

1 **Author Cover Page**

2

3 Article submission to PeerJ

4 Manuscript category: Bioinformatics Tool

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

6 Article title: **SECAPR - A bioinformatics pipeline for the rapid and user-friendly**  
7 **alignment of hybrid enrichment sequences, from raw reads to alignments**

8

9 Authors: Tobias Andermann\* <sup>(1,2)</sup>, Ángela Cano <sup>(2,3)</sup>, Alexander Zizka <sup>(1,2)</sup>, Christine  
10 Bacon <sup>(1,2)</sup>, Alexandre Antonelli <sup>(1,2,4,5)</sup>

11

12 Affiliations:

13 <sup>1</sup> Department of Biological and Environmental Sciences, University of Gothenburg, Box  
14 461, SE 405 30, Göteborg, Sweden

15 <sup>2</sup> Gothenburg Global Biodiversity Centre, Göteborg, 41319, Sweden

16 <sup>3</sup> Department of Botany and Plant Biology, University of Geneva, Geneva, Switzerland

17 <sup>4</sup> Gothenburg Botanical Garden, Göteborg, 41319, Sweden

18 <sup>5</sup> Department of Organismic and Evolutionary Biology, Harvard University, Cambridge,  
19 MA 02138 USA

20

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

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

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

## 40 Introduction

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

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

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

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

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

## Methods

### *The SECAPR pipeline in a nutshell*

SECAPR is a platform-independent pipeline written in python, and tested for full functionality on Linux and MacOS. It can be easily downloaded and installed using the conda ([www.conda.io/docs/](http://www.conda.io/docs/)) package manager (see Availability). The conda package automatically includes all software dependencies and ensures version compatibility and standardized working paths. The strength of SECAPR is that it channels the main functionalities of many commonly used bioinformatics programs and enables the user to apply these to sets of samples, rather than having to apply different software to each sample individually.

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

#### *1. Quality filtering and adapter trimming*

2. *De novo* contig assembly
3. Selection of target contigs
4. Building MSAs from contigs
5. Reference-based assembly
6. Allele phasing

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

#### *Description of the SECAPR workflow*

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

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

3. *Selection of target contigs (secapr find\_target\_contigs)*. The SECAPR function *find\_target\_contigs* identifies and extracts those contigs that represent the DNA targets of interest. This function implements the program LASTZ (formerly BLASTZ, Harris 2007)

by searching the contig files for matches with a user-provided FASTA-formatted reference library. For sequence capture datasets, a suitable reference library is the reference file that was used for synthesizing the RNA baits, which will return all contigs that match the enriched loci of interest. The *find\_target\_contigs* function identifies potentially paralogous loci (loci that have several matching contigs) and excludes these from further processing. It further allows the user to keep or exclude long contigs that match several adjacent reference loci, which can occur if the reference file contains sequences that are located in close proximity to each other on the genome (e.g. several separate exons of the same gene).

4. *Building MSAs from contigs (secapr align\_sequences)*. The SECAPR function *align\_sequences* builds multiple sequence alignments (MSAs) from the target contigs that were identified in the previous step. The function builds a separate MSA for each locus with matching contigs for  $\geq 3$  samples.

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

The *reference\_assembly* function includes different options for generating a reference library for all loci of interest:

- *--reference\_type alignment-consensus*: The user provides a link to a folder containing MSAs, e.g. the folder with the contig MSAs from the previous step, and the function calculates a consensus sequence from each alignment. These consensus sequences are then used as the reference sequence for the assembly. This function is recommended when running reference-based assembly for groups of closely related samples (e.g. samples from the same genus or family).

- *--reference\_type sample-specific*: From the MSAs, the function extracts the contigs for each sample and uses them as a sample-specific reference library. If the user decides to use this function it is recommendable to only use alignments for reference that contain sequences for all samples. This will ensure that the same loci are being assembled for all samples.
- *--reference\_type user-ref-lib*: The user can provide a FASTA file containing a custom reference library.

