

A peer-reviewed version of this preprint was published in PeerJ on 13 July 2018.

[View the peer-reviewed version](https://peerj.com/articles/5175) (peerj.com/articles/5175), which is the preferred citable publication unless you specifically need to cite this preprint.

Andermann T, Cano Á, Zizka A, Bacon C, Antonelli A. 2018. SECAPR—a bioinformatics pipeline for the rapid and user-friendly processing of targeted enriched Illumina sequences, from raw reads to alignments. PeerJ 6:e5175 <https://doi.org/10.7717/peerj.5175>

1 **Author Cover Page**

2 Article submission to PeerJ

3 Manuscript category: Bioinformatics Tool

4 Collection: “*Endless forms: Advances in evolutionary analyses of biodiversity*”

5 Article title: **SECAPR - A bioinformatics pipeline for the rapid and user-friendly**
6 **processing of Illumina sequences, from raw reads to alignments**

7

8 Authors: Tobias Andermann* ^(1,2), Ángela Cano ^(2,3), Alexander Zizka ^(1,2), Christine
9 Bacon ^(1,2), Alexandre Antonelli ^(1,2,4,5)

10

11 Affiliations:

12 ¹ Department of Biological and Environmental Sciences, University of Gothenburg, Box
13 461, SE 405 30, Göteborg, Sweden

14 ² Gothenburg Global Biodiversity Centre, Göteborg, 41319, Sweden

15 ³ Department of Botany and Plant Biology, University of Geneva, Geneva, Switzerland

16 ⁴ Gothenburg Botanical Garden, Göteborg, 41319, Sweden

17 ⁵ Department of Organismic and Evolutionary Biology, Harvard University, Cambridge,
18 MA 02138 USA

19

20 ***Corresponding author: Tobias Andermann, E-mail: tobias.hofmann@bioenv.gu.se**

21 **Abstract**

22 Evolutionary biology has entered an era of unprecedented amounts of DNA sequence
23 data, as new sequencing platforms such as Massive Parallel Sequencing (MPS) can
24 generate billions of nucleotides within less than a day. The current bottleneck is how to
25 efficiently handle, process, and analyze such large amounts of data in an automated and
26 reproducible way. To tackle these challenges we introduce the Sequence Capture
27 Processor (SECAPR) pipeline for processing raw sequencing data into multiple sequence
28 alignments for downstream phylogenetic and phylogeographic analyses. SECAPR is
29 user-friendly and we provide an exhaustive empirical data tutorial intended for users with
30 no prior experience with analyzing MPS output. SECAPR is particularly useful for the
31 processing of sequence capture (synonyms: target or hybrid enrichment) datasets for non-
32 model organisms, as we demonstrate using an empirical sequence capture dataset of the
33 palm genus *Geonoma* (Arecaceae). Various quality control and plotting functions help
34 the user to decide on the most suitable settings for even challenging datasets. SECAPR is
35 an easy-to-use, free, and versatile pipeline, aimed to enable efficient and reproducible
36 processing of MPS data for many samples in parallel.

37

38 **Keywords:** Next generation sequencing (NGS), exon capture, Illumina, FASTQ,
39 contig, allele phasing, phylogenetics, phylogeography, BAM, assembly

40

41 **Introduction**

42 An increasing number of studies apply sequence data generated by Massive Parallel
43 Sequencing (MPS) to answer phylogeographic and phylogenetic questions (e.g. Botero-
44 Castro et al. 2013; Smith et al. 2014; Faircloth et al. 2015; Heyduk et al. 2016).
45 Researchers often decide to selectively enrich and sequence specific genomic regions of
46 interest, rather than sequencing the complete genome. One reason is that enriching
47 specific markers leads to a higher sequencing depth for each individual marker, as
48 compared to the alternative whole genome sequencing. Sequencing depth is important for
49 the extraction of single nucleotide polymorphisms (SNPs) and for allele phasing
50 (Andermann et al. 2018; Bravo et al. 2018). Additionally, phylogenetic analysis software
51 usually relies on multiple sequence alignments (MSAs) with homologous sequences
52 across many taxa, which are easiest to recover when specifically enriching these
53 sequences across all samples prior to sequencing.

54 The enrichment of specific genomic regions (markers) is usually archived through
55 sequence capture (synonyms: hybrid enrichment, hybrid selection, exon capture, target
56 capture) prior to sequencing (Gnirke et al. 2009). This technique applies specific RNA
57 baits, which hybridize with the target regions and can be captured with magnetic beads.
58 Sequence capture is gaining popularity, as more bait sets for non-model organisms are
59 being developed. Some bait sets are designed to match one specific taxonomic group (e.g.
60 Heyduk et al. 2016; Kadlec et al. 2017), while others are designed to function as more
61 universal markers to capture homologous sequences across broad groups of taxa (e.g.
62 UCEs, Faircloth et al. 2012). After enrichment of targeted markers with such bait sets, the
63 enriched sequence libraries are sequenced on a MPS machine (see Reuter, Spacek, and
64 Snyder 2015).

65 Despite recent technological developments, analyzing sequencing results is a great
66 challenge due to the amount of data produced by MPS machines. An average dataset
67 often contains dozens to hundreds of samples, each with up to millions of sequencing
68 reads. Such amounts of sequence data require advanced bioinformatics skills for storing,
69 quality checking, and processing the data, which may represent an obstacle for many

70 students and researchers. This bottleneck calls for streamlined, integrative and user-
71 friendly pipeline solutions.

72 To tackle these challenges, here we introduce the Sequence Capture Processor (SECAPR)
73 pipeline, a semi-automated workflow to guide users from raw sequencing results to
74 cleaned and filtered multiple sequence alignments (MSAs) for phylogenetic and
75 phylogeographic analyses. We designed many of the functionalities of this pipeline
76 toward sequence capture datasets in particular, but it can be effectively applied to any
77 MPS dataset generated with Illumina sequencing (Illumina Inc., San Diego, CA, USA).
78 SECAPR comes with a detailed documentation in form of an empirical data tutorial,
79 which is explicitly written to guide users with no previous experience with MPS datasets.
80 To simplify the processing of big datasets, all available functions are built to process
81 batches of samples, rather than individual files. We developed SECAPR to provide the
82 maximum amount of automation, while at the same time allowing the user to choose
83 appropriate settings for their specific datasets. The pipeline provides several plotting and
84 quality-control functions, as well as more advanced processing options such as the
85 assembly of fully phased allele sequences for diploid organisms (Andermann et al. 2018).

