

RESOURCE ARTICLE

OPEN ACCESS



A Snakemake Toolkit for the Batch Assembly, Annotation and Phylogenetic Analysis of Mitochondrial Genomes and Ribosomal Genes From Genome Skims of Museum Collections

Oliver W. White | Andie Hall | Ben W. Price | Suzanne T. Williams | Matthew D. Clark

The Natural History Museum, London, UK

Correspondence: Suzanne T. Williams (s.williams@nhm.ac.uk) | Matthew D. Clark (matt.clark@nhm.ac.uk)**Received:** 5 April 2024 | **Revised:** 10 September 2024 | **Accepted:** 1 October 2024**Handling Editor:** Alana Alexander**Funding:** The authors received no specific funding for this work.**Keywords:** assembly | bioinformatics | DNA barcode | genome skim | iBOL | museomics | phylogenetics | snakemake

ABSTRACT

Low coverage ‘genome-skims’ are often used to assemble organelle genomes and ribosomal gene sequences for cost-effective phylogenetic and barcoding studies. Natural history collections hold invaluable biological information, yet poor preservation resulting in degraded DNA often hinders polymerase chain reaction-based analyses. However, it is possible to generate libraries and sequence the short fragments typical of degraded DNA to generate genome-skims from museum collections. Here we introduce a snakemake toolkit comprised of three pipelines *skim2mito*, *skim2rrna* and *gene2phylo*, designed to unlock the genomic potential of historical museum specimens using genome skimming. Specifically, *skim2mito* and *skim2rrna* perform the batch assembly, annotation and phylogenetic analysis of mitochondrial genomes and nuclear ribosomal genes, respectively, from low-coverage genome skims. The third pipeline *gene2phylo* takes a set of gene alignments and performs phylogenetic analysis of individual genes, partitioned analysis of concatenated alignments and a phylogenetic analysis based on gene trees. We benchmark our pipelines with simulated data, followed by testing with a novel genome skimming dataset from both recent and historical solariellid gastropod samples. We show that the toolkit can recover mitochondrial and ribosomal genes from poorly preserved museum specimens of the gastropod family Solariellidae, and the phylogenetic analysis is consistent with our current understanding of taxonomic relationships. The generation of bioinformatic pipelines that facilitate processing large quantities of sequence data from the vast repository of specimens held in natural history museum collections will greatly aid species discovery and exploration of biodiversity over time, ultimately aiding conservation efforts in the face of a changing planet.

1 | Introduction

Natural history collections are home to more than 1 billion, expert-verified specimens worldwide (Bartolozzi, Bettison-Varga, and Chernetsov 2023) as well as large numbers of unclassified and bulk samples, and as such represent a vast repository

of historical biological data that remains largely untapped as genetic resources. Challenges associated with such material include poor preservation, the use of unknown preservatives, and ongoing DNA degradation and contamination. Despite these challenges, early studies were able to isolate DNA from museum collections of extinct species. Notably, Higuchi et al. (1984)

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2024 The Author(s). *Molecular Ecology Resources* published by John Wiley & Sons Ltd.

isolated and sequenced short mitochondrial DNA sequences from a 140-year-old museum collection of the quagga, an extinct subspecies of the zebra, using a plasmid cloning approach. In addition, Krause et al. (2006) sequenced the entire mitochondrial genome of the woolly mammoth using multiplex polymerase chain reactions (PCRs). The degraded nature of DNA from museum collections can make direct recovery of intact gene sequences by PCR impossible, without the amplification of many short overlapping fragments (D'Ercole, Prosser, and Hebert 2021). Fortunately, advances in novel laboratory techniques (Ruane and Austin 2017; Straube et al. 2021) and next generation sequencing (NGS) technology make it possible to obtain DNA sequences from many historical specimens, unlocking the potential for wide-ranging genomic analyses. Using natural history collections provides the opportunity to easily work on many species, even if they are now extinct, rarely collected, or from areas of the world that are poorly sampled. Given that almost all known species have vouchers in one or more natural history collections, the use of these specimens could rapidly fill gaps in DNA reference libraries, greatly accelerating biodiversity discovery and DNA-based monitoring of the environment.

'Genome skimming' is a term referring to the generation of low coverage whole genome sequence data, first coined by Straub et al. (2012). Although genome skimming does not generate data with sufficient coverage to assemble the entire nuclear genome, there are sufficient reads to assemble sequences that are present in the genome in multiple copies and are therefore highly represented in the sequence data. Common targets for genome skimming studies include organelle genomes (a cell has one nucleus but many organelles) and nuclear ribosomal genes (there are typically 100s or more copies of nuclear rRNA genes). Organelle and ribosomal genes have been widely used for phylogenetics due to their discriminatory power and the availability of 'universal' primers and continue to be employed as 'barcode' genes for identification. Specifically, the mitochondrial gene *cox1* for animals, chloroplast genes *matK* and *rbcL* for plants, 16S rRNA for bacteria and 18S or ITS for fungi are used now and dominate barcode sequence databases. In addition, novel genome skimming approaches are increasingly being developed, to use all sequence data to assign species identity via a 'DNA-mark' (Bohmann et al. 2020) or 'varKode' (De Medeiros et al. 2024).

When working with historical specimens in particular, genome skimming offers many advantages over PCR amplification and Sanger sequencing of individual genes. Relatively long fragments of intact genomic DNA are required for PCR, whereas degraded and low amounts of DNA (~1 ng) typical of museum specimens (Mullin et al. 2023), are also suitable for short read NGS platforms. The wet lab work is now relatively straightforward for a skilled molecular biologist using ancient DNA protocols developed over many years, with key components involving DNA extraction and library methods optimised for degraded DNA. Genome skimming also has additional benefits over targeted PCR since multiple loci can be recovered at the same time without development and optimisation of multiple PCR assays, or the need to design group specific DNA capture baits for target enrichment (Call et al. 2023). With advances in bioinformatic tools, it is likely that low coverage genome skimming datasets will have even greater utility

in the future. For example, K-mer based approaches have been developed for genome skims to investigate phylogenetic relationships (Sarmashghi et al. 2017) and genome properties including genome length and repetitiveness (Sarmashghi et al. 2021). Finally, genome skimming is increasingly cost effective as the cost of NGS sequencing continues to decrease. In the light of these advantages, genome skimming is seen as a hugely scalable process that is suitable for batch recovery of useful genomic data from museum collections.

However, few bioinformatic pipelines are available to assist with the assembly of large numbers of organelle and nuclear ribosomal sequences from batches of genome skimming data. Notable exceptions include MitoZ (Meng et al. 2019) and NOVOWrap (Wu et al. 2021) for the assembly and annotation of mitochondrial genomes. MitoGeneExtractor can be used to extract mitochondrial protein coding genes from NGS datasets (Brasseur et al. 2023). In addition, plastomatic (W. Chen, Achakkagari, and Strömvik 2022) is available for chloroplast assembly and annotation and PhyloHerb (Cai, Zhang, and Davis 2022) can be used for the assembly of chloroplast and nuclear ribosomal repeats without annotation. However, these tools were not designed with historical and/or degraded samples in mind nor issues such as contamination and the undesirable assembly of non-target sequences. In addition, these tools do not implement phylogenetic analysis of the annotated genes identified. Other targeted assembly approaches are available including Orthoskim (Pouchon et al. 2022) for chloroplast, mitochondrial and ribosomal sequences, and this pipeline has been benchmarked with a large genome skimming dataset, including libraries created from herbarium collections. Despite the presence of other tools, there is a gap for methods utilising a pipeline approach that assemble and annotate multicopy sequences from museum collections in a repeatable, portable, scalable and bioinformatically robust way. Utilising a snakemake framework allows for a repeatable and robust pipeline.

Here we introduce a snakemake toolkit comprised of three pipelines: *skim2mito*, *skim2rrna* and *gene2phylo*. These pipelines are designed to unlock the potential of historical museum specimens when using genome skimming. Specifically, *skim2mito* and *skim2rrna* perform the batch assembly, annotation and phylogenetic analysis of mitochondrial genomes and nuclear ribosomal genes, respectively. The third pipeline, *gene2phylo*, takes a set of gene alignments and performs phylogenetic analysis of individual genes, partitioned analysis of concatenated alignments and a phylogenetic analysis based on gene trees. The pipelines wrap 12 published bioinformatic tools as well as custom Python and R scripts into a single user-friendly pipeline designed to cope with more challenging data from historical collections, permitting large scale genome skimming studies from museum specimens. These pipelines have many advantages for such studies, for example they (1) run on a single machine or in parallel on a High Performance Computing cluster, (2) process batches of samples, (3) assemble both mitochondrial and nuclear ribosomal sequences using GetOrganelle (Jin et al. 2020), (4) perform basic assembly checking, for example, for contamination and non-target sequences and (5) generate phylogenetic gene trees based on annotated genes. GetOrganelle was selected for sequence assembly because it can be used for mitochondrial,

chloroplast and nuclear genes. In addition, an independent systematic comparison (Freudenthal et al. 2020), highlighted that this tool was the best performing assembly method in their tests. Whilst it would be possible to wrap all steps into one single pipeline, it is necessary to implement mitochondrial and ribosomal assembly individually, and user feedback suggested it was necessary to check outputs manually, for example, to check for possible contamination, before more detailed phylogenetic analysis (see Section 3 and Section 4 for more details). Note that as the pipelines are written in Snakemake, a relatively accessible option for writing pipelines, it is possible for users to adjust parameters and steps in the pipeline as necessary.