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

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

### *Benchmarking with empirical data*

We demonstrate the functionalities of SECAPR on a novel dataset of target sequencing reads of *Geonoma*, one of the most species-rich palm genera of tropical Central and South America. (Dransfield et al. 2008) (Henderson 2011). Our data comprised newly generated Illumina sequence data for 17 samples of 14 *Geonoma* species (Supplementary Table S1), enriched through sequence capture. The bait set for sequence capture was designed specifically for palms by Heyduk et al. (2016) to target 176 genes with in total 837 exons. More detailed information about the generation of the sequence data can be

found in Appendix 1 (Supplemental Material). All settings and commands used during processing of the sequence data can be found in the SECAPR documentation on our GitHub page.

## Results and Discussion

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

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

However, reference-based assembly is very sensitive toward the user providing orthologous reference sequences that are similar enough to the sequencing reads of the studied organisms. If the reference sequences are too divergent from the sequenced organisms, only a small fraction of the existing orthologous sequencing reads will be successfully assembled for each locus. In contrast, when relaxing similarity thresholds and other mapping parameters too much (e.g. to increase the fraction of reads included in the assembly) there is higher a risk of assembling non-orthologous reads, which can lead to chimeric sequences being assembled. This can be a problem, particularly in cases of datasets containing non-model organisms, since suitable reference sequences for all loci usually do not exist.

For this reason, the SECAPR workflow encourages the user to use these two different assembly approaches in concert (Figure 1). Our general suggestion is to first assemble contig MSAs for all regions of interest, resulting from *de novo* assembly and then use these MSAs to build a reference library for reference-based assembly. In that case



SECAPR produces a reference library from the sequencing data itself, which is specific for the taxonomic group of interest or even for the individual sample.

A common approach is to stop data processing after the *de novo* assembly step and then use the contig MSAs for phylogenetic analyses (e.g. Faircloth et al. 2012; B. T. Smith et al. 2014; Faircloth 2015). However, there are several reasons to further process the data. These further steps include generating new reference libraries for all samples, and using the raw sequence data for a reference-based assembly. There are several reasons for carrying out these additional steps:

1. Sensitivity: In order to identify *de novo* contigs that are orthologous to the loci of interest, the user is usually forced (because of the lack of availability) to use a set of reference sequences for many or all loci that are not derived from the studied group. Additionally these reference sequences may be more similar to some sequenced samples than to others, which can introduce a bias in that the number of recovered target loci per sample is based on how divergent their sequences are to the reference sequence library. In other words, the 'one size fits all' approach for recovering contig sequences is not the preferred option for most datasets and may lead to taxonomic biases. For this reason it is recommended to generate family, genus or even sample-specific reference libraries using the recovered contigs, and use these to re-assemble the sequencing reads.
2. Intron/exon structure: Another reason for creating a new reference library from the data is that available reference sequences often constitute exons, omitting the interspersed intron sequences (as in the case of using bait sequences as the reference library). The more variable introns in between exons are usually not suitable for designing baits, they are too variable, but are extremely useful for most phylogenetic analyses because they have more parsimony informative sites. There is a good chance that the assembled contigs will contain parts of the trailing introns or even span across the complete intron, connecting two exon sequences (e.g. Bi et al. 2012). This is why it is preferable to use these usually longer and more complete contig sequences for reference-based assembly, rather than the shorter exon sequences from the bait sequence file, in order to capture all reads that match either the exon or the trailing intron sequences at a locus.

3. Allelic variation: Remapping the reads in the process of reference-based assembly will identify the different allele sequences at a given locus. This can also aid in the evaluation of the ploidy level of samples and in identifying loci potentially affected by paralogy.
4. Coverage: Reference-based assembly will give the user a better and more intuitive overview over read-depth for all loci. There are excellent visualization softwares (such as Tablet Milne et al. 2013) that help interpret the results.