86

87 **Material & Methods**

88 *The SECAPR pipeline in a nutshell*

89 SECAPR is a platform-independent pipeline written in python, and tested for full
90 functionality on Linux and MacOS. It can be easily downloaded together with all its
91 dependencies as a virtual environment, using the conda package manager (see
92 Availability). The strength of SECAPR is that it channels the main functionalities of
93 many commonly used bioinformatics programs and enables the user to apply these to sets
94 of samples, rather than having to apply different software to each sample individually.

95 The basic SECAPR workflow (Figure 1) includes the following steps:

- 96 1. *Quality filtering and adapter trimming*
- 97 2. *De novo contig assembly*
- 98 3. *Selection of target contigs*

99 4. *Building MSAs from contigs*

100 5. *Reference-based assembly*

101 6. *Allele phasing*

102 SECAPR automatically writes summary statistics for each processing step and sample to
103 a log-file (*summary_stats.txt*, Table 1). The pipeline includes multiple visualization
104 options (e.g. Figure 2 and Figure 4) to gauge data quality and, if necessary, adapt
105 processing settings accordingly. SECAPR comes with a detailed documentation and data
106 tutorial (see Availability).

107

108 *Description of the SECAPR workflow*

109 1. *Quality filtering and adapter trimming (secapr clean_reads)*. The SECAPR
110 *clean_reads* function applies the software Trimmomatic (Bolger, Lohse, and Usadel
111 2014) for removing adapter contamination and low quality sequences from the raw
112 sequencing reads (FASTQ-format). An additional SECAPR plotting function summarizes
113 FASTQC (Babraham Institute) quality reports of all files and produces a visual overview
114 of the whole dataset (Figure 2). This helps to gauge if the files are sufficiently cleaned or if
115 the *clean_reads* function should be rerun with different settings.

116 2. *De novo contig assembly (secapr assemble_reads)*. The SECAPR function
117 *assemble_reads* assembles overlapping FASTQ reads into longer sequences (*de novo*
118 contigs) by implementing the *de novo* assembly software Abyss (Simpson et al. 2009).
119 Abyss has been identified as one of the best-performing DNA sequence assemblers
120 currently available (Hunt et al. 2014). SECAPR produces one file for each sample
121 (FASTA-formatted) that contains all assembled contigs for that sample.

122 3. *Selection of target contigs (secapr find_target_contigs)*. The SECAPR function
123 *find_target_contigs* identifies and extracts those contigs that represent the DNA targets of
124 interest. This function implements the program LASTZ (formerly BLASTZ, Harris 2007)
125 by searching the contig files for matches with a user-provided FASTA-formatted
126 reference library. For sequence capture datasets, a suitable reference library is the
127 reference file that was used for synthesizing the RNA baits, which will return all contigs

128 that match the enriched loci of interest. The *find_target_contigs* function identifies
129 potentially paralogous loci (loci that have several matching contigs) and excludes these
130 from further processing. It further allows the user to keep or exclude long contigs that
131 match several adjacent reference loci, which can occur if the reference file contains
132 sequences that are located in close proximity to each other on the genome (e.g. several
133 separate exons of the same gene).

134 *4. Building MSAs from contigs (secapr align_sequences).* The SECAPR function
135 *align_sequences* builds multiple sequence alignments (MSAs) from the target contigs that
136 were identified in the previous step. The function builds a separate MSA for each locus
137 with matching contigs for ≥ 3 samples.

138 *5. Reference-based assembly (secapr reference_assembly).* The SECAPR
139 *reference_assembly* function applies the BWA mapper (Li and Durbin 2010) for
140 reference-based assembly of FASTQ reads and Picard (broadinstitute.github.io/picard/
141 for removing duplicate reads. The function saves the assembly results as BAM files
142 (Figure 3) and generates a consensus sequence from the read variation at each locus.
143 These consensus sequences have several advantages over the *de novo* contig sequences
144 (see Discussion) and can be used for building MSAs with the SECAPR *align_sequences*
145 function

146 The *reference_assembly* function includes different options for generating a reference
147 library for all loci of interest:

- 148 • *--reference_type alignment-consensus:* The user provides a link to a folder
149 containing MSAs, e.g. the folder with the contig MSAs from the previous step,
150 and the function calculates a consensus sequence from each alignment. These
151 consensus sequences are then used as the reference sequence for the assembly.
152 This function is recommended when running reference-based assembly for groups
153 of closely related samples (e.g. samples from the same genus or family).
- 154 • *--reference_type sample-specific:* From the MSAs, the function extracts the
155 contigs for each sample and uses them as a sample-specific reference library. If
156 the user decides to use this function it is recommendable to only use alignments

157 for reference that contain sequences for all samples. This will ensure that the same
158 loci are being assembled for all samples.

159 • *--reference_type user-ref-lib*: The user can provide a FASTA file containing a
160 custom reference library.

161 An additional SECAPR function (*locus_selection*) allows the user to select a subset of the
162 data consisting of only those loci, which have the best read-coverage across all samples
163 (Figure 4b).

164 6. *Allele phasing (secapr phase_alleles)*. The SECAPR *phase_alleles* function can be
165 used to sort out the two phases (reads covering different alleles) at a given locus. This
166 function applies the phasing algorithm as implemented in SAMtools (Li et al. 2009),
167 which uses read connectivity across multiple variable sites to determine the two phases of
168 any given diploid locus (He et al. 2010). After running the phasing algorithm, the
169 *phase_alleles* function outputs a separate BAM-file for each allele and generates
170 consensus sequences from these allele BAM-files. This results into two sequences at each
171 locus for each sample, all of which are collected in one cumulative sequence file
172 (FASTA). This sequence file can be run through the SECAPR *align_sequences* function
173 in order to produce MSAs of allele sequences.