To benchmark the utility of our pipelines, we generated a simulated data for 25 species from Papilionoidea, a superfamily of butterflies with reference genomes and known taxonomy. We then analysed a novel genome skimming dataset for the gastropod family Solariellidae (hereafter solariellid gastropods). This group was selected as it represents many of the challenges associated with genome skimming museum collections. Solariellids are small marine snails found predominantly in deep-water. Many member species are rare, and as a family they are poorly represented in museum collections worldwide, with few live-collected specimens: many species are known only from a single, dry and often damaged shell (Williams et al. 2020). Although solariellid gastropods have been the focus of previous molecular phylogenetic studies (Sumner-Rooney et al. 2016; Williams et al. 2013, 2022), these studies have relied on partial sequence from only four genes, which have not fully resolved relationships among genera. As such, our understanding of solariellid evolution would greatly benefit from increasing the number of gene sequences used, furthermore there are no published reference genomes for the group with limited sequence data on public databases. Here we demonstrate how a genome skimming approach and our snakemake toolkit can be used to improve our understanding of phylogenetic relationships even in this challenging case.

2 | Materials and Methods

2.1 | Simulated Data

Simulated data for mitochondrial and ribosomal sequences were generated for 25 species from the butterfly superfamily Papilionoidea. Simulated data were generated using gargamel (Renaud et al. 2017) with the following parameters: read length 150, location 5, scale 0.5, target coverage 20×, sequencing system HiSeq2500 and no contaminant sequences. A full list of taxa and reference datasets used to generate the simulated data are presented in Table S1. Note that the simulated data and config files for this analysis are available with all pipelines as test data.

2.2 | Solariellid Sample Selection and Sequencing

A total of 25 samples were selected, with representatives from 18 genera, encompassing the diversity of the solariellid family (Table 1; Figure 1). Samples differ in several ways that likely

affected DNA quality and yield (Table S2), for example, time since collection (1967–2015) and preservation method (dry shell with dehydrated body tissues or live-collected snail preserved in 70%–99% ethanol). In addition, some shells were cracked, allowing the rapid penetration of ethanol, which is particularly important as snails can seal their bodies inside their shells by closing their operculum, effectively excluding ethanol. Samples where shells have not been cracked generally have very degraded DNA. Samples also differ in time between sequencing and when DNA was extracted (0–10 years; Table S2). DNA was isolated using Qiagen DNeasy blood and tissue kit and quantified using a Qubit fluorimeter and High Sensitivity assay kit. A TapeStation 2200 (Agilent Technologies, Santa Clara, USA) was also used to assess DNA integrity prior to library preparation. We did not use our specialist ancient DNA clean room for this work, which uses many additional steps to limit contamination (Fulton and Shapiro 2019). We expected contamination from prior handling in collection and curation, with which we could investigate the ability of the pipeline to detect contaminants and perform acceptably despite them. PCR amplification and Sanger sequencing of mitochondrial (*cox1*, 16S and 12S) and ribosomal genes (28S) were attempted for each sample at the time of DNA extraction and these results are compared with our genome skimming approach. Illumina Libraries were prepared using a SparQ DNA Frag and Library Prep kit (QuantaBio, Beverly, USA) and sparQ PureMag Beads (QuantaBio), with SparQ Adaptor Barcode sets A and B (QuantaBio). Libraries were normalised and pooled equally before being sent to Novogene (Cambridge, UK) for sequencing. The single indexed libraries were sequenced on an Illumina Novaseq on an S4 300 cycle flowcell using 150bp paired reads (see Data and code availability statement).

Additional sequence data for '*Solariella*' *varicosa* were provided by Andrea Waeschenbach (Natural History Museum London, UK.). Raw sequence data for two outgroups from the family Turbinidae were also analysed, including: *Turbo cornutus* (Kim et al. 2022; SRR15496837) and unpublished raw data for *Lunella* aff. *cinerea* (mitochondrial genome published in Williams, Foster, and Littlewood 2014). These outgroup sequences provide the possibility of comparing published mitochondrial genomes for *Turbo cornutus* (National Center for Biotechnology Information [NCBI] GenBank accession NC_061024.1) and *Lunella* aff. *cinerea* (KF700096.1) with the results from our pipeline using the same raw sequence data. Specifically, a blast search was implemented to compare sequence homology. In addition, sequence depth variation, GC content and repeat content were visualised using custom Circos plots (Krzywinski et al. 2009; https://github.com/o-william-white/circos_plot_organelle; accessed 08/2023).

2.3 | Human Contamination

Prior to running the pipelines, the extent of human contamination was evaluated. Raw reads were adapter trimmed and quality filtered with fastp (S. Chen et al. 2018) and mapped to the human reference genome (NCBI GRCh38.p14) using bwa mem (Li 2013) and the number of mapped reads quantified with samtools (Li et al. 2009).

TABLE 1 | Sample details for 25 solariellid gastropod species and two outgroup species used in this study with museum registration numbers or NCBI sequence read archive number for sequence data (*Turbo cornutus* only), ocean of origin, latitude, longitude and depth of collection location.

| Species | Specimen voucher ^a | Ocean | Latitude | Longitude | Depth (m) |
|--|-------------------------------|----------------------------|----------|-----------|------------|
| <i>Archiminolia oleacea</i> | AMS C.133269 | Indo-West Pacific | −24.375 | 153.285 | 192–229 |
| <i>Arxellia herosae</i> | MNHN-IM-2009-28739 | Indo-West Pacific | −24.717 | 168.167 | 298–324 |
| <i>Bathymophila gravida</i> | NMNZ M.299691 | Indo-West Pacific | −36.146 | 178.202 | 712–924 |
| ' <i>Bathymophila</i> ' sp. 18 | MNHN-IM-2009-23103 | Indo-West Pacific | −22.317 | 171.333 | 925 |
| <i>Bathymophila</i> -Like sp. 12 | MNHN-IM-2009-28741 | Indo-West Pacific | −19.667 | −178.167 | 314–377 |
| <i>Chonospeira nuda</i> | SMNH 127100 | North East Pacific | 36.367 | −122.417 | 999 |
| Clade D sp. d | MNHN-IM-2013-59648 | Indo-West Pacific | 22.050 | 119.067 | 1306–1756 |
| <i>Elaphriella wareni</i> | MNHN-IM-2013–45,837 | Indo-West Pacific | −8.617 | 151.783 | 705–817 |
| <i>Ilanga whitechurchi</i> | NMSA W9631 | South West Indian Ocean | −33.167 | 28.033 | 90 |
| <i>Lamellitrochus</i> sp. 6 | MNHN-IM-2013-60491 | Caribbean | 16.350 | −60.900 | 111–162 |
| ' <i>Lamellitrochus</i> ' <i>carinatus</i> | MNHN-IM-2009-31169 | Caribbean | 16.360 | −61.579 | 29 |
| <i>Microgaza rotella</i> | MNHN-IM-2013-8023 | Caribbean | 16.400 | −61.550 | 130 |
| <i>Phragmomphalina tenuiseptum</i> | NMNZ M299700 | Indo-West Pacific | −31.867 | 172.433 | 780–790 |
| <i>Solariella amabilis</i> | NHMUK 20180166 | North Atlantic | 62.191 | 5.567 | 150–200 |
| <i>Solariella</i> sp. 7 | MNHN-IM-2019-12000 | Indo-West Pacific | −24.800 | 168.150 | 250–270 |
| ' <i>Solariella</i> ' <i>carvalhoi</i> | MNHN-IM-2013-61297 | Caribbean | 15.800 | −61.467 | 379–428 |
| ' <i>Solariella</i> ' <i>obscura</i> | NHMUK 20230529 | North Atlantic | 69.803 | 30.693 | 04-Dec |
| ' <i>Solariella</i> ' <i>varicosa</i> | NHMUK 20120235 | North Atlantic | 70.067 | 29.200 | 10–174 |
| <i>Spectamen</i> cf. <i>bellulum</i> | NHMUK 20110452 | Indo-West Pacific | −26.943 | 153.404 | 31 |
| ' <i>Spectamen</i> ' <i>franciscanum</i> | NMSA V1091 | South West Indian Ocean | −34.783 | 23.983 | 171 |
| <i>Suavotrochus lubricus</i> | MNHN-IM-2013-61096 | Caribbean | 16.033 | −61.233 | 266–388 |
| ' <i>Suavotrochus</i> ' sp. 2 | MNHN-IM-2013-61502 | Caribbean | 15.783 | −61.200 | 550–562 |
| ' <i>Zetela</i> ' <i>alphonsi</i> | SMNH 10387 | South East Pacific | −36.361 | −73.725 | 865 |
| <i>Zetela kopua</i> | NMNZ M.131532 | Indo-West Pacific | −45.403 | 173.980 | 1386 |
| <i>Zetela textilis</i> | NMNZ M.035478 | Indo-West Pacific | −42.637 | 176.283 | 256–311 |
| Outgroups | | | | | |
| <i>Lunella</i> aff. <i>cinerea</i> | NHMUK 20100448 | Indo-West Pacific | −12.554 | 130.876 | Intertidal |
| <i>Turbo cornutus</i> | SRR15496837 | Indo-West Pacific | 33.454 | 126.949 | Unknown |

Abbreviations: (Specimen voucher) AMS, Australian Muséum; MNHN, Muséum National D'histoire Naturelle; MNSA, KwaZulu-Natal Museum; NHMUK, Natural History Museum, London; NMNZ, Museum of New Zealand Te Papa Tongarewa; SMNH, Swedish Museum of Natural History.

^aNames correspond to those used in previous studies (Williams et al. 2020, 2022). Inverted commas around generic names indicates uncertainty about generic assignment based on this or previous studies. Previously published data for *Turbo cornutus* (Kim et al. 2022) and *Lunella* aff. *cinerea* (Williams, Foster, and Littlewood 2014) were also included in this study.