### *Empirical evaluation*

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

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

During reference-based assembly (*secapr reference\_assembly*) we mapped the reads against the consensus sequence of the contig MSAs for all loci. We found an average of 439 exon loci (stdev=82) per sample that were covered by more than three reads (average coverage across complete locus, Figure 4a). Hence, our approach of mapping FASTQ reads to libraries compiled from the data leads to an increase of recovered loci per sample, from 323 resulting from *de novo* assembly to 439 from the referenced-based assembly (36% increase). Further, the number of loci that were recovered with sufficient coverage for all samples increased by 116%, from 45 after the *de novo* assembly, to 97 after the reference-based assembly (Supplementary Table S4). We extracted the 50 loci with the best coverage across all samples (*secapr locus\_selection*), as shown in Figure 4b. In cases of irregular read-coverage across samples (as in our sample *Geonoma* data), we strongly recommend the use of the *locus\_selection* function before further processing the data, as demonstrated in our tutorial (see

Availability).

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

## Novelty

Several pipelines and collections of bioinformatics tools exist for processing sequencing reads generated by MPS techniques, e.g. PHYLUCE (Faircloth 2015), GATK (McKenna et al. 2010) and ‘reads2trees’ (Heyduk et al. 2016). In contrast to some of these existing

pipelines, SECAPR i) is targeted towards assembling full sequence data (as compared to only SNP data, e.g. GATK); ii) is intended for general use (rather than project specific, e.g. reads2trees); iii) is optimized particularly for non-model organisms and non-standardized sequence capture datasets (as compared to specific exon sets, e.g. PHYLUCE); iv) allows allele phasing and selection of the best loci based on read coverage, which to our knowledge are novel to SECAPR. This is possible due to the approach of generating a clade- or even sample-specific reference library from the sequencing read data, which is then used for reference-based assembly; v) offers new tools and plotting functions to give the user an overview of the sequencing data after each processing step.

### Acknowledgements

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

### Funding

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

## Competing Interests

The authors declare there are no competing interests.

## Author contributions

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

## Availability

The SECAPR pipeline is open source and freely available from [http://www.github.com/AntonelliLab/seqcap\\_processor](http://www.github.com/AntonelliLab/seqcap_processor). SECAPR and all software dependencies can be downloaded and installed with the conda package manager using the command 'conda create -c bioconda -n secapr\_env secapr' in the bash-shell command line. A detailed tutorial of all SECAPR functions using the *Geonoma* sample data can be found at [http://github.com/AntonelliLab/seqcap\\_processor/blob/master/documentation.ipynb](http://github.com/AntonelliLab/seqcap_processor/blob/master/documentation.ipynb). The raw sequencing read data of the *Geonoma* sample data is available at <https://www.ncbi.nlm.nih.gov/sra/SRP131660>. All other empirical data produced in this project is available from Zenodo (<https://doi.org/10.5281/zenodo.1162653>).

## References

- Andermann, Tobias, Alexandre M. Fernandes, Urban Olsson, Mats Topel, Bernard Pfeil, Bengt Oxelman, Alexandre Aleixo, Brant C. Faircloth, and Alexandre Antonelli. 2018. "Allele Phasing Greatly Improves the Phylogenetic Utility of Ultraconserved Elements." *bioRxiv*, January. Cold Spring Harbor Laboratory, 255752. doi:10.1101/255752.
- Bi, Ke, Dan Vanderpool, Sonal Singhal, Tyler Linderoth, Craig Moritz, and Jeffrey M

