

1

2

3 **METAWORKS: A flexible, scalable bioinformatic pipeline for multi-marker**

4 **biodiversity assessments**

5

6 Teresita M. Porter^{1*}, Mehrdad Hajibabaei¹

7

8 ¹ Centre for Biodiversity Genomics @ Biodiversity Institute of Ontario & Department of

9 Integrative Biology, University of Guelph, Guelph, ON, CANADA

10

11 * Corresponding author:

12 T.M. Porter

13 terrimporter@gmail.com

14

15

16 **Abstract**

17

18 **Background:** Multi-marker metabarcoding is increasingly being used to generate
19 biodiversity information across different domains of life from microbes to fungi to
20 animals such as in ecological and environmental studies. Current popular bioinformatic
21 pipelines support microbial and fungal marker analysis, while ad hoc methods are used
22 to process animal metabarcode markers from the same study. The purpose of this
23 paper is to introduce MetaWorks, a ‘meta’barcode pipeline that does ‘the works’ and
24 supports the bioinformatic processing of various metabarcoding markers including rRNA
25 and their spacers as well as protein coding loci.

26 **Results:** MetaWorks provides a Conda environment to quickly gather most of the
27 programs and dependencies for the pipeline. MetaWorks is automated using
28 Snakemake to ensure reproducibility and scalability. We have supplemented existing
29 RDP-trained classifiers for SSU (prokaryotes), ITS (fungi), and LSU (fungi) with trained
30 classifiers for COI (eukaryotes), rbcL (diatoms or eukaryotes), SSU (diatoms or
31 eukaryotes), and 12S (fish). MetaWorks can process rRNA genes, but it can also
32 properly handle ITS spacers by trimming flanking conserved rRNA gene regions, as well
33 as handle protein coding genes by removing obvious pseudogenes.

34 **Conclusions:** As far as we are aware, MetaWorks is the first flexible multi-marker
35 metabarcode pipeline that can accommodate rRNA genes, spacer, and protein coding
36 markers in the same pipeline. This is ideal for large-scale, multi-marker studies to
37 provide a harmonized processing environment, pipeline, and taxonomic assignment
38 approach. Updates to MetaWorks will be made as needed to reflect advances in the

39 underlying programs, reference databases, or hidden Markov model (HMM) profiles for
40 pseudogene filtering. Future developments will include support for additional
41 metabarcode markers, RDP trained reference databases, and HMM profiles for
42 pseudogene filtering.

43

44 **Keywords**

45 Metabarcoding, conda, snakemake, COI, rbcL, rRNA gene, ITS

46

47 **Background**

48

49 Marker gene sequencing, metabarcoding, or metasytematics are interrelated
50 techniques that involves extracting environmental DNA from bulk samples such as soil,
51 water, or individuals collected from traps without having to isolate any individual
52 specimens followed by enrichment of a signature DNA region to identify biological
53 community composition using bioinformatics [1–3]. In different fields of study, from
54 microbial ecology to animal biodiversity studies, different signature DNA regions are
55 targeted for their ability to identify target taxa. For example, in prokaryotes, the 16S
56 rRNA region is often used for genus level taxonomic assignments [4, 5]. In animals,
57 plants, and fungi, COI, rbcL, or ITS are commonly targeted to identify metabarcode
58 reads to the species rank, respectively, based on the availability of reference sequences
59 [6–9]. Other markers commonly used for phylogenetic analyses are also popular
60 targets, such as SSU for eukaryotes, arbuscular mycorrhizal fungi, and diatoms; or 12S
61 for fish [10–15]. For each of these markers, reference sequences are often housed in

62 their own separate databases where they can be analyzed by custom- or built-in tools
63 [16–22].

64 Existing well-developed and popular pipelines such as QIIME and MOTHUR
65 were initially developed to support the microbial ecology community [23–25]. Methods
66 incorporated into these pipelines largely support the analysis of 16S rRNA genes, but
67 they do also support the analysis of popular fungal markers such as ITS and LSU [5, 19,
68 26, 27]. QIIME also supports a data flow to properly isolate the ITS spacer regions from
69 conserved flanking regions [28]. In our study of biodiversity genomics, we also had a
70 need for a pipeline that could handle other metabarcode regions such as COI, rbcL, and
71 12S. For protein-coding markers, we wanted the option to filter for pseudogenes as an
72 additional way to remove these non-target sequences from the dataset. Pseudogenes
73 are duplicated copies of a gene that are not functional, may be truncated, under relaxed
74 selection pressure means these sequences may accumulate insertions and deletions
75 that may introduce premature stop codons or frame shifts. This may be problematic for
76 biodiversity studies that use metabarcoding if these pseudogenes are amplified and
77 result in inflated richness estimates. Also, since sequencing technologies are often
78 changing, with datasets getting larger, and active development of metabarcode
79 sequence handling programs are often updated, we also wanted a pipeline with the
80 flexibility to update underlying programs frequently at the start of every new project.

81 For analyzing rRNA genes such as 16S, 18S, and ITS, there exist numerous
82 projects that offer the ability to identify metabarcodes based on a comparison to
83 sequences classified using morphology-based taxonomy, phylogenetic relatedness, or
84 some mixture of both and these are supported by the major pipelines [16, 17, 19, 20].

85 For analyzing animal or plant markers, such as COI and rbcL, the BOLD database has
86 a number of in-house tools that are appropriate for analyzing single sequences and is
87 best used as a curation tool [18] but the underlying database is not fully available in a
88 format that existing popular pipelines can use as is. We also identified a need for a
89 pipeline that would use a proven taxonomic assignment method, as opposed to a
90 sequence similarity-based method like BLAST, to identify our metabarcodes and reduce
91 false positive rates [29–31]. Thus, we also needed RDP-trained reference sequence
92 datasets to support the assignment of 18S, 12S, COI, and rbcL datasets, with options to
93 support the identification of target taxa such as diatoms or fish [29, 30, 32]. Another
94 important consideration when working with protein coding genes, is the ability to easily
95 filter out obvious pseudogenes, such as the nuclear encoded mitochondrial sequences
96 (NuMTs) co-sequenced with COI primers during PCR [33, 34].

97 For large scale or multi-year projects such as ecological and biomonitoring
98 programmes, e.g. Baird and Hajibabaei, 2012, the question of version control and
99 reproducibility is of utmost concern [35]. Coordinating analyses across different labs
100 focusing on different target taxa or utilizing different markers can be complicated, with
101 researchers each using their own bioinformatic pipelines as appropriate. Here is where
102 we also felt the need to use tools to ensure environment as well as pipeline
103 reproducibility. Though the target taxa or markers may differ across labs, we saw a
104 need for a unifying bioinformatic pipeline that could offer a degree of standardization,
105 consistency, and reproducibility.

