

1 **MetaWorks: A flexible, scalable bioinformatic pipeline for high-throughput multi-**  
2 **marker biodiversity assessments**

3

4 Teresita M. Porter<sup>1\*</sup>, Mehrdad Hajibabaei<sup>1</sup>

5

6 <sup>1</sup> Centre for Biodiversity Genomics @ Biodiversity Institute of Ontario & Department of  
7 Integrative Biology, University of Guelph, Guelph, ON, CANADA

8

9 \* Corresponding author:

10 T.M. Porter

11 [terrimporter@gmail.com](mailto:terrimporter@gmail.com)

12

13

## 14 **Abstract**

15

16 Multi-marker metabarcoding is increasingly being used to generate biodiversity  
17 information across different domains of life from microbes to fungi to animals such as for  
18 molecular ecology and biomonitoring applications in different sectors from academic  
19 research to regulatory agencies and industry. Current popular bioinformatic pipelines  
20 support microbial and fungal marker analysis, while ad hoc methods are often used to  
21 process animal metabarcode markers from the same study. MetaWorks provides a  
22 harmonized processing environment, pipeline, and taxonomic assignment approach for  
23 demultiplexed Illumina reads for all biota using a wide range of metabarcoding markers  
24 such as 16S, ITS, and COI. A Conda environment is provided to quickly gather most of  
25 the programs and dependencies for the pipeline. Several workflows are provided such  
26 as: taxonomically assigning exact sequence variants, provides an option to generate  
27 operational taxonomic units, and facilitates single-read processing. Pipelines are  
28 automated using Snakemake to minimize user intervention and facilitate scalability. All  
29 pipelines use the RDP classifier to provide taxonomic assignments with confidence  
30 measures. We extend the functionality of the RDP classifier for taxonomically assigning  
31 16S (bacteria), ITS (fungi), and 28S (fungi), to also support COI (animals), rbcL  
32 (eukaryotes, land plants, diatoms), 12S (fish), 18S (eukaryotes, diatoms) and ITS (fungi,  
33 plants). MetaWorks properly handles ITS by trimming flanking conserved rRNA gene  
34 regions as well as protein coding genes by providing two options for removing obvious  
35 pseudogenes. MetaWorks is available at <https://github.com/terrimporter/MetaWorks>

36 along with quick-start instructions using test data, detailed workflow descriptions, and a  
37 tutorial for new users.

38

### 39 **Keywords**

40 Metabarcoding, Conda, Snakemake, ITS extraction, pseudogene removal, rRNA, ITS,  
41 COI, rbcL, 12S

42

### 43 **Background**

44

45         Marker gene sequencing, metabarcoding, or metasytematics are interrelated  
46 techniques that involves extracting DNA from bulk samples such as soil, water, or  
47 mixtures of individuals collected from traps, without having to isolate or identify  
48 individual specimens, followed by enrichment of a signature DNA region to identify  
49 biological community composition using bioinformatics [1–3]. In microbial ecology to  
50 animal biodiversity studies, different signature DNA regions are targeted for their ability  
51 to identify target taxa. For example, in prokaryotes, the 16S small subunit (SSU)  
52 ribosomal RNA (rRNA) region is often used for genus level taxonomic assignments  
53 [4,5]. Other popular markers include cytochrome c oxidase (COI) for animals; ribulose  
54 bisphosphate large subunit (rbcL) for plants and diatoms; the internal transcribed spacer  
55 (ITS) for fungi and plants; 18S SSU for eukaryotes, arbuscular mycorrhizal fungi, and  
56 diatoms; and 12S mitochondrial SSU for fish [6–14].

57         Existing pipelines such as QIIME2 and DADA2 were initially developed to  
58 support the microbial ecology community [15,16]. In comprehensive, multi-trophic,

59 multi-marker studies, there is a need for a pipeline that can handle rRNA genes, spacer  
60 regions, as well as protein-coding markers in a single harmonized environment [17,18].  
61 For the ITS region, we needed a pipeline that could remove the conserved flanking  
62 rRNA genes as this has been shown to improve taxonomic assignment accuracy [19].  
63 For protein-coding regions, we needed a pipeline that could remove putative  
64 pseudogenes [20–23]. We also wanted the ability to generate high quality exact  
65 sequence variants (ESVs) for any marker (not just 16S or ITS) for the additional level of  
66 genetic and taxonomic resolution SVs can provide [24–26]. For taxonomic assignment,  
67 we wanted to use a classifier that would provide a measure of confidence for  
68 assignments to reduce false-positive assignments [27–29].