174

175 *Benchmarking with empirical data*

176 We demonstrate the functionalities of SECAPR on a novel dataset of target sequencing
177 reads of *Geonoma*, one of the most species-rich palm genera of tropical Central and
178 South America. (Dransfield et al. 2008) (Henderson 2011). Our data comprised newly
179 generated Illumina sequence data for 17 samples of 14 *Geonoma* species (Supplementary
180 Table S1), enriched through sequence capture. The bait set for sequence capture was
181 designed specifically for palms by Heyduk et al. (2016) to target 176 genes with in total
182 837 exons. More detailed information about the generation of the sequence data can be
183 found in Appendix 1 (Supplemental Material). All settings and commands used during
184 processing of the sequence data can be found in the SECAPR documentation on our
185 GitHub page (see Availability).

186

187 **Results**

188 The newly generated *Geonoma* data used for benchmarking constitute an empirical
189 example of a challenging dataset, characterized by irregular read coverage and multiple
190 haplotypes. Despite these challenges, the SECAPR workflow provides the user all the
191 necessary functions to filter and process datasets into MSAs for downstream phylogenetic
192 analyses.

193 After *de novo* assembly (*secapr assemble_reads*) we recovered an average of 323
194 (stdev=14) contigs per sample (*secapr find_target_contigs*) that matched sequences of the
195 837 targeted exons (Table 1, Figure 4a, Supplementary Table S2). In total 45 exons were
196 recovered for all samples. Many of the recovered target contigs spanned several reference
197 exons (all samples: mean=100, stdev=25) and hence were flagged as contigs matching
198 multiple loci (Supplementary Table S3). Since these contigs may be phylogenetically
199 valuable, as they contain the highly variable interspersed introns, we decided to keep
200 these sequences. We extracted these longer contigs together with all other non-duplicated
201 contigs that matched the reference library (*secapr find_target_contigs*) and generated
202 MSAs for each locus that could be recovered in at least three *Geonoma* samples (*secapr*
203 *align_sequences*). This resulted in contig alignments for 593 exon loci (center line in
204 Figure 4a).

205 During reference-based assembly (*secapr reference_assembly*) we mapped the reads
206 against the consensus sequence of the contig MSAs for all loci. We found an average of
207 439 exon loci (stdev=82) per sample that were covered by more than three reads (average
208 coverage across complete locus, Figure 4a). Hence, our approach of mapping FASTQ
209 reads to libraries compiled from the data leads to an increase of recovered loci per
210 sample, from 323 resulting from *de novo* assembly to 439 from the referenced-based
211 assembly (36% increase). Further, the number of loci that were recovered with sufficient
212 coverage for all samples increased by 116%, from 45 after the *de novo* assembly, to 97
213 after the reference-based assembly (Supplementary Table S4). We extracted the 50 loci
214 with the best coverage across all samples (*secapr locus_selection*), as shown in Figure 4b.
215 In cases of irregular read-coverage across samples (as in our sample *Geonoma* data), we

216 strongly recommend the use of the *locus_selection* function before further processing the
217 data, as demonstrated in our tutorial (see Availability).

218 The results of the reference-based assembly also revealed that our sample data showed
219 more than two haplotypes for many loci. Future research may clarify whether this is the
220 result of various paralogous loci in the dataset or if it is the result of a recent genome
221 duplication or hybridization event in the ancestry of our *Geonoma* samples. Due to the
222 presence of more than two haplotypes at various loci, the results of the allele-phasing step
223 (*secapr phase_alleles*) are to be viewed critically, since the algorithm is built for phasing
224 the read data of diploid organisms or loci only. All phased BAM files and the compiled
225 allele MSAs are available online (see Availability).

226

227 **Discussion**

228 *De novo assembly vs. reference-based assembly*

229 There are several ways of generating full sequences from raw FASTQ-formatted
230 sequencing reads. The SECAPR pipeline contains two different approaches, namely *de*
231 *novo* assembly and reference-based assembly (Figure 1). *De novo* assembly can be
232 directly applied to any raw read data while reference-based assembly requires the user to
233 provide reference sequences for the assembly. We find for the *Geonoma* example data
234 that reference-based assembly results into recovering more target sequences per sample
235 (Figure 4) and provides the user a better handle on quality and coverage thresholds. It is
236 also computationally much less demanding in comparison to *de novo* assembly.

237 However, reference-based assembly is very sensitive toward the user providing
238 orthologous reference sequences that are similar enough to the sequencing reads of the
239 studied organisms. If the reference sequences are too divergent from the sequenced
240 organisms, only a small fraction of the existing orthologous sequencing reads will be
241 successfully assembled for each locus. In contrast, when relaxing similarity thresholds
242 and other mapping parameters too much (e.g. to increase the fraction of reads included in
243 the assembly) there is higher a risk of assembling non-orthologous reads, which can lead
244 to chimeric sequences being assembled. This can be a problem, particularly in cases of

245 datasets containing non-model organisms, since suitable reference sequences for all loci
246 usually do not exist.

247 For this reason, the SECAPR workflow encourages the user to use these two different
248 assembly approaches in concert (Figure 1). Our general suggestion is to first assemble
249 contig MSAs for all regions of interest, resulting from *de novo* assembly and then use
250 these MSAs to build a reference library for reference-based assembly. In that case
251 SECAPR produces a reference library from the sequencing data itself, which is specific
252 for the taxonomic group of interest or even for the individual sample.

253 A common approach is to stop data processing after the *de novo* assembly step and then
254 use the contig MSAs for phylogenetic analyses (e.g. Faircloth et al. 2012; B. T. Smith et
255 al. 2014; Faircloth 2015). Here we take additional processing steps, including generating
256 new reference libraries for all samples and using these for reference-based assembly.