106 As multi-marker studies are carried out on phylogenetically divergent taxa, such
107 as in biodiversity or trophic studies, there is a need for more generic pipelines where

108 different markers can be analyzed using similar dataflows with 3rd party programs
109 instead of being limited to database-specific pipelines and tools [36, 37]. Our objective
110 was to develop a flexible bioinformatic pipeline suitable for processing multi-marker
111 metabarcode datasets generated using paired-end Illumina sequencing with the
112 following considerations: 1) reproducibility with respect to the computational
113 environment used as well as the pipeline itself, 2) scalability to leverage multi-core
114 processors to speed up the analysis of large datasets, 3) support the use widely used
115 rRNA gene, spacer, or protein-coding markers. To support the last consideration, we
116 have curated multiple reference sets to ensure consistently defined taxonomic lineages,
117 trained the RDP naive Bayesian classifier to make rapid and accurate flexible-rank
118 taxonomic assignments. Our pipeline also supports the processing of protein-coding
119 markers, by implemented steps to help remove obvious pseudogenes. Finally, as
120 collaborators in several large-scale metabarcoding projects involving various
121 stakeholders, this pipeline has evolved to address the practical needs of various
122 projects from graduate student research to national biomonitoring programmes.

123

124 **Implementation**

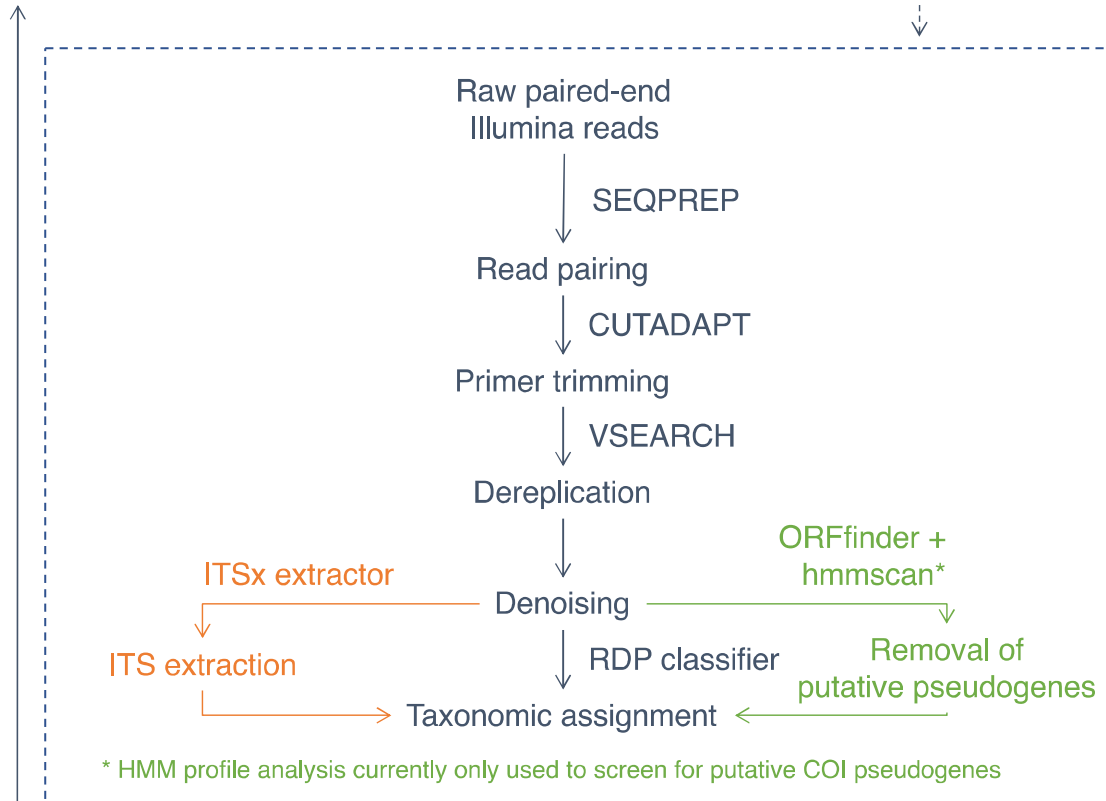
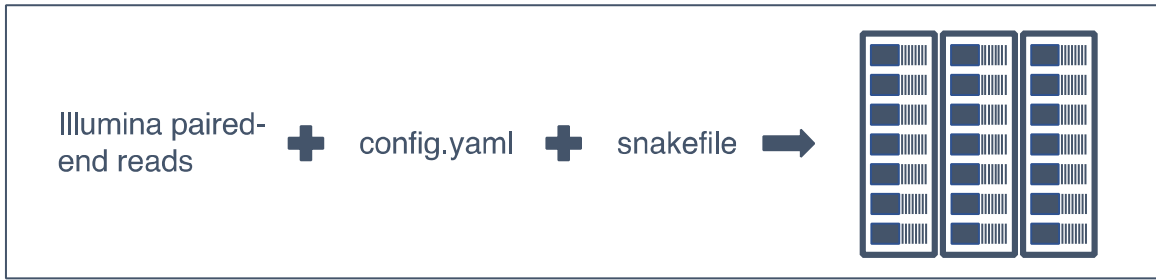
125 MetaWorks is a multi-marker ‘meta’-barcode pipeline that does ‘the works’ by
126 supporting the bioinformatic processing of popular markers including rRNA genes,
127 spacers, and protein coding genes. The standard pipeline is shown in Figure 1, as well
128 as accommodations for processing ITS spacers, such as the removal of flanking rRNA
129 gene sequences using the ITSx program [38], and protein coding genes, such as
130 pseudogene filtering. More detailed data flows used to process different markers are

131 shown in the supplement (Fig S1 - S4). This pipeline is meant to be run in a Linux
132 environment at the command-line. The most recent version is available from
133 <https://github.com/terrimporter/MetaWorks>.

134

135 **Figure 1. MetaWorks v1 pipeline overview.** The pipeline can be run in a conda
136 environment, providing raw paired-end Illumina reads, a configuration file, and a
137 snakefile. Snakemake can be directed to run parallel jobs across many CPUs in a high
138 performance computing environment. A summary of the steps carried out in the
139 pipeline are shown in the dashed box. The standard pipeline is shown along with the
140 variation needed to process ITS sequences (orange) and the variations needed to
141 screen out putative pseudogenes (green). We note with the asterisk that the
142 pseudogene removal step is currently different for rbcL and COI: for rbcL, longest ORF
143 lengths are screened for outliers; whereas for COI, longest ORFs are further subjected
144 to hidden Markov model (HMM) profile analysis and HMM scores are used to screen for
145 outliers. For each exact sequence variant (ESV), for each sample, read counts and
146 taxonomic assignments are provided along with bootstrap support values. An example
147 of the taxonomic assignment output is shown in the table.

conda environment



Rank	Taxon	Bootstrap support
Kingdom	Metazoa	1.0
Phylum	Arthropoda	1.0
Class	Insecta	1.0
Order	Lepidoptera	1.0
Family	Coleophoridae	1.0
Genus	Coleophora	1.0
Species	<i>Coleophora duplicis</i>	0.56*

* For a ~ 200 bp fragment, we need a bootstrap support value ≥ 0.30 for 95% confidence in the assignment at this rank

149

150

151 The first feature of MetaWorks is the use of a conda environment [39]. We
152 provide a conda environment file (environment.yml) and when activated ensures that
153 most of the required software programs and their dependencies are available for the
154 pipeline to call and ensures a consistent processing environment is created for all users
155 who run the pipeline. As the pipeline is updated, the environment may also be updated
156 to contain the newest versions of the underlying programs. Using the provided
157 environment also ensures that the correct version of each program is used. If anaconda
158 isn't already installed on the user's system, a stripped-down version of anaconda,
159 'miniconda' is available from <https://docs.conda.io/en/latest/miniconda.html> . We
160 provided instructions on how to install and use conda on the MetaWorks README page
161 on GitHub <https://github.com/terrimporter/MetaWorks>.