69 As multi-marker studies are carried out on phylogenetically divergent taxa, such  
70 as in biodiversity or trophic studies, there is a need for more generic pipelines where  
71 different markers can be analyzed using similar dataflows with 3rd party programs  
72 instead of being limited to database-specific pipelines and tools [17,30]. We developed  
73 MetaWorks with the following objectives: 1) reproducibility with respect to the  
74 computational environment used as well as the pipeline itself, 2) scalability to leverage  
75 high performance computer clusters to speed up the analysis of large datasets, 3) naive  
76 Bayes classifier support for popular metabarcode markers; and 4) to support marker-  
77 specific processing steps such as ITS extraction and pseudogene-removal for protein-  
78 coding markers. MetaWorks was designed for data analysts who are comfortable using  
79 Linux command-line tools but would like a single harmonized environment and pipeline  
80 to process multi-marker metabarcode datasets.

81

## 82 **Implementation & Workflow**

83

### 84 *Implementation*

85         MetaWorks is a multi-marker ‘meta’-barcode pipeline that does ‘the works’ by  
86 supporting the bioinformatic processing of popular markers including rRNA genes,  
87 spacers, and protein coding genes generating taxonomically assigned sequence  
88 variants or operational taxonomic units (OTUs). To facilitate reproducibility, scalability,  
89 and shareability of workflows we use the Conda package manager to facilitate the  
90 download of most programs and dependencies and the Snakemake workflow manager  
91 to automate pipelines and utilize computational resources efficiently [31–33].  
92 Snakemake supports re-entrancy and automatic deployment of multiple parallel jobs,  
93 both ideal for high performance computing environments where many cores are  
94 available to speed up the analysis of large datasets.

95         We provided instructions on how to install and use Conda in the online  
96 documentation. One additional program not available as a Conda package, ORFfinder,  
97 may need to be downloaded separately if pseudogene-filtering will be conducted and  
98 instructions are provided in the online documentation. MetaWorks is available from  
99 <https://github.com/terrimporter/MetaWorks> and a suite of trained classifiers for  
100 taxonomic assignment are also available from GitHub ([Table 1](#)). Depending on the  
101 DNA metabarcode marker(s) the user will be processing, these can be individually  
102 downloaded from GitHub and instructions are provided in the online documentation.

