divergence. A similar pattern emerged when we used a mouse
274 reference to seed assembly of a mammoth mitochondrial genome. A mere 223 reads mapped on the first
275 iteration, which was sufficient to seed assembly of an almost full-length mitochondrial sequence
276 assembled from 4,507 reads.

277 Repetitive sequences and excess coverage are well-known issues, which increase memory usage
278 and slow assembly (Li 2012; Miller et al. 2010). Although ARC partially addresses this problem by
279 breaking the full set of reads into smaller subsets before assembly, it can still encounter issues with very
280 high coverage libraries, or when a target includes repetitive sequence and recruits a large numbers of
281 similar reads. For example, when testing ARC's ability to handle diverse mitochondrial references,
282 assemblies did not complete for specimen S10 using any of the 11 reference target sequences. In this case
283 the sequence depth was ~1500x for the mitochondrial genome; this depth is not suited for the Newbler
284 assembler, which performs pairwise comparisons of every read and works best when coverage is closer to
285 an expected depth of 60x. The excess coverage led to long assembly times and an eventual timeout.
286 Although the iterative ARC process did not run to completion in this case, intermediate contigs are still
287 reported and contained the full, although fragmented, mitochondrial genome.

288 ARC has a number of built in mechanisms to mitigate problems caused by repetitive sequences
289 and excess coverage. These include a masking algorithm that inhibits recruitment of reads from simple
290 tandem repeats, as well as tracking of read recruitment patterns that quits assembly if an unexpectedly
291 large number of reads are recruited between iterations, and an assembly timeout parameter that terminates
292 assemblies that run beyond a specified limit. In addition to these strategies there is also an option to
293 subsample reads in cases of known very high sequence depth. Subsampling was not used in any of the
294 tests described here, but may have improved results for samples such as S10. During testing and
295 development, we observed improved behavior with each of these measures on large datasets while
296 minimizing the impact of excess sequencing coverage and repeat elements. Implementing them has
297 allowed ARC to run more quickly and efficiently; however, it is clear that in some cases, recruitment of

298 excess reads and repeat elements can still cause problems for some targets or samples. In all completed
299 assemblies, the resulting set of reads and contigs were either identical or nearly so, providing strong
300 evidence that ARC is able to assemble high quality, unbiased contigs using even very divergent
301 references. This capability makes ARC a very useful tool when analyzing sequence data from non-model
302 organisms or when the identity of a sample is in question.

303 We tested ARC's ability to assemble contigs with short, low quality reads recovered from ancient
304 mammoth DNA and found that read length and quality did not impact ARC's ability to assemble full
305 length genomes. The resulting mitochondrial genome assemblies appear to be as good as or even better
306 than the published assembly for this sample despite using a divergent reference for ARC. Assembly of the
307 M1 mammoth sequence by Gilbert et al. (2007) was achieved through mapping against another mammoth
308 mitochondrial sequence published by Krause et al. (2006) that was generated using a laborious PCR-
309 based strategy. Because ancient DNA sequencing projects are often targeted at extinct organisms (Knapp
310 and Hofreiter 2010) there is rarely a high quality reference from the same species that can be aligned and
311 mapped to. This makes ARC an excellent choice for this type of data, where a target sequence from a
312 related, extant organism is likely to successfully seed assembly. Even in the case where no closely related
313 organism exists, a more distance reference may still be available, as was demonstrated by the assembly of
314 two large contigs representing ~96% of the mammoth mitochondrial genome using a mouse
315 mitochondrial genome for a reference. Additionally, ARC can be configured to use multiple reference
316 sequences as a single target. In cases where specimens cannot be identified, the user can select a set of
317 potentially homologous targets from many phylogenetically diverse taxa so that all sequences may serve
318 as references in order to seed assembly.

319 Analysis of HTS data can be computationally intensive, and time and memory requirements can
320 become serious limitations, especially with larger datasets (Zhang et al. 2011). With ARC, we have
321 attempted to reduce these requirements using a 'divide and conquer' approach that breaks large HTS
322 datasets up into many smaller problems, each of which can be solved quickly and with reduced resources.

323 In the large, 55 sample, 7,627 target dataset, ARC completed over 1.3 million assemblies, averaging
324 seven assemblies per second, in less than 78 hours. This approach allows the user to control memory
325 usage simply by changing the number of CPU cores available to ARC as shown in Figure 4. Less than 3
326 Gb of RAM was required when using 10 cores, despite processing a 21.9 Gbp dataset that would have
327 required many times this amount of memory using traditional assembly methods. Of course, using fewer
328 CPUs comes with the cost of a longer run time, so ARC can be tuned to the resources available.

329 It is useful to think of the DNA sequence mapping problem as a trade-off between sensitivity and
330 specificity (Fonseca et al. 2012). To avoid mapping reads to multiple loci throughout the reference,
331 mapping parameters must be tuned for high specificity. However, when divergent loci exist within the
332 reference sequence, high specificity limits the sensitivity of the mapper, leaving reads unmapped.
333 Assembly, on the other hand, can be seen as mapping reads against themselves, thereby removing
334 difficulties associated with divergent reference loci, but incurring the burden of pairwise read
335 comparisons that is significant in large datasets. ARC circumvents these problems by removing reference
336 bias through an iterative mapping and assembly process. As the intermediate reference is improved, more
337 reads can be recruited without sacrificing specificity, allowing both specificity and sensitivity to remain
338 high. At the same time, because only a small subset of reads is assembled, the all-by-all comparisons are
339 less burdensome. This process is carried out in an automated, easily configured manner, with standardized
340 output that simplifies additional analysis, or integration into existing sequence analysis pipelines.

341 **METHODS**

342 The ARC algorithm proceeds through a number of stages, which have been outlined below and
343 are presented in Figure 1. This algorithm consists of four steps: mapping, splitting, assembling, and
344 finishing. A graphical representation of the algorithm is presented in Figure 1, while an example
345 illustrating the ARC process from the perspective of reads and contigs is provided in Figure 2.