- 354 Good. 2012. "Transcriptome-Based Exon Capture Enables Highly Cost-Effective  
355 Comparative Genomic Data Collection at Moderate Evolutionary Scales." *BMC*  
356 *Genomics* 13 (1). BioMed Central: 403. doi:10.1186/1471-2164-13-403.
- 357 Bolger, Anthony M, Marc Lohse, and Bjoern Usadel. 2014. "Trimmomatic: A Flexible  
358 Trimmer for Illumina Sequence Data." *Bioinformatics* 30 (15): 2114–20.  
359 doi:10.1093/bioinformatics/btu170.
- 360 Botero-Castro, Fidel, Marie Ka Tilak, Fabienne Justy, François Catzefflis, Frédéric  
361 Delsuc, and Emmanuel J P Douzery. 2013. "Next-Generation Sequencing and  
362 Phylogenetic Signal of Complete Mitochondrial Genomes for Resolving the  
363 Evolutionary History of Leaf-Nosed Bats (Phyllostomidae)." *Molecular*  
364 *Phylogenetics and Evolution* 69 (3). Elsevier Inc.: 728–39.  
365 doi:10.1016/j.ympev.2013.07.003.
- 366 Bravo, Gustavo A, Alexandre Antonelli, Christine D Bacon, Krzysztof Bartoszek, Mozes  
367 Blom, Stella Huynh, Graham Jones, et al. 2018. "Embracing Heterogeneity:  
368 Building the Tree of Life and the Future of Phylogenomics," January. PeerJ Inc.  
369 doi:10.7287/peerj.preprints.26449v3.
- 370 Dransfield, J, NW Uhl, CB Asmussen, and WJ Baker. 2008. "Genera Palmarum." *Royal*  
371 *Botanic Gardens*, 410–42.
- 372 Faircloth, Brant C. 2015. "PHYLUCE Is a Software Package for the Analysis of  
373 Conserved Genomic Loci." *Bioinformatics* 32 (5): 786–88.  
374 doi:10.1093/bioinformatics/btv646.
- 375 Faircloth, Brant C, Michael G Branstetter, Noor D White, and Seán G Brady. 2015.  
376 "Target Enrichment of Ultraconserved Elements from Arthropods Provides a  
377 Genomic Perspective on Relationships among Hymenoptera." *Molecular Ecology*  
378 *Resources* 15 (3): 489–501. doi:10.1111/1755-0998.12328.
- 379 Faircloth, Brant C, John E McCormack, Nicholas G Crawford, Michael G Harvey, Robb  
380 T Brumfield, and Travis C Glenn. 2012. "Ultraconserved Elements Anchor  
381 Thousands of Genetic Markers Spanning Multiple Evolutionary Timescales."  
382 *Systematic Biology* 61 (5): 717–26. doi:10.1093/sysbio/sys004.

- Gnirke, Andreas, Alexandre Melnikov, Jared Maguire, Peter Rogov, Emily M LeProust, William Brockman, Timothy Fennell, et al. 2009. "Solution Hybrid Selection with Ultra-Long Oligonucleotides for Massively Parallel Targeted Sequencing." *Nature Biotechnology* 27 (2). Nature Publishing Group: 182–89. doi:10.1038/nbt.1523.
- Harris, R.S. 2007. "Improved Pairwise Alignment of Genomic DNA." The Pennsylvania State University.
- He, D., A. Choi, K. Pipatsrisawat, A. Darwiche, and E. Eskin. 2010. "Optimal Algorithms for Haplotype Assembly from Whole-Genome Sequence Data." *Bioinformatics* 26 (12). Oxford University Press: i183–90. doi:10.1093/bioinformatics/btq215.
- Henderson, Andrew James. 2011. *A Revision of Geonoma (Arecaceae)*.
- Heyduk, Karolina, Dorset W Trapnell, Craig F Barrett, and Jim Leebens-Mack. 2016. "Phylogenomic Analyses of Species Relationships in the Genus Sabal (Arecaceae) Using Targeted Sequence Capture." *Biological Journal of the Linnean Society* 117: 106–20.
- Hunt, Martin, Chris Newbold, Matthew Berriman, and Thomas D Otto. 2014. "A Comprehensive Evaluation of Assembly Scaffolding Tools." *Genome Biology* 15 (3). BioMed Central: R42. doi:10.1186/gb-2014-15-3-r42.
- Kadlec, Malvina, Dirk U. Bellstedt, Nicholas C. Le Maitre, and Michael D. Pirie. 2017. "Targeted NGS for Species Level Phylogenomics: 'made to Measure' or 'one Size Fits All'?" *PeerJ* 5: e3569. doi:10.7717/peerj.3569.
- Li, Heng, and Richard Durbin. 2010. "Fast and Accurate Long-Read Alignment with Burrows-Wheeler Transform." *Bioinformatics* 26 (5): 589–95. doi:10.1093/bioinformatics/btp698.
- Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, and Richard Durbin. 2009. "The Sequence Alignment/Map Format and SAMtools." *Bioinformatics* 25 (16): 2078–79. doi:10.1093/bioinformatics/btp352.
- McKenna, Aaron, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis,