2.4 | Pipeline Descriptions

2.4.1 | skim2mito

As input, the *skim2mito* pipeline requires two files to be provided by the user: (1) a config file (in YAML format) and (2) a sample list file (in CSV format). The config file outlines the

main parameters including the GetOrganelle reference method, adapter sequences used, reference databases used for annotation, alignment trimming methods and outgroup samples for the phylogenetic analyses. The samples.csv file is a list of the samples included in the analysis with the sample names, paths to forward and reverse reads, NCBI taxonomy IDs for searches of reference sequences on NCBI or paths to manually generated

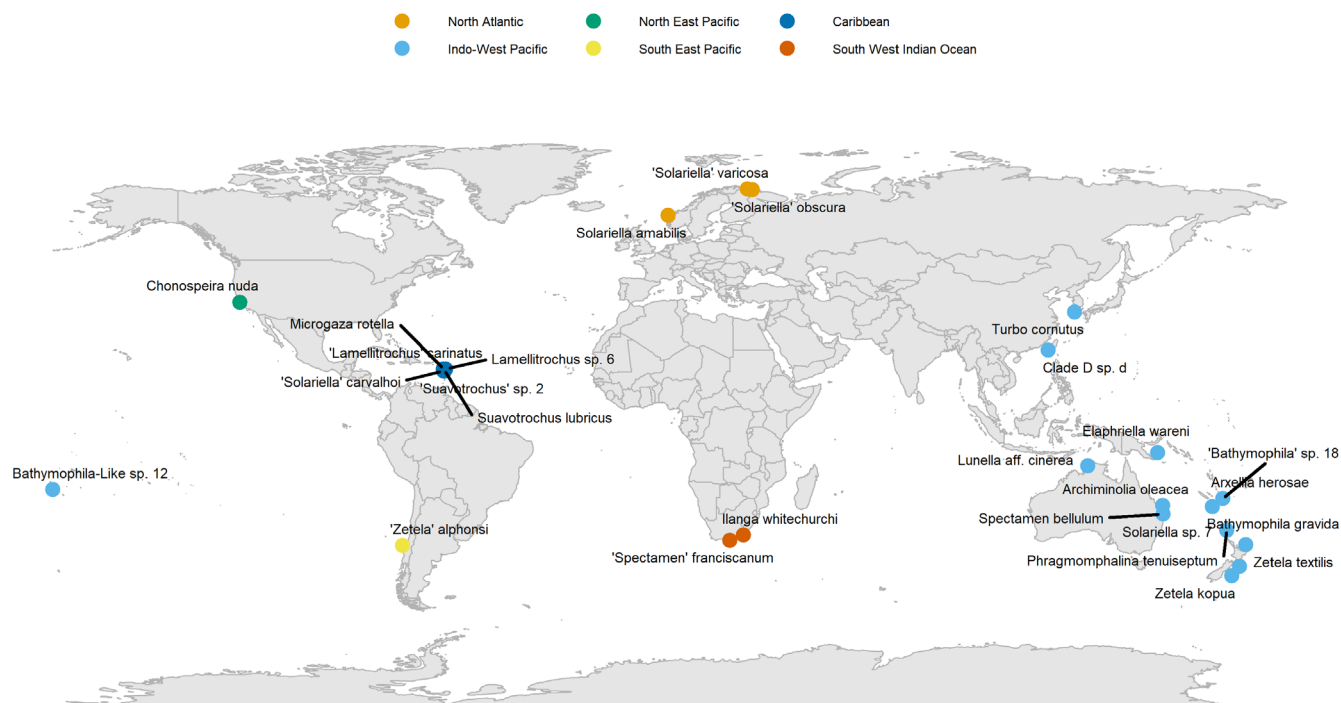


FIGURE 1 | Map showing collection localities for solariellid gastropods samples used in this study.

gene and seed databases required by GetOrganelle. Note that example config.yaml and samples.csv files are provided with all pipelines (see data and code availability statement). The pipeline accepts NGS data from short read platforms (e.g., Illumina) in demultiplexed paired fastq format.

The *skim2mito* pipeline (Figure 2) starts by processing the data from each sample using fastp (Chen et al. 2018) to detect and remove adapter sequences, trim low-quality sequences with parameters for forward and reverse adapter sequences and optionally duplication (--dedup) specified. Fastqc is implemented on raw and quality filtered reads (Andrews 2010). GetOrganelle (Jin et al. 2020) is then used to assemble the target sequence of interest, using seed and gene reference databases to identify and assemble target reads. Although GetOrganelle is provided with default seed and gene databases, our initial benchmarking highlighted that using custom reference databases from closely related taxa minimised the likelihood of assembling contaminant sequences. Therefore, the pipeline will either generate a seed and gene database using the python script go_fetch.py version 1.0.0 (https://github.com/o-william-white/go_fetch), or it will use custom databases provided by the user. The go_fetch.py script takes a NCBI taxonomy ID provided by the user in the samples.csv file, searches the NCBI nucleotide database for mitochondrial reference sequences that are as close as possible to the target taxonomy by working back up the NCBI taxonomy hierarchy until sufficient references are found (minimum 5, maximum 10), and then downloads and formats the reference data for use with GetOrganelle. GetOrganelle is implemented with the following additional parameters: --reduce-reads-for-coverage inf -max-reads inf -R 20. Sequences assembled by GetOrganelle are typically named based on the output of SPAdes (Prjibelski et al. 2020), which can produce long sequence names. Therefore, sequences are renamed to <sample_name>_contig<n> if there

are multiple contigs or <sample_name>_circular if a single circular sequence is found. Note that GetOrganelle can produce more than one assembled sequence where there are different possible paths through the same assembly graph, for example, mitochondrial genomes containing repeats. However, the pipeline simply selects the first assembled sequence for downstream analyses as the main outputs are the annotated gene sequences and the correct orientation of repeat regions is not necessary. Note that GetOrganelle also generates an assembly graph (.GFA) which can be viewed with bioinformatic tools including Bandage (Wick et al. 2015). Basic assembly statistics are summarised using SeqKit (Shen et al. 2016).

Next, the assembly quality is evaluated using a blastn search (Camacho et al. 2009) against a temporary blast database created as part of the workflow. For mitochondrial sequences, a blast database generated from the NCBI mitochondrion RefSeq database is used (<https://ftp.ncbi.nlm.nih.gov/refseq/release/mitochondrion/>). Quality filtered reads generated by fastp are mapped to the assembled sequence using minimap2 to estimate sequence coverage (Li 2018). The blast and mapping outputs are summarised using blobtools (Laetsch, Blaxter, and Leggett 2017) and the likely taxonomy of the assembled sequence is defined using the taxrule 'bestsumorder'. Following the assembly quality check, assembled sequences are annotated using MITOS2 (Bernt et al. 2013). Following assembly and annotation, a plot is created using a custom python script to visualise the location of annotated genes, coverage and proportion of mismatches in mapped reads.

Once the sequences are assembled and annotated, the checkpoint function of snakemake is used to recover all annotated gene sequences assembled across samples (Figure 2). Note that only protein coding mitochondrial genes are used by *skim2mito* for downstream analyses, with tRNA sequences excluded, as