346 **Initialization**

347 During the initialization stage a configuration file is processed and a number of checks are carried
348 out to ensure that data and prerequisite applications specified in the configuration file are available. If any
349 checks fail, ARC will report an informative error message providing details about the problem and then
350 exit. If all checks pass successfully the initialization process continues by creating internal data structures
351 to store information about the experiment and pipeline progress. Working directories and read index files
352 are created for each sample, and names that are file-system safe are assigned to each reference target
353 sequence. Finally, the job manager is started (including job queues and workers), and read recruitment
354 jobs are added to the job queue for each sample. With initialization complete, ARC begins the iterative
355 part of the pipeline.

356 **Read recruitment: reads are recruited by mapping against a set of reference target sequences**

357 In the first iterative stage, ARC recruits reads by mapping them against a set of reference targets
358 using one of the two currently supported mappers, BLAT (Kent, 2002) or Bowtie 2 (Langmead and
359 Salzberg, 2012), which is specified in the configuration file. In all subsequent iterations, the reference
360 targets consist of contigs assembled from the previous iteration and are therefore highly similar and no
361 longer represent a divergent reference sequence since they were derived from the sample reads.

362 BLAT is a fast, seed-and-extend sequence alignment tool that supports gapped alignments and
363 has proven effective at recruiting reads even in cases where global sequence identity is as low as 70%. In
364 the first iteration, BLAT is run using default parameters (minIdentity=90, minScore=30) but on all
365 subsequent iterations mapping stringency is increased (minIdentity=98, minScore=40) to reduce
366 recruitment of less similar reads. BLAT reports all alignments that meet the minimum score criteria, so it
367 is possible to use the same read multiple times if it aligns successfully against more than one target. One
368 drawback of using BLAT is that it does not support the FASTQ format. All current HTS platforms
369 produce base quality information for reads and this information is typically encoded in FASTQ format.
370 To facilitate usage of ARC and FASTQ formatted data we include a code patch for BLAT that adds

371 support for FASTQ files. Instructions for applying this patch can be found in the online manual
372 (<http://ibest.github.io/ARC/>).

373 Bowtie 2 is another fast, gapped, read aligner that was specifically designed for mapping HTS
374 reads (Langmead and Salzberg, 2012). Bowtie 2 is ran in ARC under local alignment mode (--local
375 option) which enables the recruitment of reads that partially map to the ends of contigs and in low-
376 homology regions. Additionally, the option to report up to five valid alignments (-k 5) is used by default.
377 This setting can be modified based on the user's expectations by setting the bowtie2_k parameter in the
378 ARC configuration file. Setting bowtie2_k=1 will cause Bowtie 2 to run in default local-alignment mode
379 where only the best alignment found is reported.

380 **Split reads into bins: reads are split into subsets based on mapping results**

381 In the second iterative stage, ARC splits reads into bins based on the mapping results. The
382 supported mappers, BLAT and Bowtie 2 generate PSL or SAM (Li et al. 2009) formatted output files,
383 respectively. ARC processes each sample's mapping output file and reads are split by reference target.
384 This is accomplished by creating a series of FASTQ files corresponding to reads which map to each
385 reference target; allowing for the assembly of each target's reads independently from the others. Splitting
386 requires fast random access to the read files, which is facilitated by storing read offset values in a SQLite
387 database as implemented in the Biopython SeqIO module (Cock et al. 2009). Two special considerations
388 are taken into account during splitting. First, since the Newbler assembler uses pre CASAVA 1.8 Illumina
389 read identifiers to associate paired reads, it is necessary to reformat the read identifier to ensure
390 compatibility with Newbler paired-end detection. This is performed by ensuring that the read identifier is
391 made up of five fields separated by a colon and ending in a sixth field indicating the pair number, a
392 format compatible with most modern day assemblers. Identifiers for single-end reads are similarly
393 reformatted, except that the sixth field, which indicates pair number, is left blank. Secondly, regardless of
394 whether one or both of a read pair map to a target, both members of the pair are recruited as long as at

395 least one of them was mapped to the target sequence. Recruiting paired reads in this way takes advantage
396 of the information stored in paired reads, and allows for faster extension of targets.

397 Despite using a fast strategy for random accessing of read files, splitting is limited by system
398 input/output latency and to a single CPU core per sample. To optimize CPU use on modern multi-core
399 systems, ARC immediately adds an assembly job to the job queue as soon as all reads associated with a
400 target have been split. This allows assemblies to proceed concurrently with the read splitting process.

401 **Assemble each bin: targets are assembled using either the Spades or Newbler assemblers**

402 Because the read splitting process is carried out sequentially across mapping reference targets, an
403 assembly job for a target can be launched as soon as all reads associated with the target have been written.
404 As soon as resources become available, assembly jobs are started, allowing ARC to run read splitting and
405 assembly processes concurrently. Two assemblers are currently supported, the Roche GS *de novo*
406 Assembler (also known as Newbler; Margulies et al., 2005), and SPAdes (Bankevich et al., 2012).
407 Assemblies within ARC are always run with a timeout in order to gracefully handle the cases where the
408 assembler crashes, does not exit properly, or takes longer than expected to run. This allows ARC to
409 continue running efficiently on large projects where a small number of targets might be problematic (e.g.,
410 due to recruiting reads from repetitive elements). The timeout value can be controlled using the assembly
411 timeout setting in the configuration file.

412 Newbler was originally designed to assemble reads generated from the Roche 454
413 pyrosequencing platform (Margulies et al. 2005), but recent versions have added support for Illumina
414 paired-end reads and Newbler can be run using only Illumina reads. The ARC configuration file supports
415 two Newbler specific parameters that can sometimes improve assembly performance. These are to set
416 `urt=True`, which instructs Newbler to “use read tips” in assemblies, and `rip=True`, which instructs
417 Newbler to place reads in only one contig and to not to break and assign reads across multiple contigs.
418 We have found that setting `urt=True` can reduce the number of ARC iterations necessary to assemble a
419 target.

420 The second assembler supported in ARC is SPAdes (Bankevich et al., 2012). SPAdes is an easy
421 to use de Bruijn graph assembler that performed well in a recent evaluation of bacterial genome
422 assemblers (Magoc et al., 2013). SPAdes performs well in the ARC pipeline, but is not as fast as Newbler
423 for small target read sets (data not shown). This may partly be because SPAdes implements a number of
424 steps in an attempt at improving the often-fragmented de Bruijn graph assembly results seen in large
425 eukaryotic genomes. These steps include: read error correction, multiple assemblies using different k-mer
426 sizes, and merging of these assemblies. In ARC, SPAdes is run using the default set of parameters.