Andrew Kernysky, Kiran Garimella, et al. 2010. “The Genome Analysis Toolkit: A MapReduce Framework for Analyzing next-Generation DNA Sequencing Data.” *Genome Research* 20 (9). Cold Spring Harbor Laboratory Press: 1297–1303. doi:10.1101/gr.107524.110.

Milne, Iain, Gordon Stephen, Micha Bayer, Peter J A Cock, Leighton Pritchard, Linda Cardle, Paul D Shaw, and David Marshall. 2013. “Using Tablet for Visual Exploration of Second-Generation Sequencing Data.” *Briefings in Bioinformatics* 14 (2): 193–202. doi:10.1093/bib/bbs012.

Reuter, Jason A, Damek V Spacek, and Michael P Snyder. 2015. “High-Throughput Sequencing Technologies.” *Molecular Cell* 58 (4). NIH Public Access: 586–97. doi:10.1016/j.molcel.2015.05.004.

Simpson, Jared T, Kim Wong, Shaun D Jackman, Jacqueline E Schein, Steven J M Jones, and Inanç Birol. 2009. “ABYSS: A Parallel Assembler for Short Read Sequence Data.” *Genome Research* 19 (6): 1117–23. doi:10.1101/gr.089532.108.

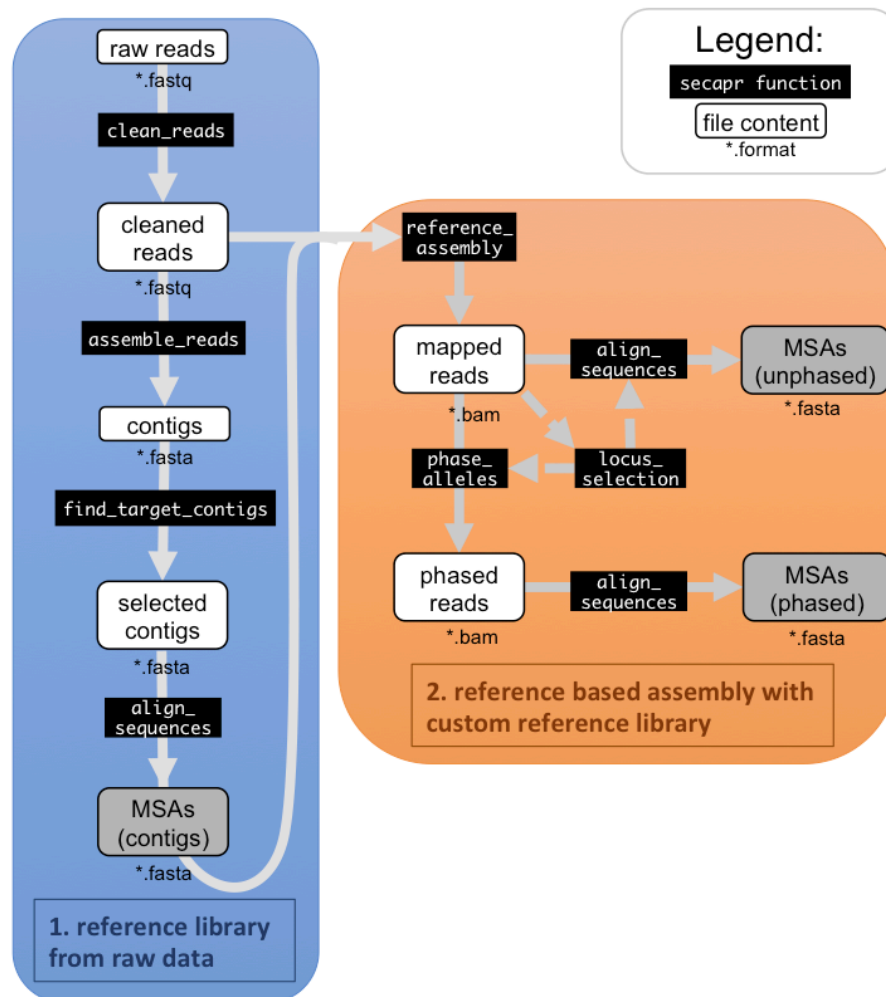
Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014. “Target Capture and Massively Parallel Sequencing of Ultraconserved Elements for Comparative Studies at Shallow Evolutionary Time Scales.” *Systematic Biology* 63 (1): 83–95. doi:10.1093/sysbio/syt061.