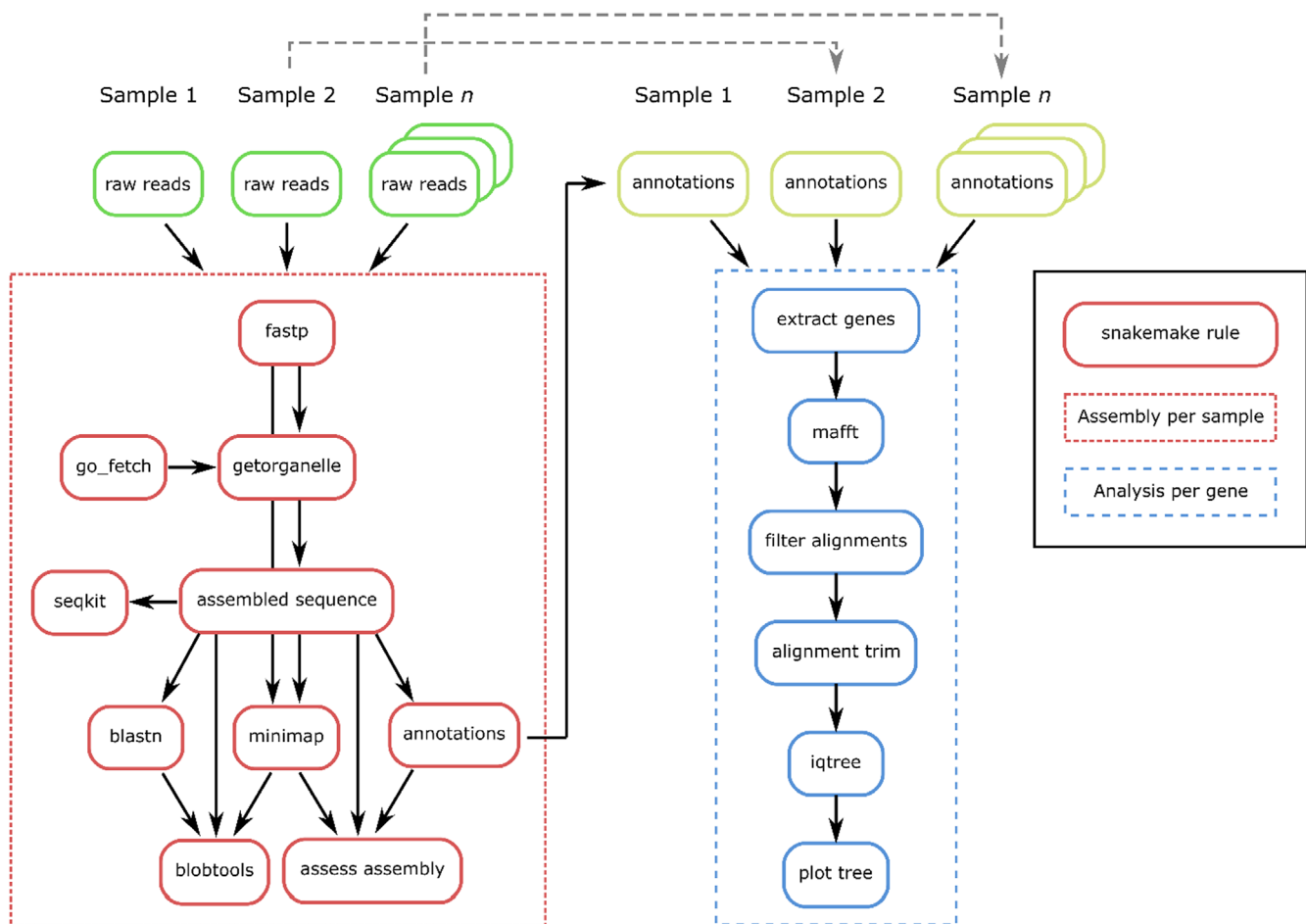


FIGURE 2 | Schematic diagram of *skim2mito* and *skim2rrna*. Both pipelines initially process raw read data from individual samples, assembling target sequences and assessing assembly quality. Each pipeline then aligns, processes and implements phylogenetic analysis for all annotated genes found across assembled sequences. The *skim2mito* and *skim2rrna* pipelines follow a similar workflow, except for GetOrganelle parameters, and homology search and annotation reference databases.

protein coding sequences were typically longer and have more robust alignments. For each annotated gene recovered, mafft (Kato and Standley 2013) is used to align sequences with the following parameters: `--maxiterate 1000 --globalpair --adjustdirectionaccurately`. To avoid the inclusion of annotated gene sequences with largely missing data in the alignment, individual sample sequences with $\geq 50\%$ missing data relative to the alignment are removed. Note that the threshold for missing data can be adjusted by the user in the config.yaml input, which may be of particular value when working with very rare specimens. The alignments are trimmed using either Gblocks (Castresana 2000) or Clipkit (Steenwyk et al. 2020) as specified in the config.yaml file. Phylogenetic analysis is then implemented with IQ-TREE 2 (Minh et al. 2020) using 1000 ultrafast bootstraps and consensus trees are plotted in R using the ggtree package (R Core Team, 2020; Yu et al. 2017). Note that phylogenetic analysis is only implemented if there are at least five sequences in the alignment.

2.4.2 | skim2rrna

The *skim2rrna* pipeline (Figure 2) requires the same input data and follows similar steps to the *skim2mito* pipeline described

above, except for the parameters and tools used for the assembly, homology search and annotation. Specifically, for the assembly with GetOrganelle, the following parameters are used: `-F anonym -reduce-reads-for-coverage inf -max-reads inf -R 10 -max-extending-len 100 -P 0`. For the homology search with blast, a blast database generated from the SILVA 138 database is used (Quast et al. 2013). Finally, barrnap (<https://github.com/tseemann/barrnap>) is used for annotation of ribosomal sequences.

2.4.3 | Assessing Assembly and Annotation Results

After running *skim2mito* and/or *skim2rrna*, it is necessary to check the results for evidence of contamination across samples or individual genes. This is especially important for museum samples. Specifically, the blobtools (Laetsch, Blaxter, and Leggett 2017) output should be checked for sequences with unusual (e.g., taxonomically divergent) blast hits and the gene alignments and gene trees reviewed by taxonomic experts to identify incongruent relationships. Users may check the summary counts of genes recovered across samples, which can be useful for identifying samples with large amounts of missing data for downstream analyses. Putative contaminant

sequences, genes with large amounts of missing sequences in the alignment or samples with large amounts of missing genes can then be removed using the supplementary python script *format_alignments.py* which removes sequences from alignments based on sequence names and formats alignments for use with *gene2phylo*.

2.4.4 | gene2phylo

After running *skim2mito* and *skim2rrna* and removing putative contaminants, the user may wish to implement further phylogenetic analyses of the filtered alignments. To assist with this, a final pipeline *gene2phylo* is provided to reanalyse assembled genes (Figure 3). As input, the *gene2phylo* pipeline only requires a config file to specify the main parameters for the phylogenetic analysis. The *gene2phylo* pipeline accepts

input data from multiple individual genes (aligned or unaligned) in a single directory. The alignments must be named after the gene names (e.g., 'cox1.fasta'), and the sample names used for sequences across alignments must be consistent. The *gene2phylo* pipeline begins by optionally re-aligning, removing individual sample sequences with $\geq 50\%$ missing data relative to the alignment, and trimming poorly aligned regions of alignments, using a similar approach to the pipelines described above. Note, this is only necessary if the input alignments were edited for the removal of contaminant sequences, which is likely to influence alignment characteristics. For each input alignment, phylogenetic analysis is implemented with IQ-TREE2 (Minh et al. 2020) using 1000 ultrafast bootstraps. Input alignments are also combined into a partitioned alignment and a partitioned phylogenetic analysis using IQ-TREE 2 is implemented with 1000 ultrafast bootstrap replicates. In addition, individual gene trees are used to infer phylogenetic

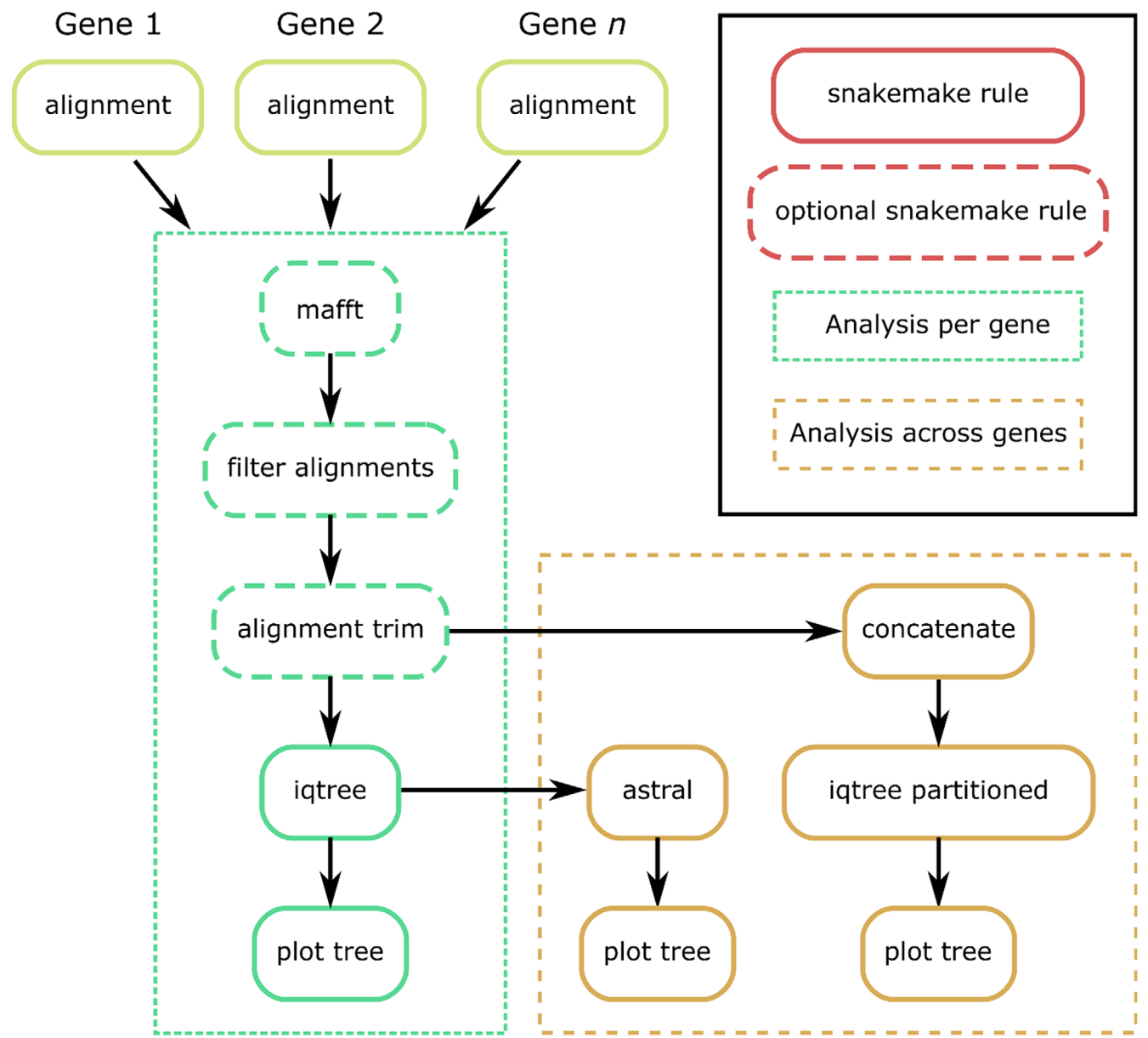


FIGURE 3 | Schematic diagram of *gene2phylo*. Gene2phylo can optionally re-align, remove sequences with too much missing data and trim poorly aligned regions of the alignments. *Gene2phylo* implements IQTREE 2 phylogenetic analysis for each annotated gene, an IQTREE 2 partitioned analysis for all assembled genes and an astral analysis across all individual gene trees.

relationships using astral (Zhang et al. 2018). Phylogenetic trees from each analysis are plotted using the ggtree package (R Core Team 2020; Yu et al. 2017).