427 In some cases, the reference targets may be very divergent from the sequenced specimen and,
428 therefore, only a small number of reads are recruited in the first iteration. If too few overlapping reads are
429 recruited, the assemblers have very little data to work with, and in the case of SPAdes, may fail to
430 assemble any contigs. In an attempt to address this specific situation, we provide a final pseudo-assembly
431 option that skips assembly on the first iteration and treats any recruited reads as contigs. These reads are
432 then used as mapping reference targets in the second iteration. This option can be enabled by setting
433 `map_against_reads=True` in the ARC configuration file. In some cases using reads as mapping targets
434 results in recruiting large numbers of reads from repeat regions, causing the assembly to timeout and fail.
435 For this reason we only recommend using this approach after testing ARC with standard settings.

436 **Finisher: assembled contigs are written as a new set of mapping targets or to finished output**

437 Once all assemblies are completed for a given sample, the final iterative stage in the ARC
438 pipeline is initiated. During this stage each target is evaluated; if stopping conditions are met, the contigs
439 are written to the final output file; and if not the contigs are written to a temporary file where they are
440 used as reference targets in the next iteration (see the section Folder structure: outputs and logging for
441 details). Stopping conditions within ARC are defined as follows: 1) iterations have reached their
442 maximum allowable number as defined by the `numcycle` parameter in the ARC configuration file; 2) no
443 additional reads have been recruited (i.e., delta read count between iterations is zero); 3) detection of an
444 assembly that was halted, or killed will result in no further attempts at assembling this target, and any

445 contigs produced on the previous iteration will be written to the output file; or 4) a sudden spike in read
446 counts. Occasionally a target will be flanked by repeated sequence in the genome that can cause a sudden
447 spike in the number of recruited reads. The max_incorporation parameter in the ARC configuration file
448 controls sensitivity to this situation and by default is triggered if five times the previous number of reads
449 are recruited.

450 During output, target contig identifiers are modified to reflect their sample, original reference
451 target, and contig number separated by the delimiter “_:_” (e.g. sample_:_original-reference-
452 target_:_contig). Contigs are also masked of simple tandem repeats in all but the final iteration, using an
453 approach that relies on frequency of trinucleotides in a sliding window. Repeats are masked by setting
454 them to lower case for Blat support, or by modifying the repeat sequence to the IUPAC character 'N' for
455 Bowtie 2 support. All target contigs in their final iteration are written to the final output file, and all
456 corresponding reads are written to the final read files, however their description field is modified to
457 reflect which reference target they are assigned to.

458 For any targets that remain unfinished (i.e., stopping conditions have not been met), those
459 reference targets are iterated using the newly assembled contigs as the next mapping reference targets.

460 **Description of input files**

461 Inputs to ARC consist of three types of files: a file containing reference target sequence(s), file(s)
462 containing sequence reads for each sample, and an ARC configuration file.

463 The reference target sequence(s) file contains the sequences that are to be used as mapping
464 references during the first iteration of ARC. This file must be in standard FASTA format and should have
465 informative, unique names. It is possible to use multiple reference sequences as a single target in cases
466 where a number of potentially homologous targets are available and it is not clear which of them is most
467 similar to the sequenced sample (e.g. in the case of ancient DNA extracted from unidentified bone
468 material). This can be accomplished by naming each reference target using ARC's internal identifier
469 naming scheme made of three parts separated by “_:_” (e.g., sample_:_reference-target_:_contig). During

470 read splitting, ARC will treat all sequences that have an identical value in “reference-target” as a single
471 target.

472 Sample sequence read files are represented with up to three sequence read files; two paired-end
473 (PE) files, and one single-end (SE) file. ARC will function with only one SE file, a PE set of files, or all
474 three files if provided. If multiple sets of reads are available for a single biological sample (i.e., from
475 different sequencing runs or technologies) they should be combined into the above described three read
476 files. All reads for all samples must be in the same format (i.e., FASTA or FASTQ) and this format needs
477 to be indicated using the format parameter in the ARC configuration file. It is highly recommended that
478 reads be preprocessed to remove adapter sequences and low quality bases prior to running ARC.
479 Removing PCR duplicate reads and merging paired-end reads has also been observed to produce higher
480 quality, less fragmented ARC assemblies, particularly with capture data (data not shown).

481 The ARC configuration file is a plain text file describing the various parameters that ARC will
482 use during assembly, mapping, and output stages and the sample(s) read data data paths. By default the
483 configuration file should be named ARC_config.txt, but any name can be used as long as the -c filename
484 switch is used. The configuration file is split into three parts, denoted by the first characters in the line.
485 Lines starting with the characters “##” are treated as comments and ignored, lines starting with “#” are
486 used to set ARC parameters, and lines that don't begin with “#” indicate sample read data. The one
487 exception to this rule is the sample read data column header line, which is the first line that doesn't begin
488 with “#”, and contains column names. This line is ignored by ARC, but is expected in the configuration
489 file. An example ARC configuration file is included in the “test_data” directory that comes with ARC. A
490 comprehensive list of configuration options are presented in the online manual
491 (<http://ibest.github.io/ARC/>).

492 **Folder Structure Outputs and Logging**

493 In order to minimize memory usage and interface with assembly and mapping applications, ARC
494 relies heavily on temporary files. These files are organized into subdirectories under the path from which

495 ARC is launched. During ARC processing a pair of folders is created for each sample. These folders have
496 the prefixes “working_” and “finished_”. Temporary files used during ARC processing are stored in the
497 “working_” folders while completed results and statistics are recorded in the “finished_” folders.