257 There may be several reasons for carrying out these additional steps:

- 258 1. Sensitivity: In order to identify *de novo* contigs that are orthologous to the loci of
259 interest, the user is usually forced (because of the lack of availability) to use a set
260 of reference sequences for many or all loci that are not derived from the studied
261 group. Additionally these reference sequences may be more similar to some
262 sequenced samples than to others, which can introduce a bias in that the number
263 of recovered target loci per sample is based on how divergent their sequences are
264 to the reference sequence library. In other words, the 'one size fits all' approach
265 for recovering contig sequences is not the preferred option for most datasets and
266 may lead to taxonomic biases. For this reason it is recommended to generate
267 family, genus or even sample-specific reference libraries using the recovered
268 contigs, and use these to re-assemble the sequencing reads.
- 269 2. Intron/exon structure: Another reason for creating a new reference library from
270 the data is that available reference sequences often constitute exons, omitting the
271 interspersed intron sequences (as in the case of using bait sequences as the
272 reference library). The more variable introns in between exons are usually not
273 suitable for designing baits, they are too variable, but are extremely useful for
274 most phylogenetic analyses because they have more parsimony informative sites.

275 There is a good chance that the assembled contigs will contain parts of the trailing
276 introns or even span across the complete intron, connecting two exon sequences
277 (e.g. Bi et al. 2012). This is why it is preferable to use these usually longer and
278 more complete contig sequences for reference-based assembly, rather than the
279 shorter exon sequences from the bait sequence file, in order to capture all reads
280 that match either the exon or the trailing intron sequences at a locus.

281 3. Allelic variation: Remapping the reads in the process of reference-based assembly
282 will identify the different allele sequences at a given locus. This can also aid in
283 the evaluation of the ploidy level of samples and in identifying loci potentially
284 affected by paralogy.

285 4. Coverage: Reference-based assembly will give the user a better and more intuitive
286 overview over read-depth for all loci. There are excellent visualization softwares
287 (such as Tablet Milne et al. 2013) that help interpret the results.

288

289 *Novelty*

290 Several pipelines and collections of bioinformatics tools exist for processing sequencing
291 reads generated by MPS techniques, e.g. PHYLUCE (Faircloth 2015), GATK (McKenna
292 et al. 2010) and ‘reads2trees’ (Heyduk et al. 2016). In contrast to some of these existing
293 pipelines, SECAPR i) is targeted towards assembling full sequence data (as compared to
294 only SNP data, e.g. GATK); ii) is intended for general use (rather than project specific,
295 e.g. reads2trees); iii) is optimized particularly for non-model organisms and non-
296 standardized sequence capture datasets (as compared to specific exon sets, e.g.
297 PHYLUCE); iv) allows allele phasing and selection of the best loci based on read
298 coverage, which to our knowledge are novel to SECAPR. This is possible due to the
299 approach of generating a clade- or even sample-specific reference library from the
300 sequencing read data, which is then used for reference-based assembly; v) offers new
301 tools and plotting functions to give the user an overview of the sequencing data after each
302 processing step.

303

304 Conclusions

305 The SECAPR pipeline described here constitutes a bioinformatic tool for the processing
306 and alignment of raw Illumina sequence data. It is particularly useful for sequence
307 capture datasets and we show here how it can be applied to even challenging datasets of
308 non-model organisms.

309

310 Acknowledgements

311 We thank Estelle Proux-Wéra and Marcel Martin at the National Bioinformatics
312 Infrastructure Sweden at SciLifeLab for their support with turning the SECAPR pipeline
313 into a functioning conda package and for additional support in software development
314 questions. The code for some of the functions of the SECAPR pipeline is inspired from
315 similar functions in the PHYLUCE pipeline (Faircloth 2015).

316

317 Funding

318 This work was supported by the Swedish Research Council (B0569601), the European
319 Research Council under the European Union's Seventh Framework Programme
320 (FP/2007-2013, ERC Grant Agreement n. 331024), the Swedish Foundation for Strategic
321 Research, the Faculty of Science at the University of Gothenburg, the David Rockefeller
322 Center for Latin American Studies at Harvard University, and a Wallenberg Academy
323 Fellowship to A.A.; and a SciLifeLab Bioinformatics Long-term Support from the
324 Wallenberg Advanced Bioinformatics Infrastructure to A.A. and Bengt Oxelman.

325

326 Competing Interests

327 The authors declare there are no competing interests.

328

329 **Author contributions**

330 TA, CDB and AA conceived of this study, TA developed and implemented the pipeline
331 and analyzed the data with contribution from AZ, AC provided the empirical data. TA
332 wrote the manuscript with contributions from all authors.

333

334 **Availability**

335 The SECAPR pipeline is open source and freely available from
336 http://www.github.com/AntonelliLab/seqcap_processor. SECAPR and all software
337 dependencies can be downloaded as a virtual environment with the conda package
338 manager (<http://bioconda.github.io/recipes/secapr/README.html>). Installation
339 instructions, a detailed documentation and an empirical data tutorial with the *Geonoma*
340 sample data can be found at
341 http://github.com/AntonelliLab/seqcap_processor/blob/master/documentation.ipynb. The
342 raw sequencing data for all *Geonoma* samples is available at
343 <https://www.ncbi.nlm.nih.gov/sra/SRP131660>. All other empirical data produced in this
344 study is available from Zenodo (<https://doi.org/10.5281/zenodo.1162653>).