3 | Results

3.1 | Simulated Data

The *skim2mito* and *skim2rrna* pipeline recovered mitochondrial and ribosomal sequences for all simulated datasets, with a mean assembly size of 15,068 and 6,780 for mitochondrial and ribosomal sequences respectively. The blobtools blast hits predicted the correct taxonomic family for all assembled mitochondrial sequences. However, the correct taxonomic family was only predicted for seven assembled ribosomal sequences out of the 25 simulated datasets. This is likely due to the lack of annotated ribosomal sequences available for the butterfly superfamily Papilionoidea in the SILVA reference data used for the blast database. All annotated genes were retained for phylogenetic analysis with *gene2phylo* and the relationships identified by both the partitioned IQTREE 2 and astral analyses conformed with known subfamily relationships (Figures S1 and S2).

3.2 | Solariellid Data

The results of PCR amplification targeting four genes (*cox1*, 28S, 12S and 16S) of our solariellid specimens suggest that many DNA samples were highly degraded (Table S2). In some cases, faint bands were observed when PCR products were visualised on agarose gels, but clean Sanger sequence could not be obtained, because of low yield and noisy background. DNA quality was confirmed by recording the DNA Integrity Number (DIN) for samples, with a DIN of 10 indicating highly intact DNA fragments, whilst a DIN of 1 indicates a highly degraded DNA sample. DNA quality for the samples used in this study ranged from not detectable for the poorest samples to 6.5 for the best (Table S2).

Approximately 870M raw sequence reads were generated across all samples, with an average of 32M raw reads per sample (Table S2). On average, 95.62% of reads were retained following adapter removal and basic filtering with fastp. Sequence data contamination with reads originating from human DNA was extensive across samples, with an average of 61.86% (range 32.74%–80.02%; Table S2).

The *skim2mito* pipeline successfully recovered mitochondrial genome sequences from 24/27 samples with an average assembly size of 14,441 bp (range 347–24,670 bp). A circular mitochondrial genome was assembled for a single sample (*Zetela kopua*; Figure 4). However, no mitochondrial sequences could be assembled for *Elaphriella wareni*, '*Spectamen*' *franciscanum* or *Zetela textilis*. Assembled sequences for outgroups *Turbo cornutus* and *Lunella* aff. *cinerea* were compared to previously published sequences. Blast hits for these two specimens had 100% percentage sequence similarity when compared to previously published sequences, although the sequences assembled by our pipeline were not as complete as those published on GenBank. The assembled

sequence for *Turbo cornutus* matched the published sequence (NCBI accession NC_061024.1; Figure S3) from 1 to 13,676 and 14,107 to 17,299. The assembled sequence for *Lunella* aff. *cinerea* matched the published sequence (KF700096.1; Figure S4) from 1 to 13,973 and 14,412 to 17,670. Visualisation of circos plots for published circular sequences suggests that the assembly process broke in regions of low coverage, low GC content and/or high repeat content (Figures S3 and S4). Of the 15 mitochondrial genes annotated by MITOS2 (13 protein coding genes and two mitochondrial ribosomal subunits), an average of 11 genes were annotated across samples, with 15 of 27 samples having all protein coding and rRNA genes annotated.

The *skim2rrna* pipeline successfully recovered ribosomal gene sequences from 27 of 27 samples with an average size of 3,049 bp. Of the ribosomal genes annotated by barnnap, the 18S, 28S and 5.8S rRNA genes were annotated in 4, 24 and 4 samples respectively.

After checking the outputs of *skim2mito* and *skim2rrna*, 18S and 5.8S rDNA genes were found to have more than 50% missing data across samples and were therefore removed from downstream analyses. In addition, *atp8* was discarded as it was < 100 nucleotides in length and contained little phylogenetic information. Six samples were found to be missing more than 50% of the annotated genes and were therefore discarded from all downstream phylogenetic analysis. Excluded species were *Archiminolia oleacea*, *Arxellia herosae*, *Bathymophila grvida*, *Elaphriella wareni*, '*Spectamen*' *franciscanum*, and *Zetela textilis*. After manual checking of the blobtools output, alignments and phylogenetic trees generated by *skim2mito* and *skim2rrna*, it was determined that although all genes were annotated for *Spectamen* cf. *bellulum*, these data likely originated from another solariellid species as a contaminant, from the genus *Ilanga*. Annotations for 28S from *Phragmomphalina tenuiseptum* and *Zetela kopua* were also identified as likely non-solariellid gastropod contaminants. Duplicate gene annotations were identified in Clade D sp. d (*cox3*, *nad3*, *nad2*, *cox1*, *cox2*, *atp8*, *atp6*), *Solariella amabilis* (*nad3*) and *Zetela alphonsi* (*nad2*). Where duplicate genes were identified, only the first annotation was used in downstream analyses. Following the removal of genes and samples with too much missing data, contaminant sequences and duplicate annotations, the final dataset consisted of 15 genes for 20 specimens with an average 6.33% missing data on average across samples.

With the filtered set of genes, *gene2phylo* was implemented to re-align and analyse the individual gene alignments. Phylogenetic analysis of the partitioned alignment in IQ-TREE (Figure 5) recovered a tree with support values ranging from poor to optimal (29%–100%). Phylogenetic analysis using the individual genes trees with astral (Figure S5) recovered a tree with broadly consistent topology.

4 | Discussion

This study demonstrates the utility of a snakemake toolkit comprised of the pipelines *skim2mito*, *skim2rrna* and *gene2phylo*, for the assembly, annotation, and phylogenetic analysis of mitochondrial and ribosomal genes from genome skimming datasets. Analysis of simulated data from the butterfly superfamily

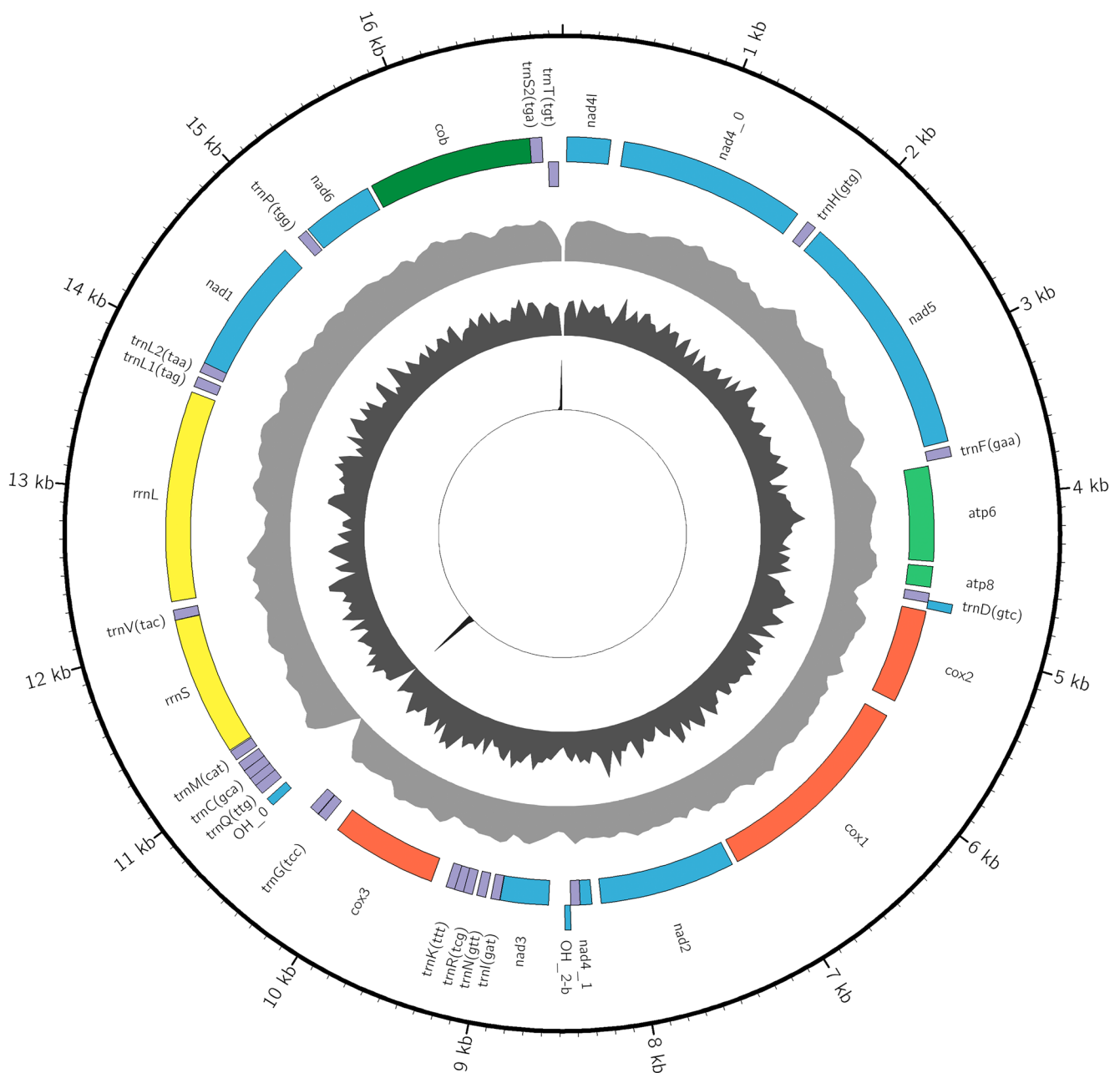


FIGURE 4 | Assembled circular sequence for *Zetela kopua* with the following attributes from outside to inside: Sequence position, annotation names, annotations on the +strand, annotations on the -strand, coverage (max=2779), GC content (max=0.6) and repeat content (max=1.0). This image was created using a custom organelle visualisation tool available on GitHub (https://github.com/o-william-white/circos_plot_organelle; accessed 08/2023).