103

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

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

			<b>included sequences</b>	<b>included taxa at all ranks (species)</b>	
COI	Eukaryotes	<a href="https://github.com/terrimp/orter/CO1Classifier">https://github.com/terrimp/orter/CO1Classifier</a>	1,221,528	154,351 (114,687)	BOLD [34], INSDC [35]
rbcL	Diatoms	<a href="https://github.com/terrimp/orter/rbcLdiatomClassifier">https://github.com/terrimp/orter/rbcLdiatomClassifier</a>	3,504	1,432 (1,023)	Diat.barcode [36]
rbcL	Land plants	<a href="https://github.com/terrimp/orter/rbcL_landPlant_Classifier">https://github.com/terrimp/orter/rbcL_landPlant_Classifier</a>	148,258	61,398 (50,778)	INSDC [35]
rbcL	Eukaryotes	<a href="https://github.com/terrimp/orter/rbcLClassifier">https://github.com/terrimp/orter/rbcLClassifier</a>	164,454	65,742 (53,344)	INSDC [35]
12S	Fish	<a href="https://github.com/terrimp/orter/12SfishClassifier">https://github.com/terrimp/orter/12SfishClassifier</a>	2,853	4,751 (2,833)	MitoFish [37]
12S	Vertebrates	<a href="https://github.com/terrimp/orter/12SvertebrateClassifier">https://github.com/terrimp/orter/12SvertebrateClassifier</a>	10,654	15,007 (9,564)	INSDC [35] and MitoFish [37]
SSU (18S)	Diatoms	<a href="https://github.com/terrimp/orter/SSUdiatomClassifier">https://github.com/terrimp/orter/SSUdiatomClassifier</a>	2,962	1,198 (828)	Diat.barcode [36]
SSU (16S)	Vertebrates	<a href="https://github.com/terrimp/orter/16SvertebrateClassifier">https://github.com/terrimp/orter/16SvertebrateClassifier</a>	72,195	21,282 (15,155)	INSDC [35]
SSU (18S)	Eukaryotes	<a href="https://github.com/terrimp/orter/18SClassifier">https://github.com/terrimp/orter/18SClassifier</a>	42,301	7,504 (5,440 genera)	SILVA [38]
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)	Deshpande et al., 2016 [39]
ITS	Fungi (UNITE 2014)	Built-in to the RDP classifier	145,019	23,222 (20,337)	Abarenkov et al., 2010 [40]
ITS	Fungi (UNITE 2021)	<a href="https://github.com/terrimp/orter/UNITE_ITSClifier">https://github.com/terrimp/orter/UNITE_ITSClifier</a>	1,393,203	376,167 (352,588)	UNITE [40]
ITS	Plants	<a href="https://github.com/terrimp/orter/PLANIITS_ITSClifier">https://github.com/terrimp/orter/PLANIITS_ITSClifier</a>	104,387	72,632 (61,693)	PLANIITS [41] and UNITE [40]
LSU	Fungi	Built-in to the RDP classifier	11,442	2,633 (1,895)	Liu et al., 2012 [42]

105

106

107

## 108 *Workflow*

109 The pipeline begins with demultiplexed Illumina paired-end reads as this is the

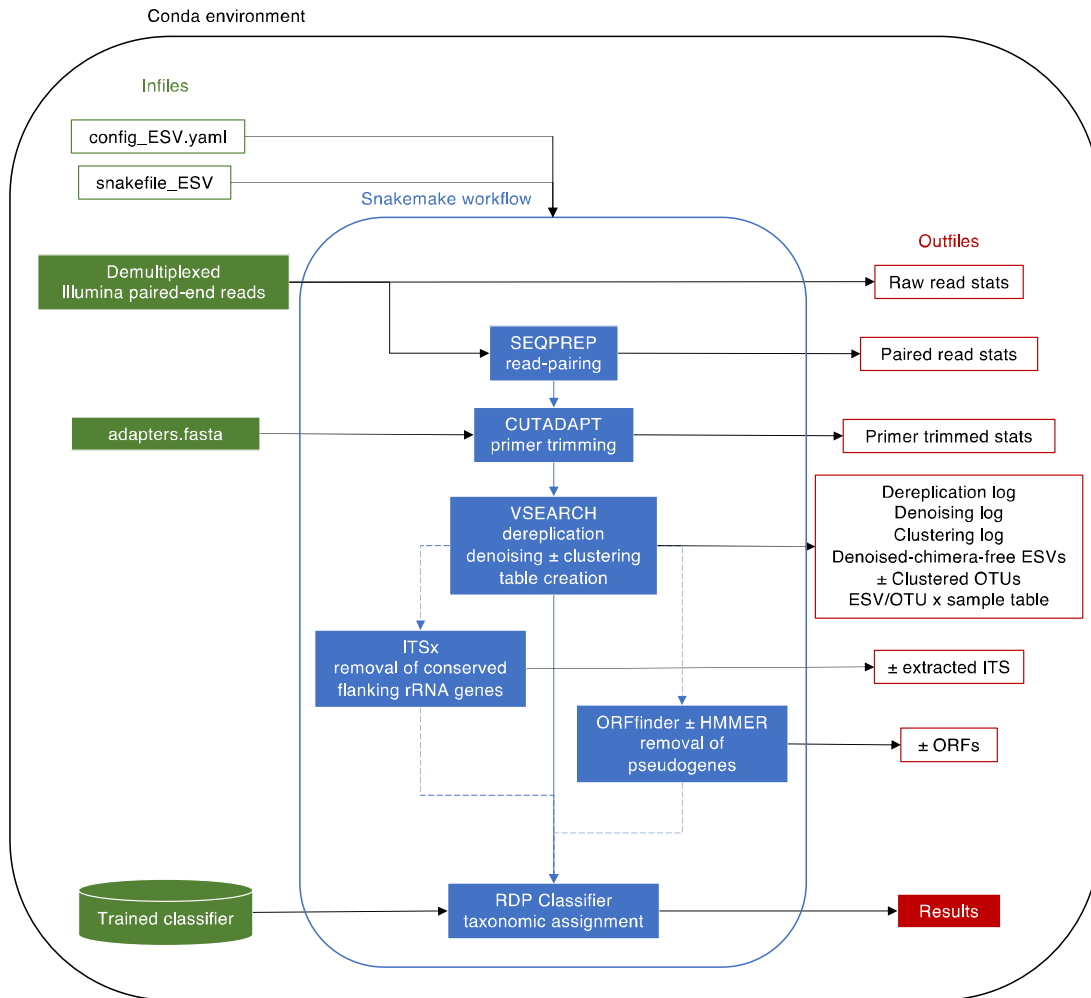
110 format most often provided by sequencing centres to their clients. Several workflows