498 The “working_” directories contain the sample contigs assembled during each iteration in a set of
499 files with file names “I00N_contigs.fasta” (where “N” corresponds to the iteration) and the latest
500 assembly directory denoted by “t_0000N” (where “N” corresponds to the numeric index of the target).
501 These directories and files can be informative in determining why an assembly failed or for examining
502 assembly statistics of a particular sample and target in more depth. Additionally, these folders provide the
503 option of manually re-running an assembly with a different set of parameters than those chosen within
504 ARC. In addition to the per iteration contigs and latest assembly directories, the “working_” folders also
505 contain the sample read indexes, which can be reused when re-running ARC with new parameters, and
506 the latest mapping log report. The “working_” folders only contain temporary files used by ARC and can
507 be safely deleted after the ARC run.

508 The “finished_” directories contain the following files: contigs.fasta, mapping_stats.tsv,
509 target_summary_table.tsv, and final read files. The contigs.fasta file contains the final set of assembled
510 contigs for each target. Contigs are named according to the three part naming scheme previously
511 described (sample_:_original-reference-target_:_contig) in order to facilitate downstream comparisons
512 between samples. The mapping_stats.tsv and target_summary_table.tsv files are tab-separated values files
513 that store information on the number of reads mapped to each target at each iteration and per target final
514 summary statistics respectively. These files can be easily loaded into a spreadsheet, or statistical program
515 such as R to generate plots or for other downstream analysis. The final read files (PE1.fasta/PE1.fastq,
516 PE2.fasta/ PE2.fastq, and SE.fasta/ SE.fastq) contain all the reads that were mapped, and consequently
517 used during assembly, on the final ARC iteration. If only pair-end or single-end files were provided then
518 only reads of this type will be reported. These files will be formatted in the same manner as the original

519 input files (FASTA or FASTQ) and have modified description fields to indicate the sample and target to
520 which they were assigned.

521 **ARC post processing and contrib scripts**

522 ARC contains a number of add on scripts in the “contrib” folder of the application, for
523 downstream processing of assembled contigs and visualization of ARC results. These scripts include R
524 functions to profile and plot memory usage and to plot data from the run log. The contrib folder also
525 contains number of Python scripts for post-processing ARC contigs for use in downstream applications
526 such as phylogenomics. Two scripts in particular are “ARC_Add_Cigar_Strings.py” and
527 “ARC_Call_and_Inject_hets.py”. The first allows users to determine the order and orientation of ARC-
528 generated contigs relative to the original reference, using the program BLAT to align assembled contigs
529 against sequences from the original reference targets sequence file. The script then generates a CIGAR
530 string in standard SAM format to describe the alignment. In situations where the contig extends beyond
531 the 5’ or 3’ ends of the target sequence, those bases are described as soft-clipped. The order of the
532 CIGAR string depends on the orientation of the contig with respect to the target (as is the case with
533 similar programs such as Bowtie2). If the contig maps to the forward strand, the CIGAR string reports the
534 matches, insertions, deletions, and soft-clipped regions of the alignment in the 5’ to 3’ direction. In
535 contrast, if the contig maps to the reverse strand, the CIGAR string reports components of the alignment
536 in the 3’ to 5’ direction. The script generates an output file (in FASTA format) that includes the contig
537 sequence from the original ARC output file, the name of the contig, the name of the target sequence the
538 contig mapped to, the start and end positions of the contig relative to the target sequence, the contig’s
539 orientation (i.e., “+” or “-” depending on whether the contig mapped to the forward or reverse strand of
540 the target), and the CIGAR string. With this information the user can ascertain the order and orientation
541 of ARC-generated contigs with respect to the reference.

542 The second script, “ARC_Call_and_Inject_hets.py”, produces both a variant call formatted file
543 (VCF) per sample and a new contigs file with ambiguity bases at heterozygous loci. This script uses

544 Bowtie 2 to map the reads recruited for each target to their respective assembled contigs. GATK and
545 Picard Tools are then used to call heterozygous SNPs and output a VCF file for each sample. Finally, the
546 script encodes the heterozygous SNP calls using their respective IUPAC ambiguity code and ‘injects’
547 those bases into the original contig sequences producing a new contigs file containing heterozygous sites.

548 **Datasets used for testing**

549 We tested ARC with two datasets. The first dataset is made up of Illumina sequence reads from
550 two chipmunk (*Tamias* sp.) exome capture experiments. This combined dataset consists of sequence reads
551 from 55 specimens, 3 of which were sequenced as part of Bi et al. (2012) while the other 52 were
552 sequenced as part of a separate study (Sarver et al. in prep). The second dataset consists of Roche
553 FLX sequence reads from a whole-genome shotgun sequencing experiment using ancient DNA extracted
554 from a mammoth hair shaft sample (Gilbert et al. 2007).

555 The first chipmunk dataset was used to investigate ARC's sensitivity to divergent references as
556 well as its utility and performance with large datasets. For all 55 specimens, libraries were captured using
557 an Agilent SureSelect custom 1M-feature microarray capture platform that contains 13,000 capture
558 regions representing the mitochondrial genome and 9,716 genes (Bi et al. 2012). Libraries were then
559 sequenced on the Illumina HiSeq 2000 platform (100bp paired-end). The 55 chipmunks represent seven
560 different species within the genus *Tamias* with representatives of *T. canipes*: 5, *T. cinereicollis*: 9, *T.*
561 *dorsalis*: 12, *T. quadrivittatus*: 1, *T. rufus*: 5, and *T. umbrinus*: 10, collected and sequenced as part of
562 Sarver et al. (in prep) and *T. striatus*: 3 collected and sequenced by Bi et al. (2012).

563 Prior to ARC analysis, reads were preprocessed through a read cleaning pipeline consisting of the
564 following steps. PCR duplicates were first removed using a custom Python script. Sequences were then
565 cleaned to remove sequencing adapters and low quality bases using the software package Seqclean
566 (Zhbannikov et al. in prep, <https://bitbucket.org/izhbannikov/seqclean>). Finally, because paired-end
567 sequencing produces two reads sequenced from either end of a single template, it is often possible to
568 overlap these reads to form a single long read representing the template in its entirety. This overlapping

569 was carried out using the Flash software package (Magoc and Salzberg, 2011). Post-cleaning, the dataset
570 consisted of 21.9 Gbp (giga base pairs) in 194,597,935 reads.