Papilionoidea is used to benchmark our pipelines, generating expected phylogenetic relationships. Further analysis of novel sequence data from the gastropod family Solariellidae, generated the first mitochondrial genomes and new ribosomal sequences for the family. Specifically, complete or partial mitochondrial genomes were obtained for 23 of 27 non-contaminated specimens, and ribosomal sequences were assembled for 24 of 27 samples.

The Solariellidae samples included in this study represent many of the issues that are typical of historical museum specimens. For example, three samples (*Solariella* *obscura*, *Solariella* *varicosa* and *Solariella* *amabilis*) were collected more than 50 years ago and preserved in low percentage

ethanol (70%) with uncracked shells. In addition, DNA was extracted more than 10 years ago from dehydrated tissue samples of another sample (*Bathymophila*-Like sp. 12). Therefore, several of the Solariellidae samples had highly degraded DNA (DIN < 2). Despite this, our pipelines were able to recover nearly complete mitochondrial genomes and ribosomal sequences for most samples.

Contamination is a common issue with historical specimens. Indeed, human contamination was extensive across all the Solariellidae samples sequenced in this study, with an average of 61.86% human reads per sample (Table S2). This is likely due to the presence of only small amounts of highly

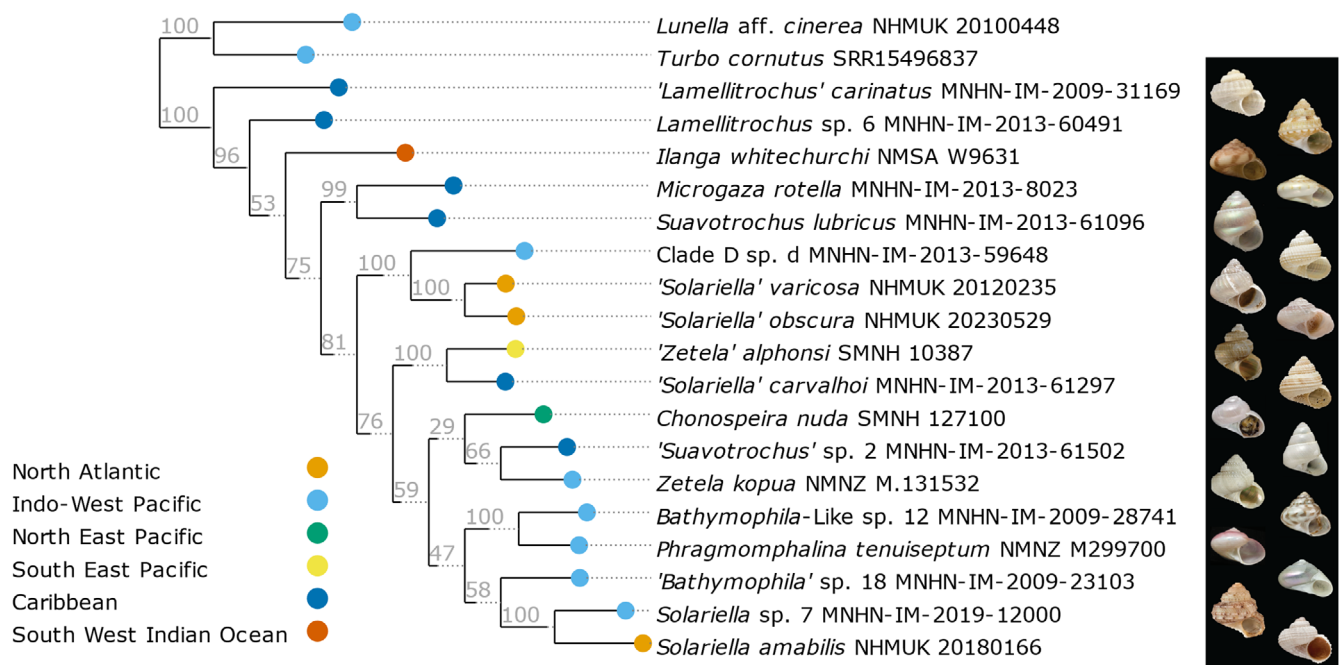


FIGURE 5 | Partitioned maximum likelihood tree of 15 genes including 12 mitochondrial protein-coding genes, two mitochondrial ribosomal genes and one nuclear ribosomal gene (28S) generated using IQ-TREE 2 and visualised using ete3. The tree is rooted on the outgroup taxa and values on branches are ultrafast bootstrap values. Images of each specimen including in the analysis are provided, except for outgroups *Lunella* aff. *cinerea* and *Turbo cornutus*.

degraded target DNA and the coextraction of contaminating human DNA accumulated from contact with samples during collection, curation, tissue sampling and laboratory work. Despite this, our pipelines avoided the assembly of human sequences by using reference databases from closely related species. The assembly of contaminant sequences from closely related species is more problematic, but these contaminants can be identified by sequence homology searches and phylogenetic analysis. From our analysis, we were able to determine that all sequences obtained for *Spectamen* cf. *bellulum* are likely contaminants from a morphologically distinct solariellid sequenced previously in the same lab. In addition, 28S annotations from *Phragmomphalina tenuiseptum* and *Zetela kopua* were also identified as likely contaminants from non-solariellid gastropods, also sequenced in the same lab. Therefore, these contaminants are likely due to laboratory contamination, which can be minimised by switching from a standard bench to clean room facilities and following protocols used for ancient DNA extractions.

The mitochondrial sequences assembled for outgroups *Turbo cornutus* and *Lunella* aff. *cinerea* showed high sequence homology to previously published sequences, although the assemblies were not as complete as those published on GenBank, with assemblies appearing to break in regions of low coverage, low GC content and/or high repeat content (Figure S3~13,676 bp; Figure S4~13,973 bp). Using a reference-based organelle assembly tool such as MITObim (Hahn, Bachmann, and Chevreux 2013) may increase the likelihood of a complete assembly, but at the risk of reference bias. However, for studies such as this, selecting a single reference may be problematic when working with a diverse range of taxa, with complex variation in organelle genome sequences. Our pipeline utilises

GetOrganelle which uses a *de novo* assembly approach and so should have a smaller reference bias, even from a diverse range of taxa. In addition, the most important output for the phylogenetic analysis is the list of annotated genes used in phylogenetic analyses, which were completed for both our outgroup samples.

Genome-skimming offers many advantages over traditional PCR amplification. Indeed, genome skimming and the *skim2mito* and *skim2rrna* pipelines were able to recover gene sequences that were previously unattainable by PCR. For example, *cox1* could be sequenced using Sanger sequencing for only 8 of 25 samples (not including outgroups), whereas *cox1* was recovered for 18 of 25 samples using genome skimming and *skim2mito*. Likewise, 28S could only be sequenced from PCR amplicons for 9 of 25 samples whereas genome skimming assembled 28S for 23 of 25 samples. The latter is likely due to DNA degradation since primer sequences come from a very conserved region, whereas some *cox1* failures may also be due to primer mismatches.

Previous phylogenetic analyses of solariellid gastropods have highlighted complex and unresolved phylogenetic relationships among genera (Williams et al. 2020). Although the tree in this study has good (>95%) ultrafast bootstrap support for most terminal clades, support for basal splits often have poorer support, suggesting that some assignments to genera require further research. One potential option would be to increase taxon sampling to test generic assignments. Likewise, while individual gene sequences from barcode genes are already available for a diverse range of solariellid gastropods, a greater sampling of gene space is required to confirm the relationships within this group.

Although this methodology will work for any sample type, *skim2mito* and *skim2rrna* were written specifically to account for many of the issues associated with historical museum samples, for example, DNA degradation and contamination. By default, *skim2mito* implements GetOrganelle using all reads as input (`--reduce-reads-for-coverage inf --max-reads inf`) and an increased number of rounds of (`-R 20`) of target read selection. For museum samples that are likely to be degraded, this maximises the inclusion of short sequencing reads. In addition, the user can specify a custom reference database for GetOrganelle using sequences from closely related taxa. This is necessary because our benchmarking of GetOrganelle using simulated datasets (White and Clark unpub. data) highlighted that a reference dataset containing closely related sequences increases the likelihood of successful assembly. Conversely, a broad reference dataset can increase the likelihood of sequence assembly from contaminated DNA. Taxonomic assignment of assembled sequences using blobtools can identify most non-target sequences, and phylogenetic analysis implemented for all annotated gene sequences may be particularly useful for identifying genetically similar contaminants not recognised by blobtools.

Although the pipelines presented simplify the bioinformatic analyses significantly, allowing for the analysis of hundreds of samples simultaneously, there is a risk to trade-off accuracy for increasing scale of the analysis. Indeed, our study highlighted that it was still important to manually check the assembled sequences for contamination or poorly annotated sequences using the standard outputs of our pipeline including blobtools, individual gene alignments and phylogenetic analyses. Although, it could be hardcoded to remove assembled sequences based on sequence homology to reference databases, this is not possible at present, because public databases are not complete for all taxa. In addition, taxonomic expertise may be necessary to identify incongruent phylogenetic relationships that can result from cross contamination from closely related taxa, highlighting the need for particular care when extracting DNA from historical specimens. An additional pipeline could be further adapted to include chloroplast organelle sequences by including a chloroplast annotation tool. However, there are few chloroplast annotation resources available as a command line tool with premade reference databases that annotate genes in a consistent and reliable way, as MITOS2 does for mitochondrial genes.