162 Unfortunately, not all of the programs we use in our pipelines are currently
163 available as conda packages and if they are not already available on the user's system,
164 they will need to be installed separately. For example, the Ribosomal Database Project
165 (RDP) classifier v2.12 is available from <https://sourceforge.net/projects/rdp-classifier/>
166 and is used to make taxonomic assignments [5]. The RDP classifier uses a naive
167 Bayesian method to taxonomically assign unknown query sequences as well as provide
168 a measure of statistical support for each assignment at each rank. We have previously
169 described and compared how this method works compared to the top BLAST hit
170 method [30]. In that paper, we showed how the classifier is faster than the top BLAST
171 hit method and helps to reduce the rate of false-positive assignments. In studies where

172 erroneously identifying a metabarcode sequence as a potential invasive species or
173 pathogen could lead to alarm, reducing the false-positive assignment rate is critical. An
174 additional program is needed if pseudogene filtering of protein-coding genes is carried
175 out. The NCBI ORFfinder is a translation program used to identify open reading frames
176 (ORFs) and is used here help screen out putative pseudogenes and is available from
177 <https://www.ncbi.nlm.nih.gov/orffinder/>. Instructions on how these programs can be
178 downloaded and installed are also provided in the MetaWorks README file on GitHub.

179 The pipeline itself is automated using Snakemake [40]. Snakemake requires
180 three sets of files to run: 1) raw paired-end Illumina sequence files, 2) a configuration
181 file that specifies file paths, as well as program and pipeline settings, and, 3) a snakefile
182 that runs each step of the bioinformatic pipeline according to the settings in the
183 configuration file. One advantage of using Snakemake over a simple shell script to
184 automate the pipeline, is the ability to resume a pipeline from where it left off if a
185 problem is encountered. Another feature is that Snakemake supports the use of
186 multiple parallel jobs and this is ideal for high performance computing environments
187 where many cores are available to speed up the analysis of large datasets. The
188 MetaWorks pipeline is also versioned and regularly updated on GitHub.

189 Snakemake also requires a configuration file. In this file, the user can specify the
190 primers used in the study and adjust the settings for the major steps of the pipeline
191 (described below). The user needs to enter the complete paths to where the raw data is
192 stored and provide the path to the RDP classifier [5]. The RDP classifier already
193 provides trained prokaryote SSU (16S), fungal LSU (28S), and fungal ITS
194 (Warcup/UNITE) datasets [5, 19, 26, 27]. To use a custom trained dataset, such as

195 SSU (diatoms or eukaryotes), rbcL (diatoms or eukaryotes), or 12S (fish) the trained
196 classifier first needs to be downloaded from GitHub and the path to this reference
197 dataset needs to be provided in the configuration file. For most users, only the
198 configuration file will need to be edited and customized and the snakefile itself does not
199 need to be changed.

200 The standard MetaWorks pipeline is ideal for processing rRNA gene markers.
201 Raw, paired-end Illumina reads are paired using SEQPREP [41]. Default settings are
202 used with the option to change the Phred quality score cutoff (default, 20) or the
203 minimum overlap (default, 25 bp). Primers are removed by aligning the primer
204 sequence to the sequence and removing this sequence region using CUTADAPT in two
205 separate steps, the first to remove the forward primer region, then successfully trimmed
206 sequences are trimmed to remove the reverse primer region [42]. Default settings are
207 used with the ability to customize the primer sequences (use N's instead of I's if
208 applicable), specify the minimum sequence retained after trimming (default, 150 bp),
209 adjust the Phred quality threshold at the 5' and 3' ends of the sequence (default, 20,
210 20), and to set the maximum number of N's in the primer sequence region (default, 3).

211 At this stage, the sequences for each sample are still in their own individual files.
212 To run a global analysis, each of these sequence files concatenated into a single new
213 file. These primer-trimmed reads are now ready for further processing in VSEARCH
214 [43]. VSEARCH is an open-source program that can be used as an alternative to
215 USEARCH [44]. USEARCH is proprietary software appropriate for processing marker
216 gene sequences / metabarcoding reads with a free 32-bit version limited to using 4Gb
217 memory or less. The 64-bit version can use all the available memory on a user's

218 system and requires a paid license. For small datasets, USEARCH can be used as is
219 with no issues, but for larger datasets that require more memory, VSEARCH may be a
220 better alternative. As a result, we use VSEARCH in this pipeline but keep in mind that
221 the algorithms used for dereplication, denoising, chimera-removal, and tracking read
222 counts are based on those originally developed in USEARCH [45, 46].

223 Primer-trimmed reads are dereplicated using the 'derep_fulllength' command in
224 VSEARCH while tracking read numbers in each sequence cluster for downstream
225 steps. This step retains just the unique sequences. At this step it is possible to have
226 two sequences that differ in length but may be identical in the overlapping region.
227 Sequences in the output file are ordered by decreasing cluster size. Dereplicated reads
228 are then denoised using the 'unoise3' algorithm. At this step, sequences with predicted
229 errors are corrected and rare reads are removed. The parameter to set the minimum
230 number of reads per cluster is set to 3 to remove just singletons and doubletons, but
231 this can be modified in the configuration file. The output can be thought of as
232 operational taxonomic units (OTUs) clustered with 100% sequence similarity, but they
233 have also been denoised and we refer to these as exact sequence variants (ESVs).
234 Similar sets of denoised reads have been previously described as amplicon sequence
235 variants (ASVs) when using the DADA2 pipeline or zero-radius OTUs (ZOTUs) in the
236 original USEARCH pipeline [45, 47]. In USEARCH, chimera removal is incorporated in
237 the unoise3 function, but in VSEARCH, chimera-removal needs to be run separately
238 using the 'uchime3_denovo' command. To map read counts to the newly generated
239 denoised ESVs, we use the 'search_exact' method. The 'search_exact' method is

240 preferred over the 'usearch_global' with the 'id 1.0' parameter because it is faster and
241 optimized to find exact matches.

242 When processing ITS sequences, we use the ITSx extractor to remove any rRNA
243 gene regions (SSU, 5.8S, or LSU) [38]. ITS is the standard fungal barcode marker and
244 one of the supplementary barcodes for plants [6, 48, 49, 8]. The ITS1 region is the
245 targeted region in the Earth Microbiome Project and the ITS2 region is targeted by
246 others including the Metagenomics-based ecosystem biomonitoring project
247 (Ecobiomics) and the ITS2 database V [4, 21, 50]. Our current pipeline is set to retain
248 the ITS2 region but can be edited to target the ITS1 region in the configuration file.