345

346 **References**

- 347 Andermann, Tobias, Alexandre M. Fernandes, Urban Olsson, Mats Topel, Bernard Pfeil,
348 Bengt Oxelman, Alexandre Aleixo, Brant C. Faircloth, and Alexandre Antonelli.
349 2018. "Allele Phasing Greatly Improves the Phylogenetic Utility of Ultraconserved
350 Elements." *bioRxiv*, January. Cold Spring Harbor Laboratory, 255752.
351 doi:10.1101/255752.
- 352 Bi, Ke, Dan Vanderpool, Sonal Singhal, Tyler Linderoth, Craig Moritz, and Jeffrey M
353 Good. 2012. "Transcriptome-Based Exon Capture Enables Highly Cost-Effective
354 Comparative Genomic Data Collection at Moderate Evolutionary Scales." *BMC*
355 *Genomics* 13 (1). BioMed Central: 403. doi:10.1186/1471-2164-13-403.
- 356 Bolger, Anthony M, Marc Lohse, and Bjoern Usadel. 2014. "Trimmomatic: A Flexible

- 357 Trimmer for Illumina Sequence Data.” *Bioinformatics* 30 (15): 2114–20.
358 doi:10.1093/bioinformatics/btu170.
- 359 Botero-Castro, Fidel, Marie Ka Tilak, Fabienne Justy, François Catzefflis, Frédéric
360 Delsuc, and Emmanuel J P Douzery. 2013. “Next-Generation Sequencing and
361 Phylogenetic Signal of Complete Mitochondrial Genomes for Resolving the
362 Evolutionary History of Leaf-Nosed Bats (Phyllostomidae).” *Molecular*
363 *Phylogenetics and Evolution* 69 (3). Elsevier Inc.: 728–39.
364 doi:10.1016/j.ympev.2013.07.003.
- 365 Bravo, Gustavo A, Alexandre Antonelli, Christine D Bacon, Krzysztof Bartoszek, Mozes
366 Blom, Stella Huynh, Graham Jones, et al. 2018. “Embracing Heterogeneity:
367 Building the Tree of Life and the Future of Phylogenomics,” January. PeerJ Inc.
368 doi:10.7287/peerj.preprints.26449v3.
- 369 Dransfield, J, NW Uhl, CB Asmussen, and WJ Baker. 2008. “Genera Palmarum.” *Royal*
370 *Botanic Gardens*, 410–42.
- 371 Faircloth, Brant C. 2015. “PHYLUCE Is a Software Package for the Analysis of
372 Conserved Genomic Loci.” *Bioinformatics* 32 (5): 786–88.
373 doi:10.1093/bioinformatics/btv646.
- 374 Faircloth, Brant C, Michael G Branstetter, Noor D White, and Seán G Brady. 2015.
375 “Target Enrichment of Ultraconserved Elements from Arthropods Provides a
376 Genomic Perspective on Relationships among Hymenoptera.” *Molecular Ecology*
377 *Resources* 15 (3): 489–501. doi:10.1111/1755-0998.12328.
- 378 Faircloth, Brant C, John E McCormack, Nicholas G Crawford, Michael G Harvey, Robb
379 T Brumfield, and Travis C Glenn. 2012. “Ultraconserved Elements Anchor
380 Thousands of Genetic Markers Spanning Multiple Evolutionary Timescales.”
381 *Systematic Biology* 61 (5): 717–26. doi:10.1093/sysbio/sys004.
- 382 Gnirke, Andreas, Alexandre Melnikov, Jared Maguire, Peter Rogov, Emily M LeProust,
383 William Brockman, Timothy Fennell, et al. 2009. “Solution Hybrid Selection with
384 Ultra-Long Oligonucleotides for Massively Parallel Targeted Sequencing.” *Nature*
385 *Biotechnology* 27 (2). Nature Publishing Group: 182–89. doi:10.1038/nbt.1523.

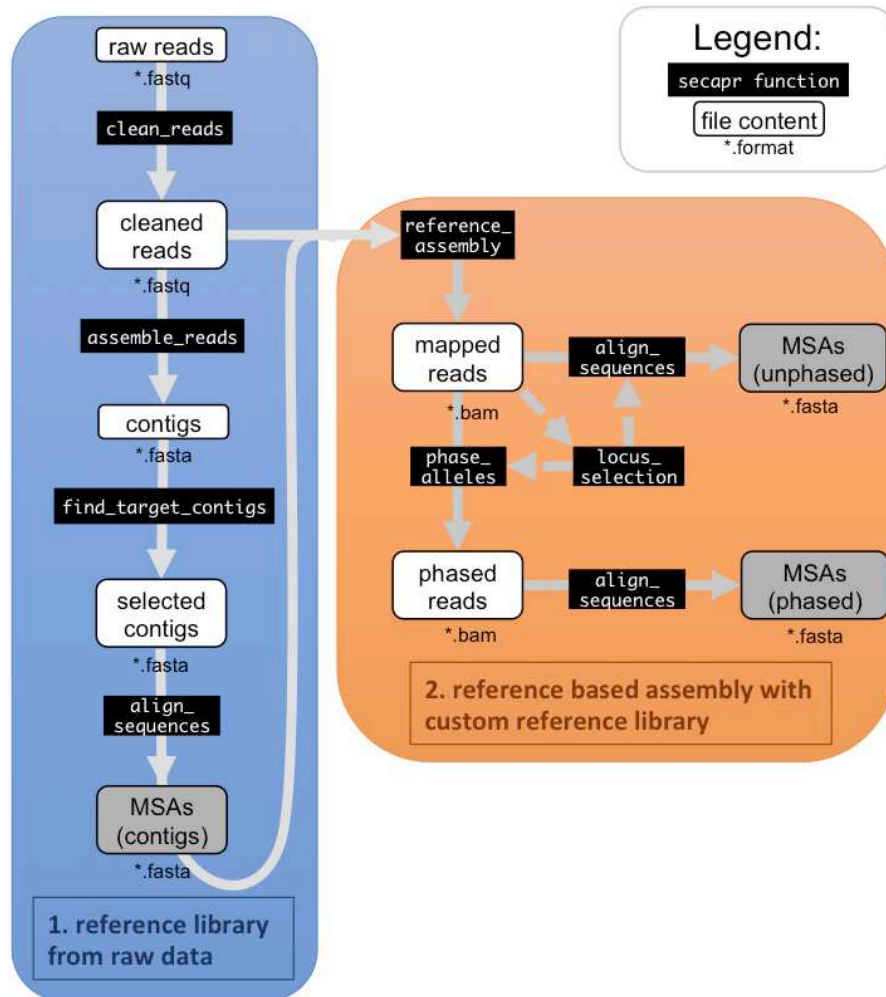
- 386 Harris, R.S. 2007. “Improved Pairwise Alignment of Genomic DNA.” The Pennsylvania
387 State University.
- 388 He, D., A. Choi, K. Pipatsrisawat, A. Darwiche, and E. Eskin. 2010. “Optimal
389 Algorithms for Haplotype Assembly from Whole-Genome Sequence Data.”
390 *Bioinformatics* 26 (12). Oxford University Press: i183–90.
391 doi:10.1093/bioinformatics/btq215.
- 392 Henderson, Andrew James. 2011. *A Revision of Geonoma (Arecaceae)*.
- 393 Heyduk, Karolina, Dorset W Trapnell, Craig F Barrett, and Jim Leebens-Mack. 2016.
394 “Phylogenomic Analyses of Species Relationships in the Genus Sabal (Arecaceae)
395 Using Targeted Sequence Capture.” *Biological Journal of the Linnean Society* 117:
396 106–20.
- 397 Hunt, Martin, Chris Newbold, Matthew Berriman, and Thomas D Otto. 2014. “A
398 Comprehensive Evaluation of Assembly Scaffolding Tools.” *Genome Biology* 15
399 (3). BioMed Central: R42. doi:10.1186/gb-2014-15-3-r42.
- 400 Kadlec, Malvina, Dirk U. Bellstedt, Nicholas C. Le Maitre, and Michael D. Pirie. 2017.
401 “Targeted NGS for Species Level Phylogenomics: ‘made to Measure’ or ‘one Size
402 Fits All’?” *PeerJ* 5: e3569. doi:10.7717/peerj.3569.
- 403 Li, Heng, and Richard Durbin. 2010. “Fast and Accurate Long-Read Alignment with
404 Burrows-Wheeler Transform.” *Bioinformatics* 26 (5): 589–95.
405 doi:10.1093/bioinformatics/btp698.
- 406 Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor
407 Marth, Goncalo Abecasis, and Richard Durbin. 2009. “The Sequence
408 Alignment/Map Format and SAMtools.” *Bioinformatics* 25 (16): 2078–79.
409 doi:10.1093/bioinformatics/btp352.
- 410 McKenna, Aaron, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis,
411 Andrew Kernytsky, Kiran Garimella, et al. 2010. “The Genome Analysis Toolkit: A
412 MapReduce Framework for Analyzing next-Generation DNA Sequencing Data.”
413 *Genome Research* 20 (9). Cold Spring Harbor Laboratory Press: 1297–1303.
414 doi:10.1101/gr.107524.110.