Smith, Brian Tilston, John E. McCormack, Andrés M. Cuervo, Michael. J. Hickerson, Alexandre Aleixo, Carlos Daniel Cadena, Jorge Pérez-Emán, et al. 2014. “The Drivers of Tropical Speciation.” *Nature* 515 (7527). Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.: 406–9. doi:10.1038/nature13687.

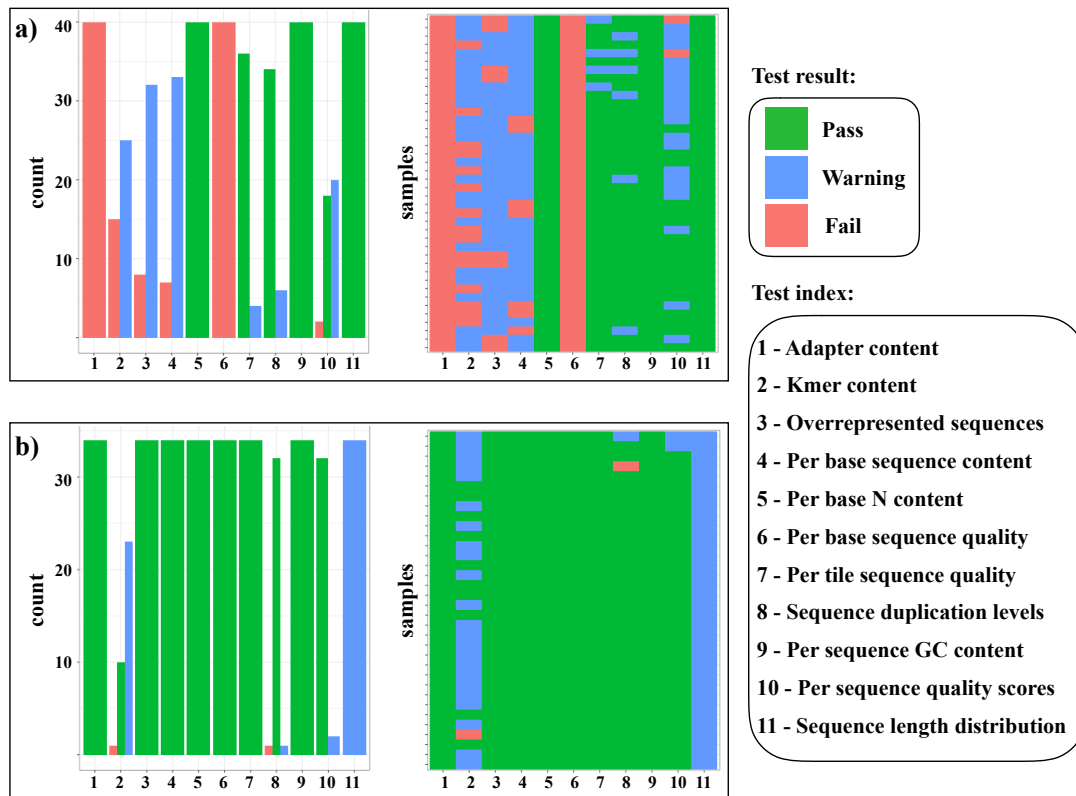


**Table 1: Summary statistics for all samples, produced by SECAPR.** Reported for each sample are the number of sequencing reads in the FASTQ sequencing files, before (1. column) and after (2. column) cleaning and trimming, the total count of assembled *de novo* contigs (3. column), the number of filtered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