249 Each denoised ESV is taxonomically assigned using the Ribosomal Database
250 Project (RDP) classifier [5]. The RDP classifier uses a naive Bayesian method for
251 taxonomic assignment and requires both reference sequences as well as a clearly
252 defined taxonomy file as input for training. Although the classifier was originally
253 developed to classify prokaryote 16S rRNA gene sequences using the Bergey's
254 taxonomy, we have retrained the classifier for a number of other rRNA gene and protein
255 coding markers including COI (eukaryotes), rbcL (diatom and eukaryote), SSU (diatom
256 and eukaryote), and 12S (fish) (Table 1) [30, 32]. Each of the trained reference sets
257 should be downloaded, as needed, for MetaWorks processing. Note that the FASTA
258 formatted files used for training the RDP classifier are also available for each of these
259 reference sets and that they could alternatively be used for custom BLAST database
260 creation if needed for comparison purposes against the top BLAST hit method.

261

262 **Table 1. RDP-trained reference sets that can be used with MetaWorks v1.**

Marker	Target taxa	Classifier availability	Number of	Number of	Source data
--------	-------------	-------------------------	-----------	-----------	-------------

			included sequence s	included taxa at all ranks (species)	
COI	Eukaryotes	https://github.com/terrimp/orter/CO1Classifier	1,221,528	154,351 (114,687)	BOLD [18], INSDC [51]
rbcl	Diatoms	https://github.com/terrimp/orter/rbclDiatomClassifier	3,504	1,432 (1,023)	R-Syst::diatom [22]
rbcl	Eukaryotes	https://github.com/terrimp/orter/rbclClassifier	164,454	65,742 (53,344)	INSDC
12S	Fish	https://github.com/terrimp/orter/12SfishClassifier	2,853	4,751 (2,833)	MitoFish [52]
SSU (18S)	Diatoms	https://github.com/terrimp/orter/SSUdiatomClassifier	2,962	1,198 (828)	R-Syst::diatom
SSU (18S)	Eukaryotes	https://github.com/terrimp/orter/18SClassifier	42,301	7,504 (5,440 genera)	SILVA [17]
SSU (16S)	Prokaryotes	Built-in to the RDP classifier*	13,212	3,247 (2,506 genera)	RDP [5]
ITS	Fungi (Warcup)	Built-in to the RDP classifier	17,878	10,621 (8,551)	[27]
ITS	Fungi (UNITE)	Built-in to the RDP classifier	145,019	23,222 (20,337)	[19]
LSU	Fungi	Built-in to the RDP classifier	11,442	2,633 (1,895)	[26]

263

264

265

266 For protein coding genes, we attempt to remove obvious pseudogenes. For
 267 rbcl, we translate our denoised ESVs using ORFfinder [53]. We retain the longest
 268 open reading frame (ORF) and screen for ORFs with outlier lengths. Outliers are
 269 identified based on their nucleotide ORF sequence length. ORFs with a sequence
 270 length less than the 25th percentile - (1.5 x interquartile range (IQR)) are removed as
 271 short outliers. ORFs with a sequence length greater than the 75th percentile length +
 272 (1.5 x IQR) are removed as long outliers. An extra step was used to filter out potential
 273 COI pseudogenes. Denoised ESVs were translated as described above and the
 274 longest ORFs were retained. Then we used hidden Markov model (HMM) profile
 275 analysis using hmmscan to compare the amino acid ORFs to an HMM profile using the
 276 program HMMER available from <http://hmmer.org/> . The COI HMM profile was built

277 using COI sequences mined from the BOLD data releases as previously described in
278 [Porter and Hajibabaei, in prep]. Amino acid ORFs with short outlier HMMER scores
279 were filtered out of the dataset as putative pseudogenes (or genuine sequences with
280 PCR/sequencing errors).

281 The final results file is a comma separated file suitable for import into R for data
282 analysis. The output contains ESVs for each sample, read counts, ESV/ORF
283 sequences along with the taxonomic assignment. Bootstrap support values are also
284 provided for the taxonomic assignments at each rank. These bootstrap support values
285 can be used to filter for assignments where assignments are likely to be correct 95-99%
286 confidence, assuming the query is present in the reference sequence database. A
287 guide for cutoff values is provided for each custom-trained classifier in the GitHub
288 README pages (Table 1). For the built-in reference sets in the RDP classifier, an 80%
289 cutoff is recommended [5] unless the query sequence is shorter than 250 bp in which
290 case a cutoff of 50% is recommended on their website at
291 <https://rdp.cme.msu.edu/classifier/classifier.jsp>. In addition, we also provide statistics
292 for the output at each major bioinformatic step (number of sequences, as well as max,
293 min, mean, median, and mode lengths). Log files are also retained for each major step
294 run in VSEARCH. If more than one amplicon was pooled prior to sample indexing, the
295 pipeline can be run multiple times updating the configuration file as needed.

296

297 **Discussion**

298 MetaWorks is meant to be run at the command line in a conda environment [39].
299 The use of a conda environment facilitates quickly obtaining most of the programs and

300 dependencies needed to run the pipeline. The choice of Snakemake to automate the
301 pipeline was meant to ensure reproducibility and provide scalability when run in a high-
302 performance computing environment [40]. The pipeline was built with flexibility in mind,
303 to support the bioinformatic processing of multiple microbial, fungal, or animal
304 metabarcode markers. There is no need to gather separate OTU x sample tables or
305 FASTA files for downstream analyses. The final results file contains taxonomic
306 assignments with bootstrap support values, for each ESV, from each sample, along with
307 read counts. The output is in a comma separated file, easily imported into R, a popular
308 open-source software environment for data analysis and visualizations [54]. At this
309 step, ESV x sample tables can be regenerated or ESV/ORF sequences can be
310 extracted for further analysis.

311 MetaWorks is a flexible, scalable bioinformatic pipeline that will process raw
312 paired-end Illumina reads for any study that uses a metabarcoding approach. We
313 envision that MetaWorks will fill a need in multi-marker metabarcoding studies that
314 target taxa from multiple different domains of life, to provide a unified processing
315 environment, pipeline, and taxonomic assignment approach for each marker from
316 ribosomal RNA genes, spacers, or protein coding genes. QIIME 2 is perhaps the most
317 popular and comprehensive platform for such work, but to date, focuses on processing
318 mainly prokaryote and fungal datasets [25]. As of yet, MetaWorks is the only
319 bioinformatic pipeline that can handle rRNA genes but that also integrates special
320 processing steps to handle ITS spacers as well as filter out obvious pseudogenes in
321 protein coding markers such as COI. For example, we have a pseudogene filtering
322 protocol for filtering out COI pseudogenes based on HMM profile analysis. This is made

323 possible due to the large amount of high-quality COI reference data available from the
324 BOLD data releases [18]. Where a comprehensive hmm profile is lacking, we provide a
325 more general pseudogene filtering pipeline based on simple translation, retention of
326 longest ORFs, and removal of ORFs with outlier lengths.