111 are available as Snakemake pipelines such as taxonomic assignment of sequence

112 variants ([Fig 1](#)), clustering of sequence variants into OTUs, or for processing single  
113 reads. For each of these workflows described below, parameter settings for each  
114 bioinformatic step can be customized in the config.yaml file. The user also needs to  
115 provide a file of primer sequences so we provide a template for the adapters.fasta file  
116 as well as a small set of raw Illumina sequences for the COI amplicon that can be used  
117 to test the installation. The online documentation provides a 'quick-start' example using  
118 the provided COI test data. We also provide a more detailed tutorial in the online  
119 documentation that walks users through the steps necessary to set up their environment  
120 to run the pipeline for the first time, assuming the user has never worked with Conda or  
121 Snakemake before.

122

123 **Figure 1. MetaWorks workflow to produce taxonomically assigned sequence**  
124 **variants.** To aid reproducibility, a Conda environment is provided (environment.yml).  
125 Although multiple Snakemake workflows are provided in MetaWorks, here we show the  
126 workflow to produce taxonomically assigned sequence variants. User input is shown in  
127 green on the left, the ESV workflow is shown in the middle in blue, and outfiles are  
128 shown in red on the right. The results are provided in a comma-separated value (CSV)  
129 file and shows each ESV/OTU per sample with read counts and taxonomic  
130 assignments. Abbreviations: internal transcribed spacer (ITS) region, open reading  
131 frame FASTA-formatted sequences (ORFs).



132

133

### 134 *Generation and taxonomic assignment of sequence variants*

135 The ESV workflow will run the pipeline shown in Fig 1. In the config\_ESV.yaml

136 file, users indicate the path to the directory that contains the demultiplexed Illumina

137 paired-end reads, specify the unique part of filenames to distinguish between samples

138 and reads, and specify the name of the directory that will contain the outfiles. Default



139 settings for each program are provided in the config\_ESV.yaml file but these can be  
140 customized by the user. SEQPREP was initially chosen for pairing the forward and  
141 reverse reads, because the program comes with the option to output the alignments for  
142 visual inspection, an option that most read-pairing programs do not have [43]. For  
143 SEQPREP read-pairing, users can specify a Phred score quality cutoff, the minimum  
144 overlap between the forward and reverse reads, the maximum fraction of mismatches  
145 allowed in the overlap region, and the minimum fraction of matching bases in the  
146 overlap region. CUTADAPT was chosen for primer-trimming because it is fast and  
147 already widely used in the metabarcoding community for this purpose, so most users  
148 will likely already be familiar with how this program works [44]. For CUTADAPT, users  
149 need to provide a FASTA-formatted primer sequence file (adapters.fasta), they can also  
150 specify the minimum sequence length to retain after primer-trimming, a Phred quality  
151 score cutoff, the maximum error rate, minimum adapter overlap, and maximum number  
152 of ambiguous bases allowed. VSEARCH was chosen to dereplicate reads (retain  
153 unique reads) and remove artefactual sequences using the UNOISE3 and UCHIME3  
154 algorithms [45,46]. We chose the open-source VSEARCH program over alternatives  
155 because the program can utilize all the available memory on a system, facilitating the  
156 analysis of large datasets on high performance computer systems. We prefer the  
157 UNOISE3 method for denoising because it performs up to 1,200 faster and uses less  
158 memory than other denoising programs [47]. To map read counts to the newly  
159 generated denoised-chimera ESVs to create an ESV x sample table, we use the  
160 'search\_exact' method because it is faster and optimized to find exact matches