- 415 Milne, Iain, Gordon Stephen, Micha Bayer, Peter J A Cock, Leighton Pritchard, Linda
416 Cardle, Paul D Shaw, and David Marshall. 2013. “Using Tablet for Visual
417 Exploration of Second-Generation Sequencing Data.” *Briefings in Bioinformatics* 14
418 (2): 193–202. doi:10.1093/bib/bbs012.
- 419 Reuter, Jason A, Damek V Spacek, and Michael P Snyder. 2015. “High-Throughput
420 Sequencing Technologies.” *Molecular Cell* 58 (4). NIH Public Access: 586–97.
421 doi:10.1016/j.molcel.2015.05.004.
- 422 Simpson, Jared T, Kim Wong, Shaun D Jackman, Jacqueline E Schein, Steven J M Jones,
423 and Inanç Birol. 2009. “ABYSS: A Parallel Assembler for Short Read Sequence
424 Data.” *Genome Research* 19 (6): 1117–23. doi:10.1101/gr.089532.108.
- 425 Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014.
426 “Target Capture and Massively Parallel Sequencing of Ultraconserved Elements for
427 Comparative Studies at Shallow Evolutionary Time Scales.” *Systematic Biology* 63
428 (1): 83–95. doi:10.1093/sysbio/syt061.
- 429 Smith, Brian Tilston, John E. McCormack, Andrés M. Cuervo, Michael. J. Hickerson,
430 Alexandre Aleixo, Carlos Daniel Cadena, Jorge Pérez-Emán, et al. 2014. “The
431 Drivers of Tropical Speciation.” *Nature* 515 (7527). Nature Publishing Group, a
432 division of Macmillan Publishers Limited. All Rights Reserved.: 406–9.
433 doi:10.1038/nature13687.
- 434

435 **Table 1: Summary statistics for all samples, produced by SECAPR.** Reported for
 436 each sample are the number of sequencing reads in the FASTQ sequencing files, before
 437 (1. column) and after (2. column) cleaning and trimming, the total count of assembled *de*
 438 *novo* contigs (3. column), the number of filtered contigs that matched target loci (4.
 439 column) and the number of sequencing reads that mapped to the new reference library
 440 generated from the contig MSAs during reference-based assembly (5. column). These
 441 summary statistics are automatically compiled and appended to a log file
 442 (*summary_stats.txt*) during different steps in the SECAPR pipeline.