327 There has been a lot of activity with respect to building new bioinformatic tools to
328 handle COI metabarcodes. Recent work, such as the BOLDigger program, has
329 attempted to make the BOLD identification engine more suitable for identifying large
330 batches of COI metabarcodes, but submission size is limited to 100 queries at a time
331 so as not to exceed the limits set by the BOLD server [55]. Unfortunately, point-and-
332 click type interfaces are not always easy to integrate into custom pipelines, especially
333 when dealing with large batches of sequences from high throughput sequencing
334 platforms, so a command-line interface would be welcome. Another python package
335 called 'Alfie' calculates k-mer frequencies and classifies COI metabarcode sequences to
336 the kingdom rank using a machine learning method [56]. A new program, called
337 NUMTdumper, has been developed as a stand-alone program meant to be incorporated
338 into bioinformatic pipelines [57]. NUMTdumper provides a method to screen for NuMTs
339 based on read counts while acknowledging the trade-offs between removing all possible
340 NuMTs while erroneously removing genuine reads. An R package called 'coil' has also
341 recently been developed that will place COI barcode and metabarcode sequences in
342 frame using profile HMM analysis [58]. In MetaWorks, we incorporate two traditional
343 approaches to identifying and removing obvious pseudogenes, one based on nucleotide
344 to open reading frame translation, and a second based on translation combined with

345 HMM profile analysis where an HMM profile is available, and we incorporate these
346 steps into the bioinformatic pipeline.

347 Our work on database curation, taxonomic assignment methods, and pipelines
348 continues. MetaWorks is currently used in the national biomonitoring project STREAM
349 (<https://stream-dna.com/>) and in various other research projects where we advance the
350 pipeline and its components regularly as needed. Updates to the underlying RDP-
351 trained classifiers described here are added to GitHub in an ongoing basis. New
352 versions of MetaWorks will be released when the underlying software packages make
353 major changes or when an existing reference database is updated or a new reference
354 database or hmm profile is created for pseudogene filtering.

355

356 **Conclusions**

357 MetaWorks is a fusion of several different marker-specific pipelines that we have
358 developed over the years. As far as we are aware, MetaWorks is the first multi-marker
359 metabarcode pipeline that not only handles rRNA genes but implements the steps to
360 rigorously analyze ITS spacers and protein coding genes. In large-scale studies that
361 use multiple metabarcode markers to target taxa across multiple domains of life,
362 MetaWorks can be used to provide a harmonized computational environment, pipeline,
363 and rigorous taxonomic assignment method. This flexible pipeline is open-source and
364 can be modified to support additional metabarcode markers, RDP trained reference
365 databases, or hmm profiles for pseudogene filtering.

366

367 **Availability and implementation**

368

369 The MetaWorks pipeline is available on GitHub at

370 <https://github.com/terrimporter/MetaWorks> as are the RDP classifier-formatted

371 reference sets and underlying FASTA files presented in Table 1 at

372 <https://github.com/terrimporter> .

373

374 **Acknowledgements**

375

376

377

378 **References**

- 379 1. Pace NR. A Molecular View of Microbial Diversity and the Biosphere. *Science*. 1997;276:734–
380 40.
- 381 2. Hajibabaei M. The golden age of DNA metasytematics. *Trends in genetics*. 2012;28:535–537.
- 382 3. Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E. Towards next-generation
383 biodiversity assessment using DNA metabarcoding. *Molecular ecology*. 2012;21:2045–2050.
- 384 4. Gilbert JA, Jansson JK, Knight R. The Earth Microbiome project: successes and aspirations.
385 *BMC biology*. 2014;12:69.
- 386 5. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid Assignment of
387 rRNA Sequences into the New Bacterial Taxonomy. *Applied and Environmental Microbiology*.
388 2007;73:5261–7.
- 389 6. CBOL Plant Working Group, Hollingsworth PM, Forrest LL, Spouge JL, Hajibabaei M,
390 Ratnasingham S, et al. A DNA barcode for land plants. *Proceedings of the National Academy of
391 Sciences*. 2009;106:12794–12797.
- 392 7. Hebert PDN, Cywinska A, Ball SL, deWaard JR. Biological identifications through DNA
393 barcodes. *Proceedings of the Royal Society B: Biological Sciences*. 2003;270:313–21.
- 394 8. Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA, et al. Nuclear ribosomal
395 internal transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi.
396 *Proceedings of the National Academy of Sciences*. 2012;109:6241–6.
- 397 9. Bruns TD, White TJ, Taylor JW. Fungal Molecular Systematics. *Annual Review of Ecology and
398 Systematics*. 1991;22:525–64.
- 399 10. Schüßler A. Glomales SSUrRNA gene diversity. *New Phytologist*. 1999;144:205–7.
- 400 11. James TY, Kauff F, Schoch CL, Matheny PB, Hofstetter V, Cox CJ, et al. Reconstructing the
401 early evolution of Fungi using a six-gene phylogeny. *Nature*. 2006;443:818–22.
- 402 12. Hibbett DS, Binder M, Bischoff JF, Blackwell M, Cannon PF, Eriksson OE, et al. A higher-level
403 phylogenetic classification of the Fungi. *Mycological Research*. 2007;111:509–47.
- 404 13. Zimmermann J, Abarca N, Enk N, Skibbe O, Kusber W-H, Jahn R. Taxonomic Reference
405 Libraries for Environmental Barcoding: A Best Practice Example from Diatom Research. *PLoS
406 ONE*. 2014;9:e108793.
- 407 14. Sato Y, Miya M, Fukunaga T, Sado T, Iwasaki W. MitoFish and MiFish Pipeline: A
408 Mitochondrial Genome Database of Fish with an Analysis Pipeline for Environmental DNA
409 Metabarcoding. *Molecular Biology and Evolution*. 2018;35:1553–5.

- 410 15. Ahmed M, Back MA, Prior T, Karssen G, Lawson R, Adams I, et al. Metabarcoding of soil
411 nematodes: the importance of taxonomic coverage and availability of reference sequences in
412 choosing suitable marker(s). *MBMG*. 2019;3:e36408.
- 413 16. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a
414 Chimera-Checked 16S rRNA Gene Database and Workbench Compatible with ARB. *Applied and*
415 *Environmental Microbiology*. 2006;72:5069–72.
- 416 17. Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, et al. SILVA: a comprehensive
417 online resource for quality checked and aligned ribosomal RNA sequence data compatible with
418 ARB. *Nucleic Acids Research*. 2007;35:7188–96.
- 419 18. Ratnasingham S, Hebert PD. BOLD: The Barcode of Life Data System ([http://www.](http://www.barcodinglife.org)
420 [barcodinglife.org](http://www.barcodinglife.org)). *Molecular ecology notes*. 2007;7:355–364.
- 421 19. Abarenkov K, Nilsson RH, Larsson K-H, Alexander IJ, Eberhardt U, Erland S, et al. The UNITE
422 database for molecular identification of fungi – recent updates and future perspectives. *New*
423 *Phytologist*. 2010;186:281–285.
- 424 20. Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, et al. Ribosomal Database Project:
425 data and tools for high throughput rRNA analysis. *Nucleic Acids Research*. 2014;42:D633–42.
- 426 21. Ankenbrand MJ, Keller A, Wolf M, Schultz J, Förster F. ITS2 Database V: Twice as Much.
427 *Molecular Biology and Evolution*. 2015;32:3030–2.
- 428 22. Rimet F, Chaumeil P, Keck F, Kermarrec L, Vasselon V, Kahlert M, et al. R-Syst::diatom: a
429 barcode database for diatoms and freshwater biomonitoring - data sources and curation
430 procedure. *INRA Report*. 2015.
- 431 23. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing
432 mothur: Open-Source, Platform-Independent, Community-Supported Software for Describing
433 and Comparing Microbial Communities. *Applied and Environmental Microbiology*.
434 2009;75:7537–41.
- 435 24. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME
436 allows analysis of highthroughput community sequencing data. *Nature Methods*. 2010;7:335–6.
- 437 25. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet C, Al-Ghalith GA, et al. Reproducible,
438 interactive, scalable, and extensible microbiome data science using QIIME 2. *Nature*
439 *Biotechnology*. 2019;37:852–7.
- 440 26. Liu K-L, Porras-Alfaro A, Kuske CR, Eichorst SA, Xie G. Accurate, Rapid Taxonomic
441 Classification of Fungal Large-Subunit rRNA Genes. *Appl Environ Microbiol*. 2012;78:1523–33.