161 compared with the 'usearch\_global' command with the 'id 1.0' parameter, but this is just  
162 an intermediate step and further filtering of this table is performed by MetaWorks.

163 If the internal transcribed spacer (ITS) region is analyzed, then the pipeline uses  
164 the ITSx program to trim away the flanking conserved rRNA gene sequences so that  
165 taxonomic assignment is based solely on the variable spacer region sequences (ITS1 or  
166 ITS2) [19]. This step has been shown to improve sensitivity of clustering and taxonomic  
167 assignments [19].

168 If a protein coding marker is being processed, the user can select a pseudogene-  
169 removal method in the config\_ESV.yaml file. We have previously described two  
170 methods for removing putative pseudogenes from DNA barcode and metabarcoding  
171 datasets [21]. The NCBI ORFfinder program is used to translate reads into all possible  
172 open reading frames (ORFs). The first pseudogene removal method retains the longest  
173 ORF for each read, calculates a distribution of ORF lengths, and removes reads with  
174 outlier lengths as putative ORFs. The second pseudogene removal method can be  
175 used if a hidden Markov model is available and is provided for processing COI  
176 arthropods. The longest ORFs are compared to the profile using HMMER available  
177 from <http://hmmer.org>. MetaWorks calculates a distribution of bit scores and removes  
178 reads with short outlier bit scores as putative ORFs. Removing noise caused by the  
179 sequencing of pseudogenes in metabarcoding datasets can help users avoid over-  
180 estimating richness in subsequent analyses, yet this step is not included in the most  
181 popular metabarcoding pipelines as they were developed to support the analysis of rRNA  
182 genes where this is not a problem.

183           One of the features of MetaWorks, is the use of a single taxonomic assignment  
184 method for any metabarcode marker that provides a measure of confidence for  
185 taxonomic assignments. We chose the RDP Classifier for this task as this method has  
186 a long-history of use in the microbial ecology literature, additionally the classifier can be  
187 customized and validated for any metabarcode marker [5]. The RDP classifier  
188 calculates k-mer frequencies and uses a naive Bayes method to taxonomically assign  
189 unknown query sequences. Bootstrapping is used to provide a measure of statistical  
190 support, or repeatability, for each assignment at each rank. We have previously  
191 described how this method works compared to the top BLAST hit method [28]. In that  
192 comparison, we showed how the RDP classifier is faster than the top BLAST hit method  
193 and helps to reduce the rate of false-positive assignments. In studies where  
194 erroneously identifying a metabarcode sequence as a potential invasive species or  
195 pathogen could lead to alarm, reducing the false-positive assignment rate is critical. We  
196 provide a suite of trained classifiers, ready for use with MetaWorks (Table 1).  
197 Additionally, we provide the training files so that users can check that key target taxa  
198 are present in the reference database, and users are free to use the FASTA-formatted  
199 sequence files to create custom BLAST databases for similarity-based searches for  
200 data exploration or to build reference sets for subsequent phylogenetic analysis. The  
201 final file is a comma-separated value file (results.csv) where the taxonomic assignment  
202 for each sequence variant is provided for each sample along with read counts. If an  
203 rRNA marker was processed, then the sequence variant sequence is provided in this  
204 file; and if a protein coding region was processed, then longest ORF is provided.  
205