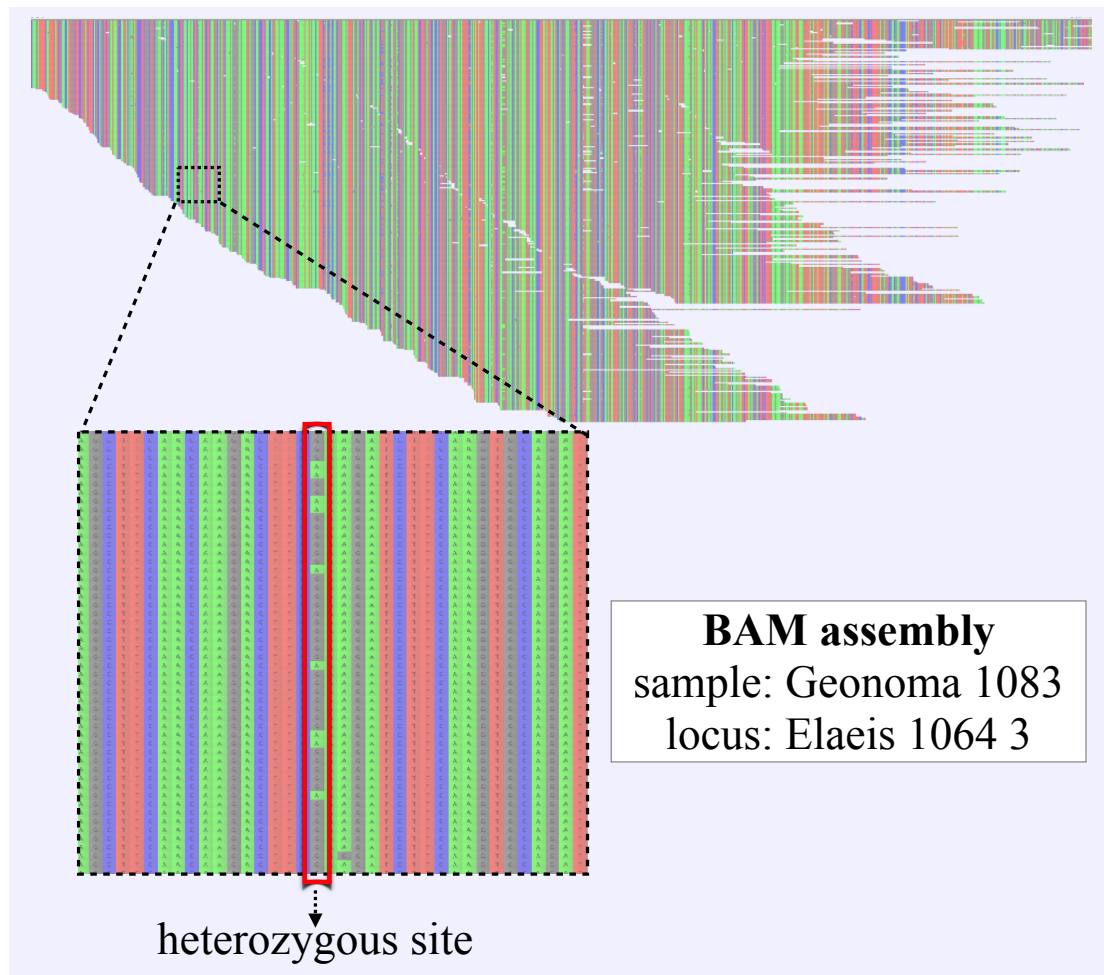


**Figure 1: SECAPR analytical workflow.** The flowchart shows the basic SECAPR functions, which are separated into two separate steps (colored boxes). Blue box (1. reference library from raw data): in this step the raw reads are cleaned and assembled into contigs (*de novo* assembly); Orange box (2.