571 ARC analysis for the first dataset was carried out using two different sets of references. To
572 determine how well ARC performs with divergent references, the mitochondrial genome of each
573 specimen was assembled against eleven different mammalian mitochondrial references (see Figure 3). We
574 also tested ARC's performance with a large number of targets by using a target set consisting of a
575 manually assembled *Tamias canipes* mitochondrial sequence plus 11,976 exon sequences comprising
576 7,627 genes. These sequences represent the unambiguous subset from the 9,716 genes that the capture
577 probes were originally designed against.

578 The second woolly mammoth dataset was used to test ARC's performance on shorter, poor
579 quality reads that are typical of ancient DNA sequencing projects. Total DNA was extracted from ancient
580 hair shafts and reads were sequenced on the Roche 454 FLX platform by (Gilbert et al. 2007). Although
581 these reads represent shotgun sequencing of both the nuclear and mitochondrial genomes, the authors
582 report a high concentration of mitochondria in hair shaft samples resulting in high levels of mitochondrial
583 reads relative to nuclear reads. Sequenced reads for *Mammuthus primigenius* specimen M1 were
584 obtained from the Short Read Archive using accession SRX001889 and cleaned with SeqyClean
585 (Zhbannikov et al. in prep, <https://bitbucket.org/izhbannikov/seqyclean>) to remove 454 sequencing
586 adapters and low quality bases. Following cleaning, this datasets contains a total of 19 Mbp in 221,688
587 reads with an average length of 86.2 bp. Although these reads were sequenced on the Roche 454 platform
588 which typically produces much longer reads (400-700bp), 75% of cleaned reads were 101bp or less in
589 length making them extremely short for this platform. ARC analysis was carried out using three
590 mitochondrial references, the published *Mammuthus primigenius* sequence from another specimen, M13,
591 Asian elephant (*Elephas maximus*) the closest extant relative of the mammoth (Gilbert et al. 2008), and a
592 divergent reference, mouse (*Mus musculus*) (accessions: EU153445, AJ428946, NC_005089 respectively).

593 **DATA ACCESS**

594 The raw data used in this study are available in NCBI Sequence Read Archives under BioProject numbers
595 SRX001889 (*Mammuthus primigenius* M1), SRA053502 (Tamias samples S10, S11, S12), and
596 SRXXXXXX (Remaining Tamias samples). Reference sequences used in this study are available in
597 NCBI Genbank under accession numbers: NC_000884.1 (guinea pig), NC_001892.1 (edible dormouse),
598 HM156679.1 (human), AJ421451.1, (ring-tailed lemur), NC_015841.1 (cape hare), KF440685.1 (eastern
599 long-fingured bat), NC_000891.1 (platypus), NC_018788.1 (tasmanian devil), NC_002369 (red squirrel),
600 NC_005089 (house mouse), EU153445 (*Mammuthus primigenius*), AJ428946 (*Elephas maxiumus*),
601 NC_005089 (*Mus musculus*).

602 **ACKNOWLEDGEMENTS**

603 We would like to thank Ilya Zhbannikov for generating the FASTQ patch to BLAT. We would
604 also like to thank Jeffery Good, John Demboski, Jack Sullivan, Dan Vanderpool, and Kayce Bell for aid
605 in collecting and sequencing chipmunk samples. Research reported in this publication was supported by
606 the National Institute Of General Medical Sciences of the National Institutes of Health under Award
607 Number P30 GM103324.

608 **DISCLOSURE DECLARATION**

609 The authors declare no competing financial interests.

610 **FIGURE LEGENDS**

611
612 **Figure 1.** An example of iteratively assembling homologous sequences using ARC. In iteration 1, a small
613 number of reads and unmapped pairs are recruited to the more highly conserved regions of the divergent
614 reference. These reads are assembled and the resulting contigs are used as mapping targets in the next
615 iteration. This process is iterated until no more reads are recruited. Mapped reads are indicated in yellow,
616 unmapped reads in orange. Paired reads are indicated with a connector. Both members of a pair are
617 recruited if only one maps.

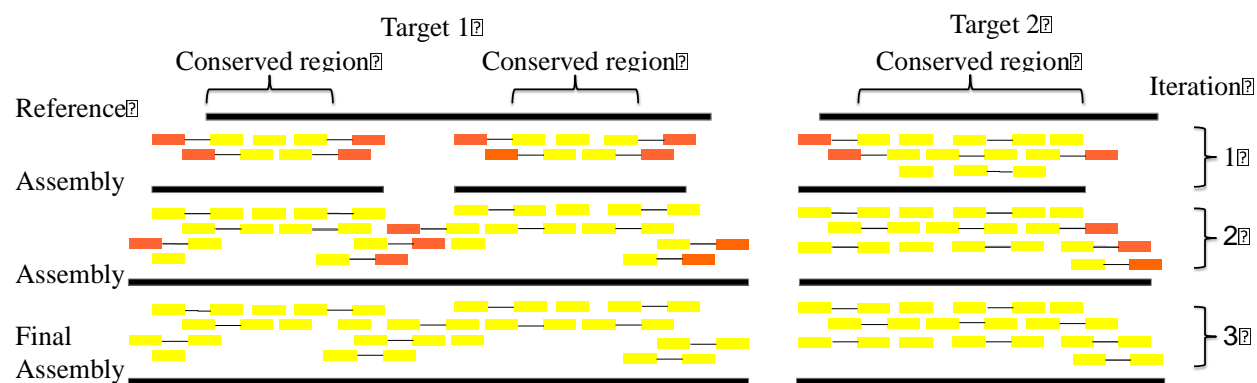
618 **Figure 2.** ARC processing stages. The ARC algorithm consists of an initialization stage, followed by four
619 steps: 1) read recruitment, 2) split reads into bins, 3) assemble each bin and 4) finisher. These steps are
620 iterated until stopping conditions are met, at which point a final set of contigs and statistics are produced.

621 **Figure 3.** Set of references used for ARC assembly of chipmunk mitochondrial genomes and their
622 respective scientific names, genome sizes, and NCBI Genbank accession numbers. Percent identity is
623 determined with respect to the Gray-Collared chipmunk (*Tamias cinereicollis*). Boxplots show the
624 variation around the number of ARC iterations for each reference species across all 55 samples, before
625 stopping conditions were met.