## 206 *Operational taxonomic units*

207           This pipeline supports the analysis of sequence variants for the additional  
208 genetic and taxonomic resolution provided by this level of analysis [24]. Though this  
209 method of analysis was initially used to process 16S rRNA genes, studies using ITS and  
210 COI have also shown that the analysis of SVs improves the detection of genetic  
211 diversity and richness, when assessing beta diversity, both SVs and OTUs tend to  
212 recover similar gradients in multivariate analyses [25,26]. Although it has been shown  
213 that for many clustering methods sequence order matters and OTU composition can  
214 change from one analysis to the next making reproducibility an issue, there are a  
215 number of reasons why a user would want to analyze OTUs. For example, it may be  
216 more advantageous to work with OTUs instead of sequence variants for network  
217 analysis to detect more co-occurrences, for legacy reasons to compare results to  
218 previous studies that used OTUs, or to approximate ‘species’ units [48].

219           After processing raw reads using the snakefile\_ESV workflow described in the  
220 previous section, users cluster sequence variants into OTUs, a method that combines  
221 the benefits of denoising with clustering using a 97% sequence similarity cutoff using  
222 the snakefile\_OTU workflow [26,49]. This method uses VSEARCH ‘cluster\_smallmem’  
223 method to cluster SVs using a 97% sequence similarity cutoff. Settings can be adjusted  
224 in the in the config\_OTU.yaml file such as pointing to the directory that contains the  
225 sequence variants and choosing a classifier for the OTUs.

226

## 227 **Results and Discussion**

228           MetaWorks has already been used in several publications for the Canadian  
229   STREAM biomonitoring program, the Government of Canada, Genomics Research and  
230   Development Initiative, Metagenomics-based ecosystem biomonitoring (Ecobiomics)  
231   project, and by Natural Resources Canada [18,50,51]. The benefits of using an  
232   automated, scalable, versioned pipeline for biomonitoring are many-fold, from the ability  
233   to share reproducible workflows with collaborators to facilitate the re-analysis of data as  
234   more samples are collected from year-to-year. Below, we describe three MetaWorks  
235   use-cases in more detail below.

236  
237   *Use case 1:* As a part of the Canadian STREAM biomonitoring Initiative, the MetaWorks  
238   pipeline has been used to process macroinvertebrates COI metabarcodes surveyed  
239   from stream sites across Canada [50]. One feature of this project is the quick 1-2  
240   month turn-around time from sampling through to the production of watershed  
241   biodiversity reports. This is an improvement over reports generated using conventional  
242   morphology-based methods that would normally take 6 - 12 months to produce. The  
243   use of a consistent bioinformatics workflow to process metabarcodes has played a key  
244   role in the reproducibility, scalability, and throughput to facilitate timely reporting [52].  
245   Generally, samples are processed in batches of 96 per sequencing run then later split  
246   into custom reports for stakeholders, processing about 500 samples per year. One  
247   feature of these reports are the taxonomic assignments made using the naive Bayesian  
248   classifier that can be filtered for identifications to the species, genus, and family ranks  
249   with 99% accuracy assuming the query is represented in the underlying sequence  
250   database [28]. This is in contrast with the use of more traditional methods for taxonomic

251 assignment, where taxa are routinely missed during subsampling and taxa detected by  
252 primary analysts and auditors may differ by up to 30% [53]. This use-case shows how  
253 MetaWorks can be used to create taxon lists for large-scale biodiversity monitoring of  
254 streams across Canada.

255  
256 *Use case 2:* Also as a part of the STREAM project described above, MetaWorks results  
257 were to analyze sequence variants from diatoms (rbcL) and arthropods (COI) sampled  
258 within and across sites of varying water quality [54]. Using the MetaWorks pipeline, two  
259 different protein-coding markers were bioinformatically processed in two runs. The first  
260 run processed the rbcL marker, using a mixture of 5 different primers in a single  
261 adapters.fasta file, and pseudogenes were removed from this dataset using the simple  
262 ORFfinder method [21]. The second run processed three COI amplicons, each  
263 targeting an approximately 200 bp length of the COI barcoding region using 6 different  
264 primers in a single adapters.fasta file, and pseudogenes were removed from this  
265 dataset using the ORFfinder+HMMER method since a COI arthropod HMM model was  
266 available [21]. The study reported a diversity assessment across sites of varying water  
267 quality using richness, effective richness, and beta diversity. Additionally, the taxonomic  
268 assignments generated from MetaWorks were used to obtain resource-consumer  
269 relationships from a global database of biotic interactions (GloBI) so that community  
270 stability using trophic and network measures could be assessed across sites with  
271 varying water quality [55]. This use-case shows how MetaWorks can handle a variety  
272 of protein-coding markers for trophic and network analyses to facilitate ecological  
273 assessments of freshwater condition.