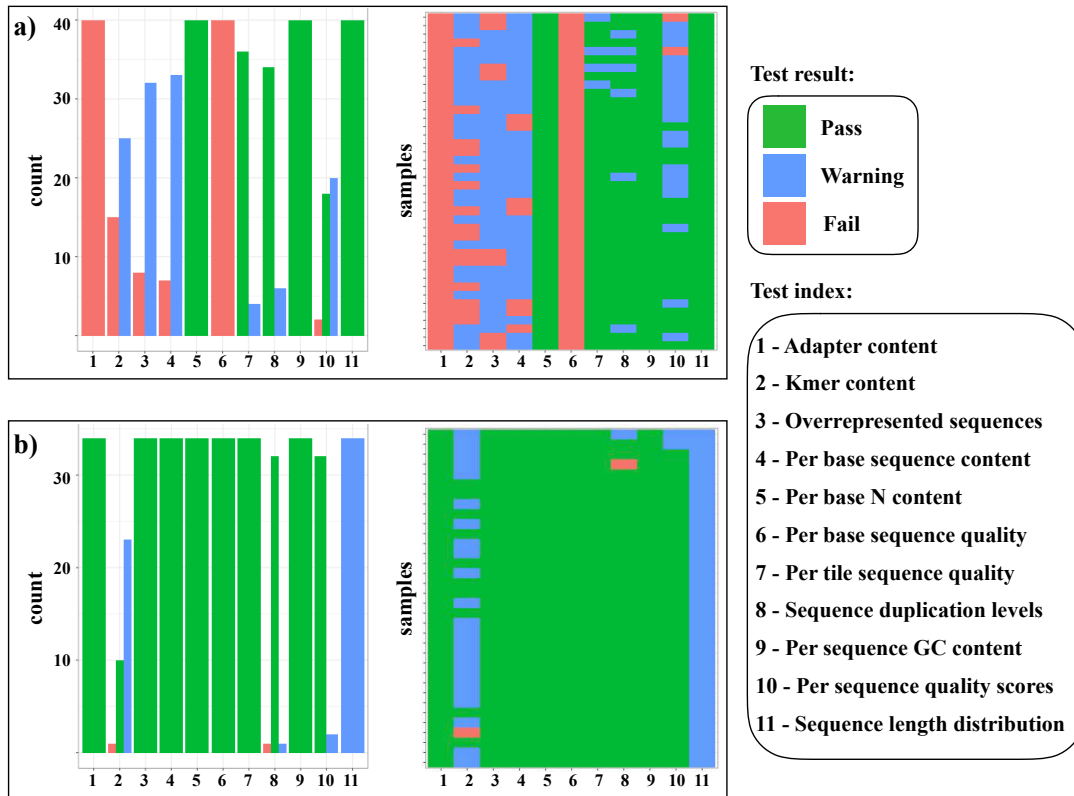
Sample ID	FASTQ read pairs (raw)	FASTQ read pairs (cleaned)	Total contig count	Recovered target contigs	Reads on target regions
1087	291089	276072	277628	562	22308
1086	244726	231326	230122	516	17969
1140	206106	192676	153377	469	18039
1083	377228	352646	309993	534	31922
1082	277999	262378	258359	556	19491
1085	307671	291377	309561	512	22030
1079	315801	298450	306369	550	13969
1061	209586	192407	177910	545	14474
1068	295402	278069	264865	563	22013
1063	354795	336356	356512	525	20439
1080	459485	434951	433954	531	41068
1065	217725	205290	204082	544	13524
1073	302798	286021	289612	529	15598
1070	295822	278011	295557	539	19288
1064	408723	384908	405080	543	21531
1074	408370	383604	398758	531	25476
1166	405667	385442	410292	544	29697

443



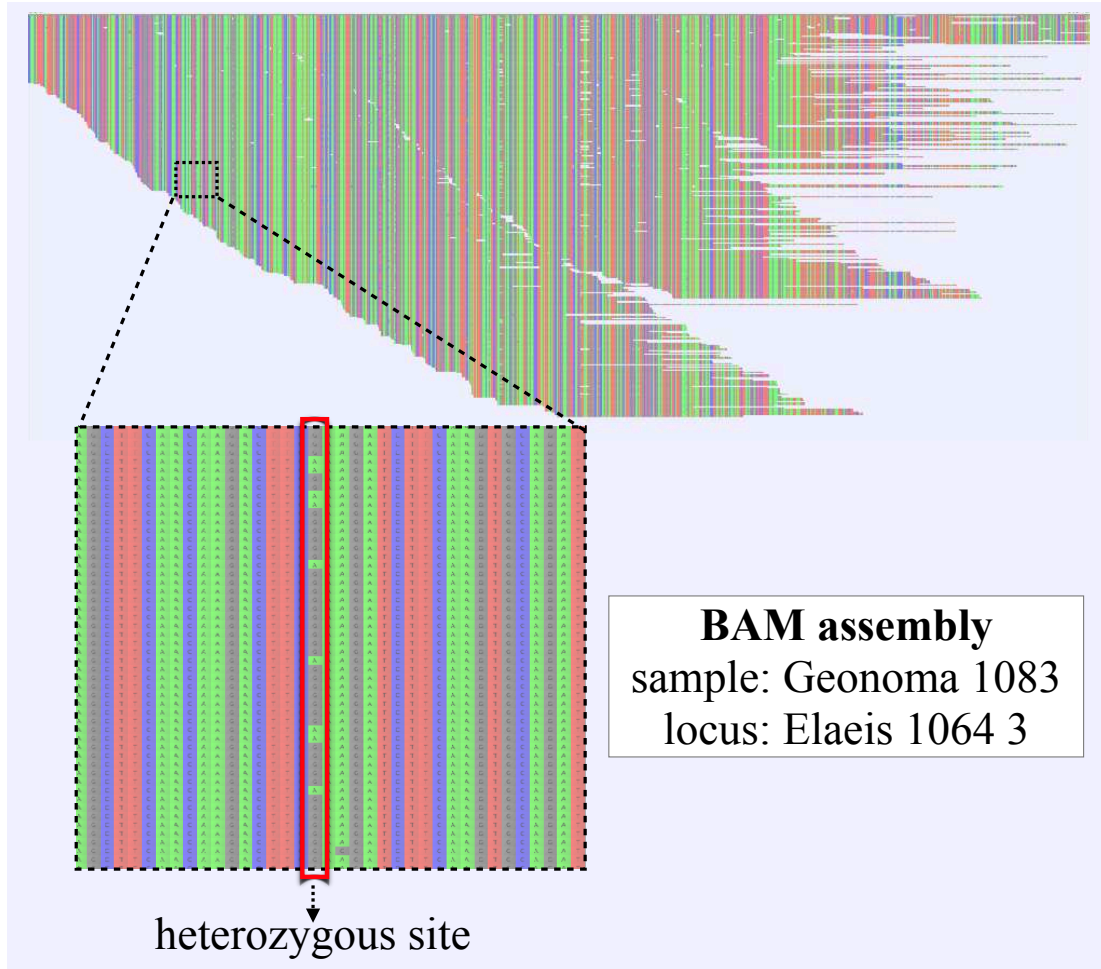
444

445 **Figure 1: SECAPR analytical workflow.** The flowchart shows the basic SECAPR
 446 functions, which are separated into two separate steps (colored boxes). Blue box (1.
 447 reference library from raw data): in this step the raw reads are cleaned and assembled into
 448 contigs (*de novo* assembly); Orange box (2. reference based assembly with custom
 449 reference library): the contigs from the previous step are used for reference-based
 450 assembly, enabling allele phasing and additional quality control options, e.g. concerning
 451 read-coverage. Black boxes show SECAPR commands and white boxes represent the
 452 input and output data of the respective function. Boxes marked in grey represent multiple
 453 sequence alignments (MSAs) generated with SECAPR, which can be used for
 454 phylogenetic inference.