- 442 27. Deshpande V, Wang Q, Greenfield P, Charleston M, Porrás-Alfaro A, Kuske CR, et al. Fungal
443 identification using a Bayesian classifier and the Warcup training set of internal transcribed
444 spacer sequences. *Mycologia*. 2016;108:1–5.
- 445 28. Rivers AR, Weber KC, Gardner TG, Liu S, Armstrong SD. ITSxpress: Software to rapidly trim
446 internally transcribed spacer sequences with quality scores for marker gene analysis. *F1000Res*.
447 2018;7:1418.
- 448 29. Porter TM, Gibson JF, Shokralla S, Baird DJ, Golding GB, Hajibabaei M. Rapid and accurate
449 taxonomic classification of insect (class Insecta) cytochrome c oxidase subunit 1 (COI) DNA
450 barcode sequences using a naïve Bayesian classifier. *Mol Ecol Resour*. 2014;14:929–42.
- 451 30. Porter TM, Hajibabaei M. Automated high throughput animal CO1 metabarcode
452 classification. *Scientific Reports*. 2018;8:4226.
- 453 31. Virgilio M, Backeljau T, Nevado B, De Meyer M. Comparative performances of DNA
454 barcoding across insect orders. *BMC bioinformatics*. 2010;11:206.
- 455 32. Maitland VC, Robinson CV, Porter TM, Hajibabaei M. Freshwater diatom biomonitoring
456 through benthic kick-net metabarcoding. *BioRxiv*.
457 2020;: <http://biorxiv.org/lookup/doi/10.1101/2020.05.25.115089>.
- 458 33. Song H, Buhay JE, Whiting MF, Crandall KA. Many species in one: DNA barcoding
459 overestimates the number of species when nuclear mitochondrial pseudogenes are
460 coamplified. *PNAS*. 2008;105:13486–91.
- 461 34. Moulton MJ, Song H, Whiting MF. Assessing the effects of primer specificity on eliminating
462 numt coamplification in DNA barcoding: a case study from Orthoptera (Arthropoda: Insecta):
463 DNA BARCODING. *Molecular Ecology Resources*. 2010;10:615–27.
- 464 35. Baird DJ, Hajibabaei M. Biomonitoring 2.0: a new paradigm in ecosystem assessment made
465 possible by next-generation DNA sequencing. *Molecular ecology*. 2012;21:2039–2044.
- 466 36. Drummond AJ, Newcomb RD, Buckley TR, Xie D, Dopheide A, Potter BC, et al. Evaluating a
467 multigene environmental DNA approach for biodiversity assessment. *GigaSci*. 2015;4:46.
- 468 37. Adamowicz SJ, Boatwright JS, Chain F, Fisher BL, Hogg ID, Leese F, et al. Trends in DNA
469 barcoding and metabarcoding. *Genome*. 2019;62:v–viii.
- 470 38. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, et al. Improved
471 software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and
472 other eukaryotes for analysis of environmental sequencing data. *Methods in Ecology and*
473 *Evolution*. 2013;4:914–9.
- 474 39. Anaconda. Anaconda Software Distribution. 2016. <https://anaconda.com>.