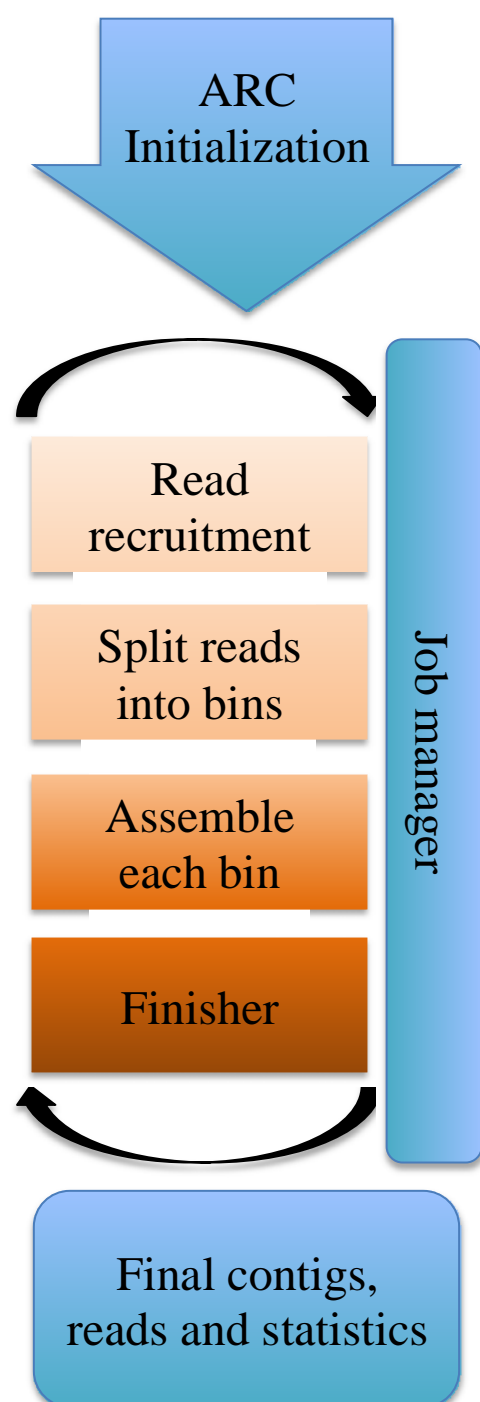
626 **Figure 4.** ARC memory requirements (y-axis) scale as function of the number of CPU cores used (x-axis).
627 A line of best fit is plotted in red.