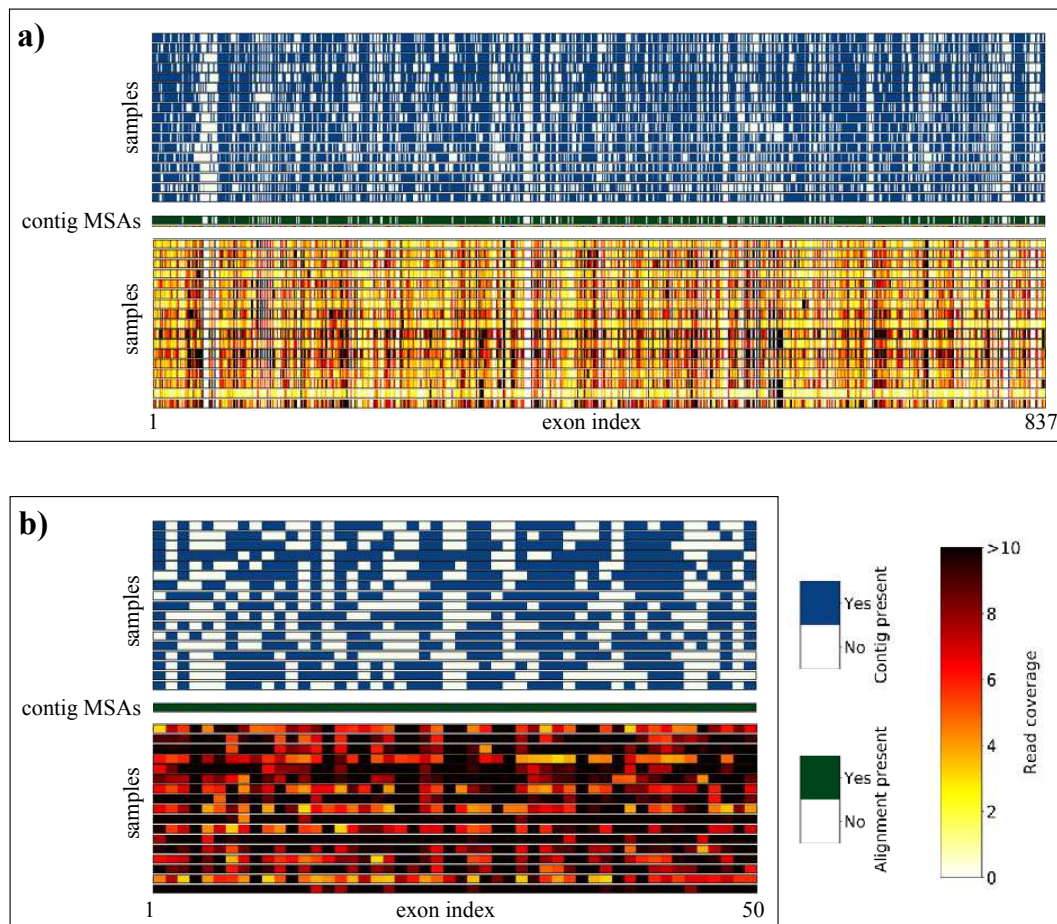
455

456 **Figure 2: Overview of FASTQC quality test results.** a) Before and b) after cleaning and
 457 adapter trimming of sequencing reads with the SECAPR function *clean_reads*. This plot,
 458 as produced by SECAPR, provides an overview of the complete dataset and helps to
 459 gauge if the chosen cleaning parameters are appropriate for the dataset. The summary
 460 plots show the FASTQC test results, divided into three categories: passed (green),
 461 warning (blue) and failed (red). The x-axis of all plots contains the eleven different
 462 quality tests (see legend). The bar-plots (left panels) represent the counts of each test
 463 result (pass, warning or fail) across all samples. The matrix plots (right panels) show the
 464 test result of each test for each sample individually (y-axis). This information can be used
 465 to evaluate both, which specific parameters need to be adjusted and which samples are
 466 the most problematic.



467

468 **Figure 3: Reference-based assembly including heterozygous sites.** BAM-assembly file
469 as generated with the SECAPR *reference_assembly* function, shown exemplarily for one
470 exon locus (1/837) of one of the *Geonoma* samples (1/17). The displayed assembly
471 contains all FASTQ sequencing reads that could be mapped to the reference sequence
472 (top panel). The reference sequence in this case is the de-novo contig that was matched to
473 the reference exon 'Elaeis 1064 3'. DNA bases are color-coded (A - green, C - blue, G -
474 black, T - red). The enlarged section (bottom panel) contains a heterozygous site, which
475 likely represents allelic variation, as we both variants A and G are found at approximately
476 equal ratio.



477

478 **Figure 4: Overview of sequence yield for *Geonoma* sample data, produced with**
 479 **SECAPR.** Each column in these matrix plots represents a separate exon locus a) for all
 480 loci targeted during sequence capture and b) for the selection of the 50 loci with the best
 481 read coverage, using the SECAPR function *locus_selection* (see Supplementary Table S5
 482 for loci-names corresponding to indices on x-axes). Top panels in a) and b) show if *de*
 483 *novo* contigs could be assembled (blue) or not (white) for the respective locus (column)
 484 and sample (row). Contig MSAs were generated for all loci that could be recovered for at
 485 least 3 samples (center row - green). The bottom panels of a) and b) show the read
 486 coverage (see legend) for each exon locus after reference-based assembly. The reference
 487 library for the assembly consisted of the consensus sequences of each contig MSA, and
 488 hence is genus specific for *Geonoma*.