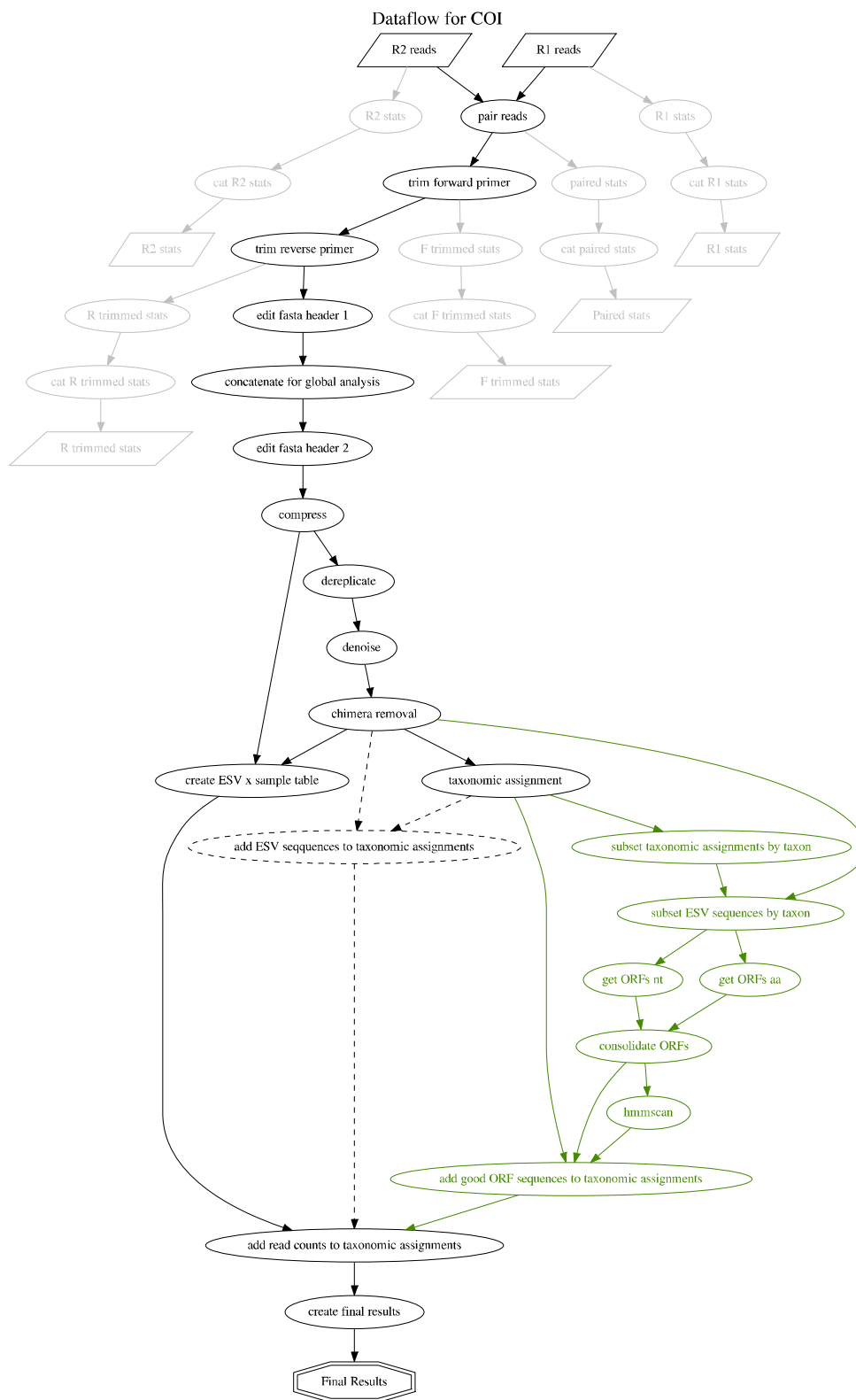
- 475 40. Koster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine.
476 Bioinformatics. 2012;28:2520–2.
- 477 41. St. John J. SeqPrep. 2016. <https://github.com/jstjohn/SeqPrep/releases>.
- 478 42. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
479 EMBnet journal. 2011;17:pp–10.
- 480 43. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for
481 metagenomics. PeerJ. 2016;4:e2584.
- 482 44. Edgar RC. Search and clustering orders of magnitude faster than BLAST. Bioinformatics.
483 2010;26:2460–1.
- 484 45. Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon
485 sequencing. bioRxiv. 2016. doi:10.1101/081257.
- 486 46. Edgar R. UCHIME2: improved chimera prediction for amplicon sequencing. bioRxiv.
487 2016;:074252.
- 488 47. Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational
489 taxonomic units in marker-gene data analysis. The ISME Journal. 2017;11:2639–43.
- 490 48. China Plant BOL Group, Li D-Z, Gao L-M, Li H-T, Wang H, Ge X-J, et al. Comparative analysis
491 of a large dataset indicates that internal transcribed spacer (ITS) should be incorporated into
492 the core barcode for seed plants. Proceedings of the National Academy of Sciences.
493 2011;108:19641–6.
- 494 49. Hollingsworth PM. Refining the DNA barcode for land plants. Proceedings of the National
495 Academy of Sciences. 2011;108:19451–2.
- 496 50. Edge TA, Baird DJ, Bilodeau G, Gagné N, Greer C, Konkin D, et al. The Ecobiomics project:
497 Advancing metagenomics assessment of soil health and freshwater quality in Canada. Science
498 of The Total Environment. 2020;710:135906.
- 499 51. Cochrane G, Karsch-Mizrachi I, Takagi T, Sequence Database Collaboration IN. The
500 International Nucleotide Sequence Database Collaboration. Nucleic Acids Research.
501 2016;44:D48–50.
- 502 52. Iwasaki W, Fukunaga T, Isagozawa R, Yamada K, Maeda Y, Satoh TP, et al. MitoFish and
503 MitoAnnotator: A Mitochondrial Genome Database of Fish with an Accurate and Automatic
504 Annotation Pipeline. Molecular Biology and Evolution. 2013;30:2531–40.
- 505 53. Wheeler DL. Database resources of the National Center for Biotechnology. Nucleic Acids
506 Research. 2003;31:28–33.

- 507 54. R Core Team. R: A Language and Environment for Statistical Computing. 2017.
508 <https://www.R-project.org/>.
- 509 55. Buchner D, Leese F. BOLDigger – a Python package to identify and organise sequences with
510 the Barcode of Life Data systems. MBMG. 2020;4:e53535.
- 511 56. Nugent CM, Adamowicz SJ. Alignment-free identification of COI DNA barcode data with the
512 Python package Alfie. preprint. Bioinformatics; 2020. doi:10.1101/2020.06.29.177634.
- 513 57. Andújar C, Creedy TJ, Arribas P, López H, Salces-Castellano A, Pérez-Delgado A, et al. NUMT
514 dumping: validated removal of nuclear pseudogenes from mitochondrial metabarcoding data.
515 preprint. Evolutionary Biology; 2020. doi:10.1101/2020.06.17.157347.
- 516 58. Nugent CM, Elliott TA, Ratnasingham S, Adamowicz SJ. coil: an R package for cytochrome C
517 oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation. bioRxiv. 2019;:35.
- 518
- 519
- 520