In conclusion, this study demonstrates that the snakemake pipelines *skim2mito*, *skim2rrna* and *gene2phylo* can cope with poor quality data from historical collections, facilitating large scale genome skimming studies from museum specimens. Given the current biodiversity crisis and lack of taxonomic expertise, it has become more important than ever to document biodiversity before it is lost. By sequencing natural history collections at scale using bioinformatic tools such as those presented here, researchers can increase the rate of phylogenetic and barcoding studies, and ultimately, species discovery.

Acknowledgements

We thank Andrea Waeschenbach for sequence data for '*Solariella*' *varicosa*; QuantaBio for providing the sparQ DNA Library Prep kit and sparQ PureMag Beads; Tore Hoisaeter for providing specimens of '*Solariella*' *obscura* and '*Solariella*' *varicosa*; Anders Warén for providing specimens of '*Solariella*' *amabilis* and loaning specimens of *Chonospeira nuda* and

'*Zetela*' *alphonsi*; Dai Herbert for loaning specimens of *Ilanga whitechurchi* and '*Spectamen*' *franciscanum*; Ian Loch for loaning specimens of *Archimnolia oleacea*; Bruce Marshall for loaning *Bathymophila gravida*, *Zetela textilis*, *Z. kopua* and the holotype of *Phragmomphalina tenuiseptum*; and Barbara Buge, Nico Puillandre and Philippe Bouchet for loaning the remaining specimens from MNHN. Thanks also to Jon Kongsrud for specimen locality data and to Harry Taylor for photos of NHMUK and MNHN specimens, Bruce Marshall for the photo of *Zetela kopua* and Sadie Mills for the photo of *Phragmomphalina tenuiseptum*. MNHN Specimens were obtained during research cruises and expeditions organised by the MNHN with Institut de Recherche pour le Développement and National Taiwan University, as part of the Tropical Deep-Sea Benthos program: ZhongSha 2015 in the South China Sea; MADEEP (<https://doi.org/10.17600/14004000>) in Papua New Guinea; BERYX 11 (<https://doi.org/10.17600/92005011>), EXBODI (<https://doi.org/10.17600/11100080>) and NORFOLK 1 (<https://doi.org/10.17600/1100050>) in New Caledonia; BORDAU 1 (<https://doi.org/10.17600/99100020>) in Fiji; KARUBENTHOS 2012 and KARUBENTHOS 2 (<https://doi.org/10.17600/15005400>) in Guadeloupe (more information can be found at expeditions.mnhn.fr). These expeditions operated under the regulations then in force in the countries in question and satisfy the conditions set by the Nagoya Protocol for access to genetic resources. Finally, we would like to acknowledge 2024 Hackathon: Enhancing Biodiversity Genomics Europe's Barcode Sequencing Stream (23-25 April, Leiden, Netherlands), which significantly improved the *skim2mito* pipeline. Specifically, we would like to acknowledge the input of Rutger Vos, Dick Groenenberg, Pierre Etienne Cholley (Naturalis Biodiversity Center, Netherlands), Douglas Lowe, Eli Chadwick (University of Manchester), Giorgia Staffoni (University of Padova), Hana Merchant and Daniel Parsons (Natural History Museum London). The hackathon was funded by Horizon Europe under the Biodiversity, Circular Economy and Environment (REA.B.3); co-funded by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number 22.00173; and by the UK Research and Innovation (UKRI) under the Department for Business, Energy and Industrial Strategy's Horizon Europe Guarantee Scheme.

Conflicts of Interest

The authors declare no conflicts of interest.

Data Availability Statement

Raw sequence data and assembled mitochondrial and ribosomal sequences are accessible from the European Nucleotide Archive (ENA) under project accession PRJEB76850, with the exception of previously published data for *Turbo cornutus* (SRR15496837; Kim et al. 2022). The snakemake pipelines are available on both GitHub and WorkflowHub: *skim2mito* <https://github.com/o-william-white/skim2mito>, *skim2rrna* <https://github.com/o-william-white/skim2rrna> and *gene2phylo* <https://github.com/o-william-white/gene2phylo>. Note that each pipeline is provided with the simulated data from 25 species of the butterfly superfamily Papilionoidea used in the manuscript, with instructions on how to run each pipeline. Configuration files and reference data used for the Solariellidae analyses are available on Data Dryad <https://doi.org/10.5061/dryad.h70rxwdt2>.

Open Research Badges



This article has earned an Open Data badge for making publicly available the digitally-shareable data necessary to reproduce the reported results. The data is available at <https://www.ncbi.nlm.nih.gov/bioproject/PRJEB76850/TheconfigurationfilesandreferencedataareavailableonDataDryad>, <https://doi.org/10.5061/dryad.h70rxwdt2>.

Benefit-Sharing Statement

Benefits Generated: The benefits from this research accrue from the sharing of our code, data and results on public databases as described above.

References

- Andrews, S. 2010. *FastQC*. Babraham Bioinformatics: Cambridge, England.
- Bartolozzi, L., L. Bettison-Varga, and N. Chernetsov. 2023. "A Global Approach for Natural History Museum Collections." *Science* 379, no. 6638: 1192–1194. <https://www.researchgate.net/publication/369480306>.
- Bernt, M., A. Donath, F. Jühling, et al. 2013. "MITOS: Improved de novo metazoan mitochondrial genome annotation." *Molecular Phylogenetics and Evolution* 69, no. 2: 313–319. <https://doi.org/10.1016/j.ympev.2012.08.023>.
- Bohmann, K., S. Mirarab, V. Bafna, and M. T. P. Gilbert. 2020. "Beyond DNA Barcoding: The Unrealized Potential of Genome Skim Data in Sample Identification." *Molecular Ecology* 29: 2521–2534. <https://doi.org/10.1111/mec.15507>.
- Brasseur, M. V., J. J. Astrin, M. F. Geiger, and C. Mayer. 2023. "MitoGeneExtractor: Efficient Extraction of Mitochondrial Genes From Next-Generation Sequencing Libraries." *Methods in Ecology and Evolution* 14, no. 4: 1017–1024. <https://doi.org/10.1111/2041-210X.14075>.
- Cai, L., H. Zhang, and C. C. Davis. 2022. "phyloHerb: A High-Throughput Phylogenomic Pipeline for Processing Genome Skimming Data." *Applications in Plant Sciences* 10, no. 3: e11475.
- Call, E., V. Twort, M. Espeland, and N. Wahlberg. 2023. "One Method to Sequence Them all? Comparison Between Whole-Genome Sequencing (WGS) and Target Enrichment (TE) of Museum Specimens From the Moth Families Epicopeidae." *BioRxiv*. <https://doi.org/10.1101/2023.08.21.553699>.
- Camacho, C., G. Coulouris, V. Avagyan, et al. 2009. "BLAST+: Architecture and Applications." *BMC Bioinformatics* 10: 1–9. <https://doi.org/10.1186/1471-2105-10-421>.
- Castresana, J. 2000. "Selection of Conserved Blocks From Multiple Alignments for Their Use in Phylogenetic Analysis." *Molecular Biology and Evolution* 17, no. 4: 540–552. <https://academic.oup.com/mbe/article/17/4/540/1127654>.
- Chen, S., Y. Zhou, Y. Chen, and J. Gu. 2018. "fastp: An Ultra-Fast All-In-One FASTQ Preprocessor." *Bioinformatics* 34, no. 17: i884–i890. <https://doi.org/10.1093/bioinformatics/bty560>.
- Chen, W., S. R. Achakkagari, and M. Strömvik. 2022. "Plastaumatic: Automating Plastome Assembly and Annotation." *Frontiers in Plant Science* 13: 1011948. <https://doi.org/10.3389/fpls.2022.1011948>.
- De Medeiros, B., L. Cai, P. J. Flynn, et al. 2024. "A Universal DNA Barcode for the Tree of Life." *EcoEvoRxiv*. <https://doi.org/10.32942/X24891>.
- D'Ercole, J., S. W. J. Prosser, and P. D. N. Hebert. 2021. "A SMRT Approach for Targeted Amplicon Sequencing of Museum Specimens (Lepidoptera)—Patterns of Nucleotide Misincorporation." *PeerJ* 9: e10420. <https://doi.org/10.7717/peerj.10420>.
- Freudenthal, J. A., S. Pfaff, N. Terhoeven, A. Korte, M. J. Ankenbrand, and F. Förster. 2020. "A Systematic Comparison of Chloroplast Genome Assembly Tools." *Genome Biology* 21, no. 1: 1–21. <https://doi.org/10.1186/s13059-020-02153-6>.
- Fulton, T. L., and B. Shapiro. 2019. "Setting up an Ancient DNA Laboratory." In *Ancient DNA: Methods and Protocols*, edited by B. Shapiro, A. Barlow, P. D. Heintzman, M. Hofreiter, J. L. A. Paijmans, and A. E. R. Soares. Methods in Molecular Biology, vol 1963. New York, NY: Humana Press. https://doi.org/10.1007/978-1-4939-9176-1_1.
- Hahn, C., L. Bachmann, and B. Chevreaux. 2013. "Reconstructing mitochondrial genomes directly from genomic next-generation sequencing reads—a baiting and iterative mapping approach." *Nucleic Acids Research* 41, no. 13: e129. <https://doi.org/10.1093/nar/gkt371>.
- Higuchi, R., B. Bowman, M. Freiburger, O. A. Ryder, and A. C. Wilson. 1984. "DNA Sequences From the Quagga, an Extinct Member of the Horse Family." *Nature* 312, no. 15: 282–284. <https://doi.org/10.1038/312282a0>.
- Jin, J. J., W. B. Yu, J. B. Yang, et al. 2020. "getOrganelle: A Fast and Versatile Toolkit for Accurate de Novo Assembly of Organelle Genomes." *Genome Biology* 21, no. 241: 241. <https://doi.org/10.1101/256479>.
- Katoh, K., and D. M. Standley. 2013. "MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability." *Molecular Biology and Evolution* 30, no. 4: 772–780. <https://doi.org/10.1093/molbev/mst010>.
- Kim, E., J. Kim, Y. Lee, G. Nah, and H. Y. Kim. 2022. "The Complete Mitochondrial Genome of *Turbo Cornutus* (Trochida: Turbinidae) and Its Phylogeny Analysis." *Mitochondrial DNA Part B: Resources* 7, no. 4: 637–639. <https://doi.org/10.1080/23802359.2022.2060764>.
- Krause, J., P. H. Dear, J. L. Pollack, et al. 2006. "Multiplex Amplification of the Mammoth Mitochondrial Genome and the Evolution of Elephantidae." *Nature* 439, no. 7077: 724–727. <https://doi.org/10.1038/nature04432>.
- Krzywinski, M., J. Schein, I. Birol, et al. 2009. "Circos: An Information Aesthetic for Comparative Genomics." *Genome Research* 19, no. 9: 1639–1645. <http://mkweb.bcgsc.ca/circos.Aninteractiveonlineversion>.
- Laetsch, D. R., M. L. Blaxter, and R. M. Leggett. 2017. "BlobTools: Interrogation of Genome Assemblies." *F1000Research* 6, no. 1287: 1–16.
- Li, H. 2013. "Aligning Sequence Reads, Clone Sequences And Assembly Contigs With BWA-MEM." *ArXiv*, 1303.3997v2. <http://arxiv.org/abs/1303.3997>.
- Li, H. 2018. "minimap2: Pairwise Alignment for Nucleotide Sequences." *Bioinformatics* 34, no. 18: 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
- Li, H., B. Handsaker, A. Wysoker, et al. 2009. "The Sequence Alignment/Map Format and SAMtools." *Bioinformatics* 25, no. 16: 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
- Meng, G., Y. Li, C. Yang, and S. Liu. 2019. "MitoZ: A Toolkit for Animal Mitochondrial Genome Assembly, Annotation and Visualization." *Nucleic Acids Research* 47, no. 11: e63. <https://doi.org/10.1093/nar/gkz173>.
- Minh, B. Q., H. A. Schmidt, O. Chernomor, et al. 2020. "IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era." *Molecular Biology and Evolution* 37, no. 5: 1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
- Mullin, V. E., W. Stephen, A. N. Arce, et al. 2023. "First Large-Scale Quantification Study of DNA Preservation in Insects From Natural History Collections Using Genome-Wide Sequencing." *Methods in Ecology and Evolution* 14, no. 2: 360–371. <https://doi.org/10.1111/2041-210X.13945>.
- Pouchon, C., F. Boyer, C. Roquet, et al. 2022. "ORTHOSKIM: In Silico Sequence Capture From Genomic and Transcriptomic Libraries for Phylogenomic and Barcoding Applications." *Molecular Ecology Resources* 22, no. 5: 2018–2037. <https://doi.org/10.1111/1755-0998.13584>.
- Prijbelski, A., D. Antipov, D. Meleshko, A. Lapidus, and A. Korobeynikov. 2020. "Using SPAdes De Novo Assembler." *Current Protocols in Bioinformatics* 70, no. 1: 102. <https://doi.org/10.1002/cpb1.102>.
- Quast, C., E. Pruesse, P. Yilmaz, et al. 2013. "The SILVA Ribosomal RNA Gene Database Project: Improved Data Processing and Web-Based Tools." *Nucleic Acids Research* 41, no. D1: D590–D596. <https://doi.org/10.1093/nar/gks1219>.
- R Core Team. 2020. "R: A language and Environment for Statistical Computing." In *R Foundation for Statistical Computing* (3.6.3). Vienna, Austria: R Foundation for Statistical Computing.