628
629

630 **FIGURES**
631

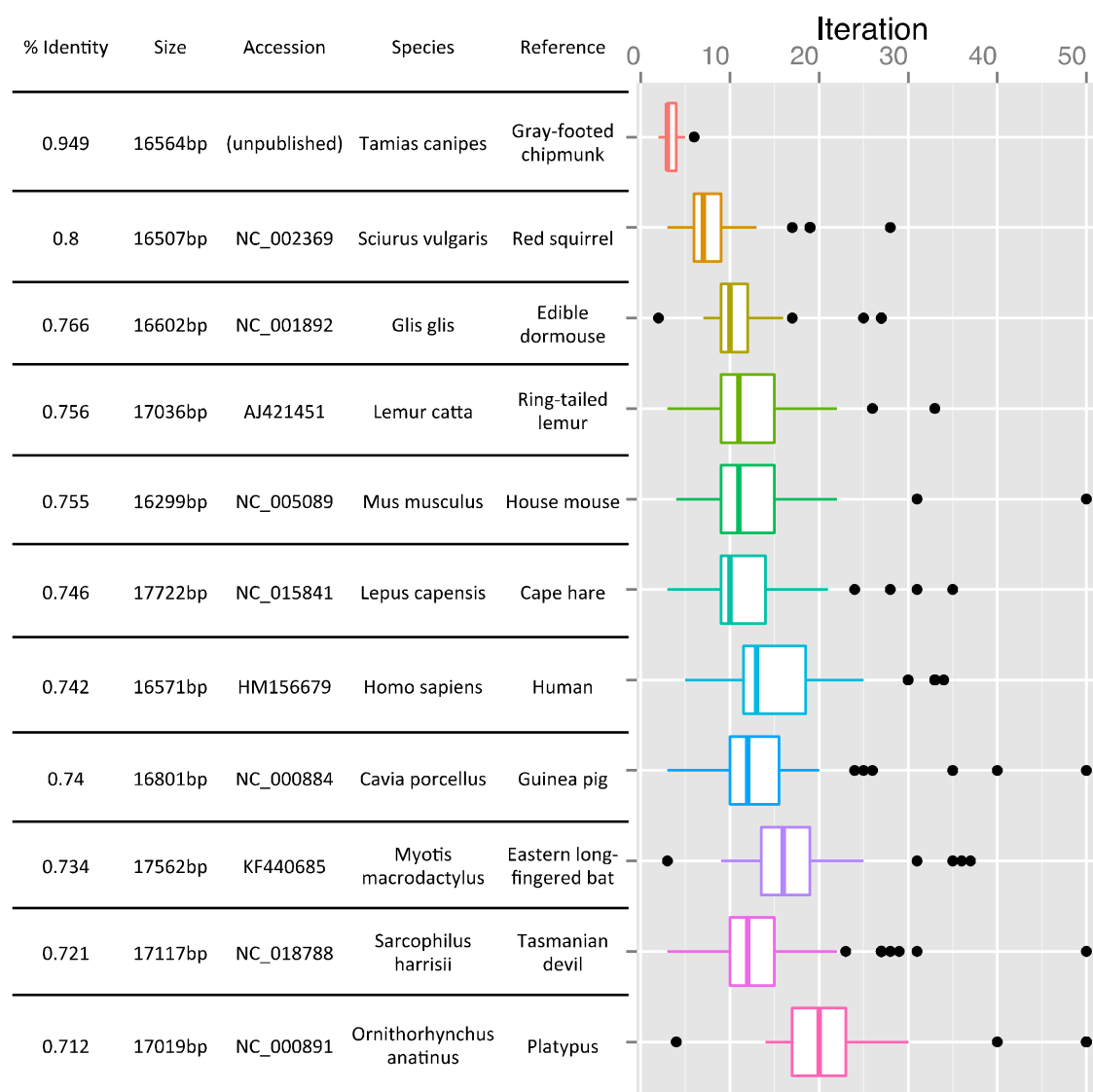


632 **Figure 1**
633
634

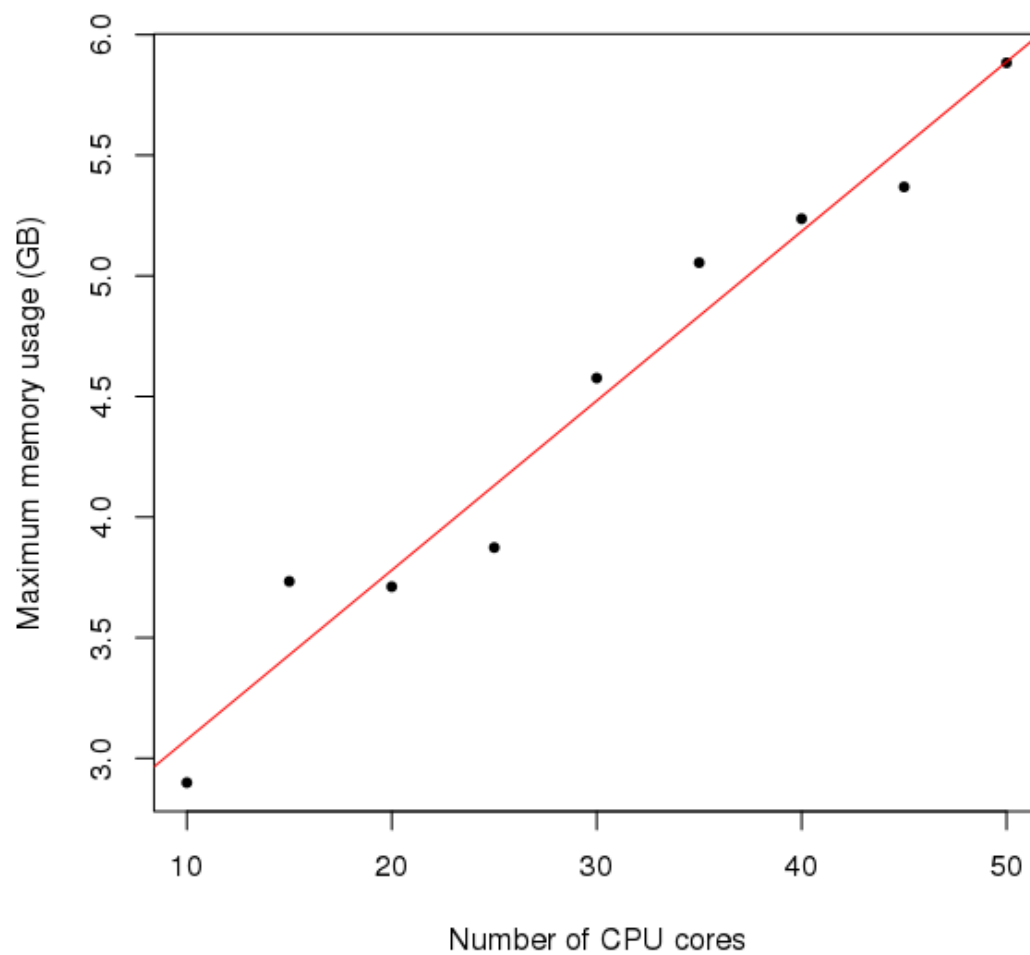


635
636
637

Figure 2



638
639 **Figure 3**
640



641
642 **Figure 4.**
643

644

645 **TABLES**

646

647 **Table 1** ARC results for assembly of ancient mammoth DNA sequences. ARC produces a small number

648 of contigs in all cases with good coverage and identity between the assembled contigs and published

649 reference.

Reference	Contig count	Total contig length (bp)	Percent coverage	Percent identity	ARC iteration	Reads recruited
<i>Mammuthus primigenius</i>	4	16,620	99.7%	98.1%	3	4633
<i>Elephas maximus</i>	4	16,603	99.7%	98.2%	5	4631
<i>Mus musculus</i>	2	15,781	95.9%	99.4%	78	4507

650

651 **Table 2** ARC assembly of 55 specimens compared to individual Roche Newbler *de novo* assemblies of
652 three specimens (S151, S152, and S223). Maximum and average memory usage (RAM) is listed in
653 gigabytes (GB). Total data processed is reported in millions of base pairs (Mbp).

	ARC	Newbler: S151	Newbler: S152	Newbler: S223
Total running time	77hr, 45min	31 min	1hr 13min	13hr 27min
Average Memory (GB)	22.78	5.847	8.337	16.36
Maximum Memory (GB)	31.19	6.71	9.967	17.54
Total assemblies performed	1,300,076	Not Applicable	Not Applicable	Not Applicable
Average assemblies per second	7.03	Not Applicable	Not Applicable	Not Applicable
Library Size (Mbp)	21,913	243	367	629