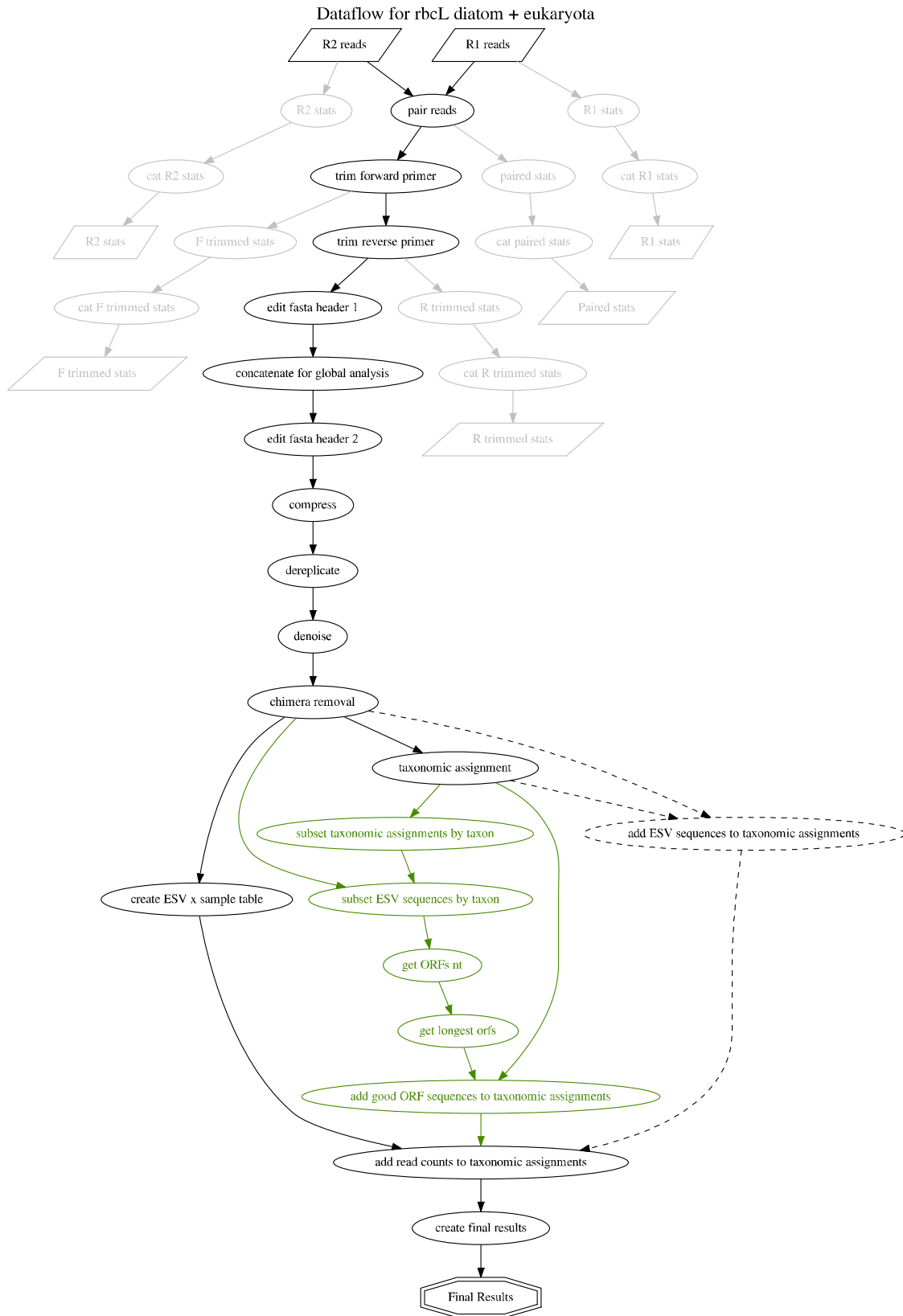
521 **Supplementary Material**

522

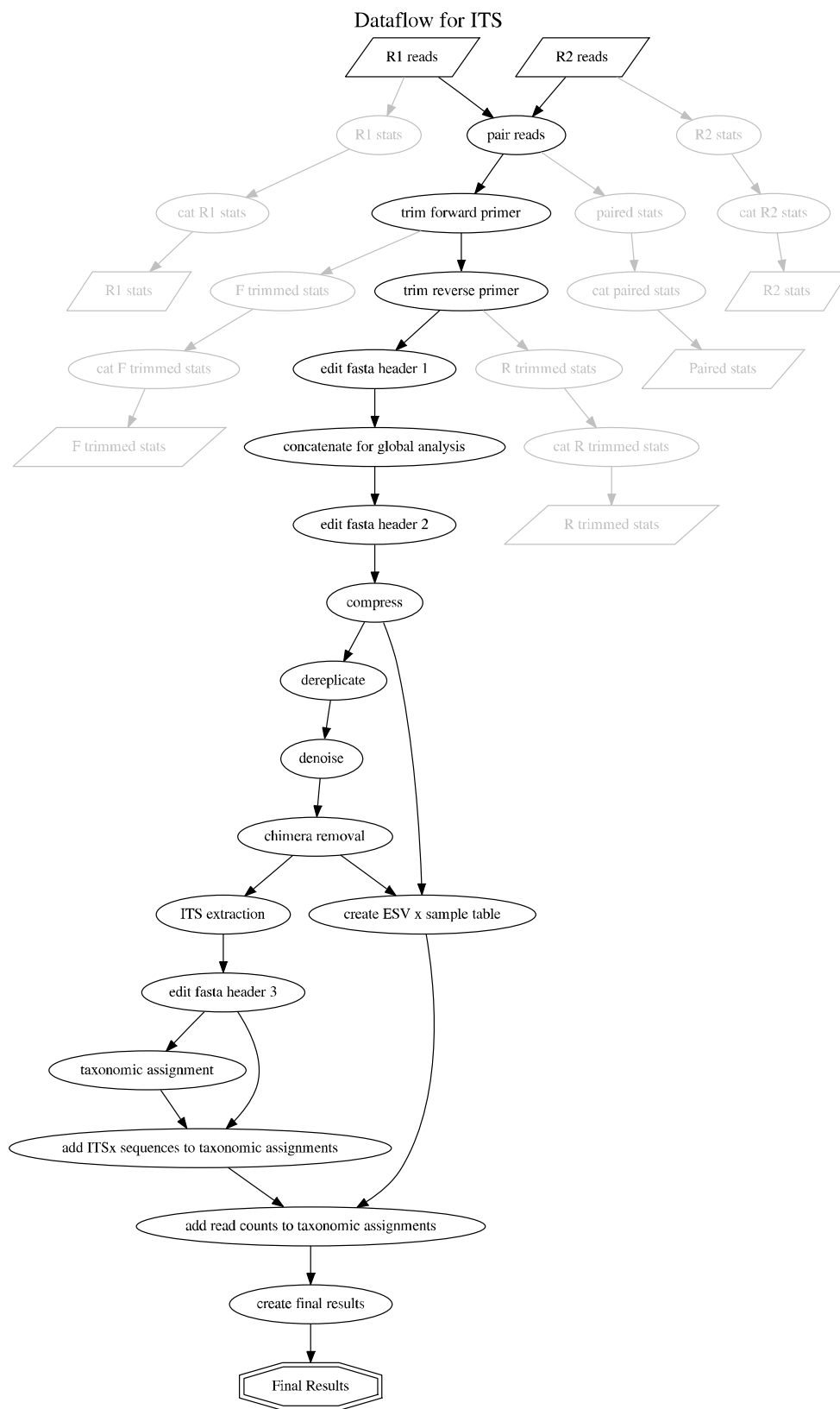
523 **Figure S1. Dataflow for COI mtDNA metabarcodes.** Input and output files are shown
524 as parallelograms. Snakemake rules or processes are shown as ovals. The final
525 results file contains ESVs, for each sample, as well as ESV/ORF sequences, read
526 counts, as well as taxonomic assignments with bootstrap support values. The main
527 dataflow is shown in black, if pseudogene filtering is selected these steps are shown in
528 green, if pseudogene filtering is not selected the dashed steps are performed. The
529 generation of various statistical reports are shown in grey.



531 **Figure S2. Dataflow for rbcL mtDNA metabarcodes.** Input and output files are
532 shown as parallelograms. Snakemake rules or processes are shown as ovals. The
533 final results file contains ESVs, for each sample, as well as ESV/ORF sequences, read
534 counts, as well as taxonomic assignments with bootstrap support values. The main
535 dataflow is shown in black, if pseudogene filtering is selected these steps are shown in
536 green, if pseudogene filtering is not selected the dashed steps are performed. The
537 generation of various statistical reports are shown in grey. The data flow is essentially
538 the same whether the rbcL eukaryota or diatom-specific classifiers are used.



540 **Figure S3. Dataflow for ITS metabarcodes.** Input and output files are shown as
541 parallelograms. Snakemake rules or processes are shown as ovals. The final results
542 file contains ESVs, for each sample, as well as ESV/ORF sequences, read counts, as
543 well as taxonomic assignments with bootstrap support values. The generation of
544 various statistical reports are shown in grey.



546 **Figure S3. Dataflow for rRNA gene metabarcodes.** Input and output files are shown
547 as parallelograms. Snakemake rules or processes are shown as ovals. The final
548 results file contains ESVs, for each sample, as well as ESV/ORF sequences, read
549 counts, as well as taxonomic assignments with bootstrap support values. The
550 generation of various statistical reports are shown in grey. The dataflow is essentially
551 the same whether the RDP classifier built-in 16S reference set is used, or the custom-
552 trained 18S eukaryote or diatom-specific classifiers are used.

Dataflow for 16S + 18S diatom + eukaryota