- Renaud, G., K. Hanghøj, E. Willerslev, and L. Orlando. 2017. "Gargammel: A Sequence Simulator for Ancient DNA." *Bioinformatics* 33, no. 4: 577–579. <https://doi.org/10.1093/bioinformatics/btw670>.
- Ruane, S., and C. C. Austin. 2017. "Phylogenomics Using Formalin-Fixed and 100+ Year-Old Intractable Natural History Specimens." *Molecular Ecology Resources* 17, no. 5: 1003–1008. <https://doi.org/10.1111/1755-0998.12655>.
- Sarmashghi, S., M. Balaban, E. Rachtman, B. Touri, S. Mirarab, and V. Bafna. 2021. "Estimating Repeat Spectra and Genome Length From Low-Coverage Genome Skims With RESPECT." *PLoS Computational Biology* 17, no. 11: e1009449. <https://doi.org/10.1371/journal.pcbi.1009449>.
- Sarmashghi, S., K. Bohmann, P. Thomas, M. Gilbert, V. Bafna, and S. Mirarab. 2017. "Skmer: Assembly-Free and Alignment-Free Sample Identification Using Genome Skims." *Genome Biology* 23, no. 34: 1–20. <https://doi.org/10.1101/230409>.
- Shen, W., S. Le, Y. Li, and F. Hu. 2016. "SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File Manipulation." *PLoS One* 11, no. 10: e0163962. <https://doi.org/10.1371/journal.pone.0163962>.
- Steenwyk, J. L., T. J. Buida, Y. Li, X. X. Shen, and A. Rokas. 2020. "ClipKIT: A Multiple Sequence Alignment Trimming Software for Accurate Phylogenomic Inference." *PLoS Biology* 18, no. 12: e3001007. <https://doi.org/10.1371/journal.pbio.3001007>.
- Straub, S. C. K., M. Parks, K. Weitemier, M. Fishbein, R. C. Cronn, and A. Liston. 2012. "Navigating the Tip of the Genomic Iceberg: Next-Generation Sequencing for Plant Systematics." *American Journal of Botany* 99, no. 2: 349–364. <https://doi.org/10.3732/ajb.1100335>.
- Straube, N., M. L. Lyra, J. L. A. Paijmans, et al. 2021. "Successful Application of Ancient DNA Extraction and Library Construction Protocols to Museum Wet Collection Specimens." *Molecular Ecology Resources* 21, no. 7: 2299–2315. <https://doi.org/10.1111/1755-0998.13433>.
- Sumner-Rooney, L., J. D. Sigwart, J. McAfee, L. Smith, and S. T. Williams. 2016. "Repeated Eye Reduction Events Reveal Multiple Pathways to Degeneration in a Family of Marine Snails." *Evolution* 70, no. 10: 2268–2295. <https://doi.org/10.1111/evo.13022>.
- Wick, R. R., M. B. Schultz, J. Zobel, and K. E. Holt. 2015. "Bandage: Interactive Visualization of de Novo Genome Assemblies." *Bioinformatics* 31, no. 20: 3350–3352. <https://doi.org/10.1093/bioinformatics/btv383>.
- Williams, S. T., P. G. Foster, and D. T. J. Littlewood. 2014. "The Complete Mitochondrial Genome of a Turbinid Vetigastropod From MiSeq Illumina Sequencing of Genomic DNA and Steps Towards a Resolved Gastropod Phylogeny." *Gene* 533, no. 1: 38–47. <https://doi.org/10.1016/j.gene.2013.10.005>.
- Williams, S. T., Y. Kano, A. Warén, and D. G. Herbert. 2020. "Marrying Molecules and Morphology: First Steps Towards a Reevaluation of Solariellid Genera (Gastropoda: Trochoidea) in the Light of Molecular Phylogenetic Studies." *Journal of Molluscan Studies* 86, no. 1: 1–26. <https://doi.org/10.1093/mollus/eyz038>.
- Williams, S. T., E. S. Noone, L. M. Smith, and L. Sumner-Rooney. 2022. "Evolutionary Loss of Shell Pigmentation, Pattern, and Eye Structure in Deep-Sea Snails in the Dysphotic Zone." *Evolution* 76, no. 12: 3026–3040. <https://doi.org/10.1111/evo.14647>.
- Williams, S. T., L. M. Smith, D. G. Herbert, et al. 2013. "Cenozoic Climate Change and Diversification on the Continental Shelf and Slope: Evolution of Gastropod Diversity in the Family Solariellidae (Trochoidea)." *Ecology and Evolution* 3, no. 4: 887–917. <https://doi.org/10.1002/ece3.513>.
- Wu, P., C. Xu, H. Chen, J. Yang, X. Zhang, and S. Zhou. 2021. "NOVOWrap: An Automated Solution for Plastid Genome Assembly and Structure Standardization." *Molecular Ecology Resources* 21, no. 6: 2177–2186. <https://doi.org/10.1111/1755-0998.13410>.
- Yu, G., D. K. Smith, H. Zhu, Y. Guan, and T. T. Y. Lam. 2017. "Ggtree: An R Package for Visualization and Annotation of Phylogenetic Trees With Their Covariates and Other Associated Data." *Methods in Ecology and Evolution* 8, no. 1: 28–36. <https://doi.org/10.1111/2041-210X.12628>.
- Zhang, C., M. Rabiee, E. Sayyari, and S. Mirarab. 2018. "ASTRAL-III: Polynomial Time Species Tree Reconstruction From Partially Resolved Gene Trees." *BMC Bioinformatics* 19: 153. <https://doi.org/10.1186/s12859-018-2129-y>.

Supporting Information

Additional supporting information can be found online in the Supporting Information section.