654

655 **REFERENCES**

- 656 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V. M.,
657 Nikolenko, S. I., Pham, S., Prjibelski, A. D., et al. (2012). SPAdes: A New Genome Assembly
658 Algorithm and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology*,
659 *19*(5), 455–477. doi:10.1089/cmb.2012.0021
- 660 Bi, K., Vanderpool, D., Singhal, S., Linderoth, T., Moritz, C., Good, J. M. (2012). Transcriptome-based
661 exon capture enables highly cost-effective comparative genomic data collection at moderate
662 evolutionary scales. *BMC Genomics*, *13*(1), 403. doi:10.1186/1471-2164-13-403
- 663 Bradnam, K. R., Fass, J. N., Alexandrov, A., Baranay, P., Bechner, M., Birol, I., Boisvert, S., Chapman,
664 J. A., Chapuis, G., Chikhi, R., et al. (2013). Assemblathon 2: evaluating de novo methods of genome
665 assembly in three vertebrate species. *GigaScience*, *2*(1), 10. doi:10.1186/2047-217X-2-10
- 666 Chevreux, B., Wetter, T., Suhai, S. (1999). Genome Sequence Assembly Using Trace Signals and
667 Additional Sequence Information. *Computer Science and Biology: Proceedings of the German
668 Conference on Bioinformatics (GCB)*, 45–56.
- 669 Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck,
670 T., Kauff, F., Wilczynski, B., et al. (2009). Biopython: Freely available Python tools for
671 computational molecular biology and bioinformatics. *Bioinformatics*, *25*(11), 1422–1423.
- 672 Fonseca, N. A., Rung, J., Brazma, A., Marioni, J. C. (2012). Tools for mapping high-throughput
673 sequencing data. *Bioinformatics*.
- 674 Gilbert, M. T. P., Tomsho, L. P., Rendulic, S., Packard, M., Drautz, D. I., Sher, A., Tikhonov, A., Dalén,
675 L., Kuznetsova, T., Kosintsev, P., et al. (2007). Whole-genome shotgun sequencing of mitochondria
676 from ancient hair shafts. *Science (New York, N.Y.)*, *317*(5846), 1927–1930.
- 677 Gilbert, M. T. P., Drautz, D. I., Lesk, A. M., Ho, S. Y. W., Qi, J., Ratan, A., Hsu, C., Sher, A., Dalén, L.,
678 Götherström, A., et al. (2008). Intraspecific phylogenetic analysis of Siberian woolly mammoths
679 using complete mitochondrial genomes. *Proceedings of the National Academy of Sciences of the
680 United States of America*, *105*(24), 8327–8332.
- 681 Hahn, C., Bachmann, L., Chevreux, B. (2013). Reconstructing mitochondrial genomes directly from
682 genomic next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic Acids
683 Research*, *41*(13), e129. doi:10.1093/nar/gkt371
- 684 Kent, W. J. (2002). BLAT - The BLAST-like alignment tool. *Genome Research*, *12*(4), 656–664.
- 685 Knapp, M., Hofreiter, M. (2010). Next Generation Sequencing of Ancient DNA: Requirements, Strategies
686 and Perspectives. *Genes*, *1*(2), 227–243. doi:10.3390/genes1020227
- 687 Krause, J., Dear, P. H., Pollack, J. L., Slatkin, M., Spriggs, H., Barnes, I., Lister, A. M., Ebersberger, I.,
688 Pääbo, S., Hofreiter, M. (2006). Multiplex amplification of the mammoth mitochondrial genome and
689 the evolution of Elephantidae. *Nature*, *439*(7077), 724–727. doi:10.1038/nature04432
- 690 Langmead, B., Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*, *9*(4),
691 357–359. doi:10.1038/nmeth.1923
- 692 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R.
693 (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*(16), 2078–
694 2079. Li, H. (2011). Improving SNP discovery by base alignment quality. *Bioinformatics*, *27*(8),
695 1157–1158.
- 696 Li, H. (2012). Exploring single-sample snp and indel calling with whole-genome de novo assembly.
697 *Bioinformatics*, *28*(14), 1838–1844.
- 698 Magoc, T., Salzberg, S. L. (2011). FLASH: Fast length adjustment of short reads to improve genome
699 assemblies. *Bioinformatics*, *27*(21), 2957–2963.
- 700 Magoc, T., Pabinger, S., Canzar, S., Liu, X., Su, Q., Puiu, D., Tallon, L. J., Salzberg, S. L. (2013). GAGE-
701 B: An evaluation of genome assemblers for bacterial organisms. *Bioinformatics*, *29*(14), 1718–1725.
- 702 Malé, P. J. G., Bardon, L., Besnard, G., Coissac, E., Delsuc, F., Engel, J., Lhuillier, E., Scotti-Saintagne,
703 C., Tinaut, A., Chave, J. (2014). Genome skimming by shotgun sequencing helps resolve the
704 phylogeny of a pantropical tree family. *Molecular Ecology Resources*, *14*(5), 966–975.

- 705 Margulies, M., Egholm, M., Altman, W. E., Attiya, S., Bader, J. S., Bemben, L. A., Berka, J., Braverman,
706 M. S., Chen, Y., Chen, Z., et al. (2005). Genome sequencing in microfabricated high-density
707 picolitre reactors. *Nature*, 437(7057), 376–380.
- 708 Miller, J. R., Koren, S., Sutton, G. (2010). Assembly algorithms for next-generation sequencing data.
709 *Genomics*.
- 710 Picardi, E., Pesole, G. (2012). Mitochondrial genomes gleaned from human whole-exome sequencing.
711 *Nature Methods*, 9(6), 523–4. doi:10.1038/nmeth.2029
- 712 Pyrkosz, A. B., Cheng, H., Brown, C. T. (2013). RNA-Seq Mapping Errors When Using Incomplete
713 Reference Transcriptomes of Vertebrates. *arXiv:1303.2411*, 1–17. Retrieved from
714 <http://arxiv.org/abs/1303.2411>
- 715 Sarver
- 716 Schbath, S., Martin, V., Zytnicki, M., Fayolle, J., Loux, V., Gibrat, J. F. (2012). Mapping Reads on a
717 Genomic Sequence: An Algorithmic Overview and a Practical Comparative Analysis. *Journal of*
718 *Computational Biology*, 19(6), 796–813. doi:10.1089/cmb.2012.0022
- 719 The 1000 Genomes Project Consortium, Abecasis, G. R., Altshuler, D., Auton, A., Brooks, L. D., Durbin,
720 R. M., Gibbs, R. A., Hurles, M. E., McVean, G. A. (2010). A map of human genome variation from
721 population-scale sequencing. *Nature*, 467(7319), 1061–73. doi:10.1038/nature09534
- 722 Zhang, W., Chen, J., Yang, Y., Tang, Y., Shang, J., Shen, B. (2011). A practical comparison of De Novo
723 genome assembly software tools for next-generation sequencing technologies. *PLoS ONE*, 6(3).
- 724 Zhbannikov, I. Y., Hunter, S. S., Foster, J. A., Settles, M. L. (2014) SeqyClean: a pipeline for high
725 throughput sequence data preprocessing. *In Prep*.
- 726