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

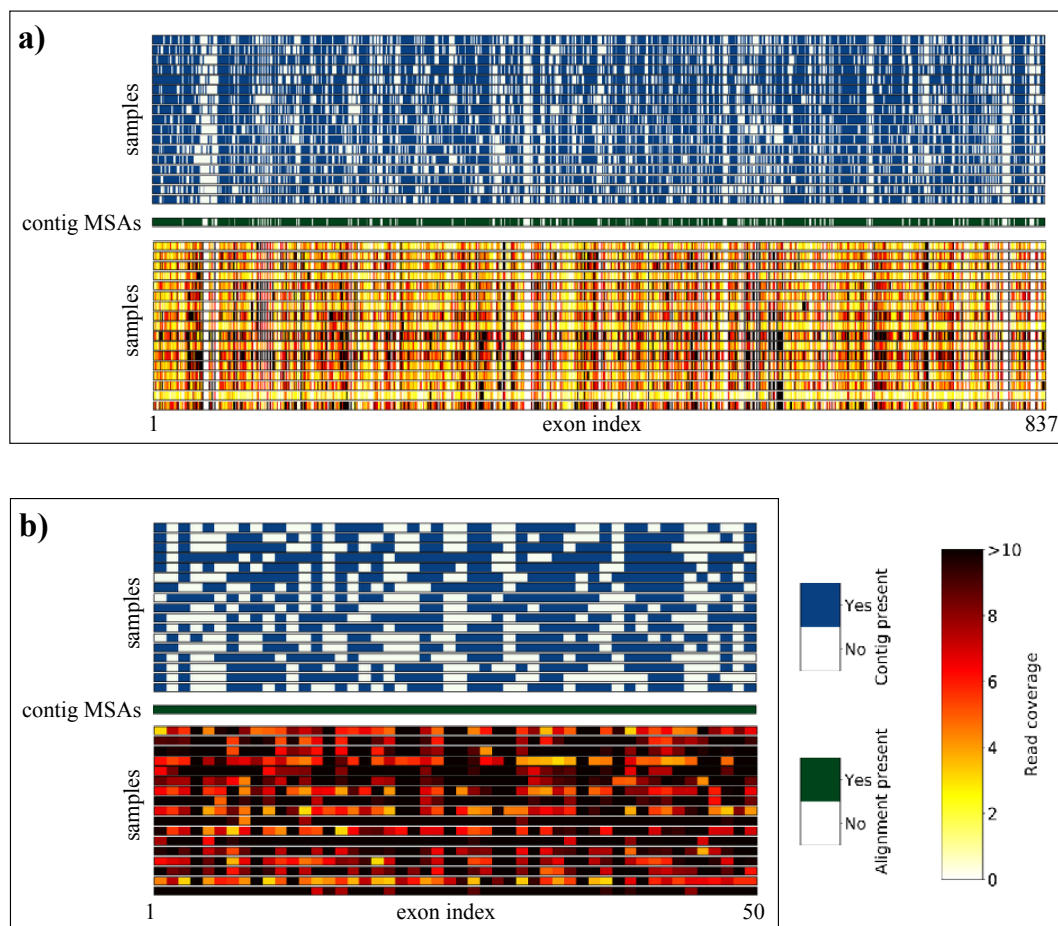


456

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



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



478

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