274  
275 *Use case 3:* As a part of collaborative work with Environment and Climate Change  
276 Canada, MetaWorks was used to assess macroinvertebrate and (non-  
277 macroinvertebrate) eukaryote taxa in an urban harbour using COI and 18S rRNA [56].  
278 Using the MetaWorks pipeline, COI metabarcodes were identified down to species rank  
279 with 99% accuracy and 18S metabarcodes were identified to genus rank with 80%  
280 accuracy using a custom-trained classifier based on the SILVA 18S release 138 [38]. In  
281 this study, conventional macroinvertebrate sampling for assessing water quality in  
282 Toronto Harbour was compared with metabarcoding methods. COI metabarcoding was  
283 found to detect more diversity at a finer level of taxonomic resolution compared with  
284 conventional approaches and was able to distinguish sites with particularly high levels  
285 of sediment contaminants. Additionally, the use of a multi-marker approach allowed  
286 microscopic eukaryote diversity to be sampled at the same time from the same  
287 samples, producing indicators that responded to gradients in both sediment  
288 contaminants and water physical-chemical features. This use-case illustrates how  
289 MetaWorks can facilitate the application of multi-marker metabarcoding approaches that  
290 target different domains of life.

291       As demonstrated in the above examples, MetaWorks supports a wide range of  
292 analysis scenarios from metabarcoding data. We envision that MetaWorks will aid  
293 broader user communities and fill a need in multi-marker metabarcoding studies that  
294 target taxa from multiple different domains of life, to provide a unified processing  
295 environment, pipeline, and taxonomic assignment approach for each marker from  
296 ribosomal RNA genes, spacers, or protein coding genes. QIIME2 is perhaps the most

297 popular and comprehensive platform for such work, but to date, focuses on processing  
298 mainly prokaryote and fungal datasets [16]. As of yet, MetaWorks is the only  
299 bioinformatic pipeline that can handle rRNA genes but that also integrates special  
300 processing steps to handle ITS spacers as well as filter out obvious pseudogenes in  
301 protein coding markers such as COI.

302       There has been a lot of activity with respect to building new bioinformatic tools to  
303 handle COI metabarcodes. Recent work, such as the BOLDigger program, has  
304 attempted to make the BOLD identification engine more suitable for identifying large  
305 batches of COI metabaracodes and has both GUI and command-line interfaces for  
306 efficient sample processing [57]. A new program, called NUMTdumper, has been  
307 developed as a stand-alone program meant to be incorporated into bioinformatic  
308 pipelines [20]. NUMTdumper provides a method to screen for NuMTs based on read  
309 counts while acknowledging the trade-offs between removing all possible NuMTs while  
310 erroneously removing genuine reads. An R package called 'coil' has also recently been  
311 developed that will place COI barcode and metabarcode sequences in frame using  
312 profile HMM analysis [58]. MetaWorks aims to extend the COI metabarcode toolkit that  
313 provides a harmonized environment where data from other organismal markers in multi-  
314 marker, multi-trophic studies can also be analyzed.

315

## 316 **Conclusion**

317

318       MetaWorks is provided as free and open software that is versioned, can be easily  
319 deployed in a Conda environment, and is supported by a suite of classifiers for popular



320 metabarcoding markers. The software comes with a small set of raw data and a quick-  
321 start pipeline to help users gain experience quickly. There is extensive online  
322 documentation available at the link above included detailed explanations of the pipeline,  
323 available workflows, and a tutorial for new users who have never used Conda or  
324 Snakemake before. MetaWorks generates a CSV file that lists all sequence clusters,  
325 for each sample, with associated read counts, taxonomic assignments, and bootstrap  
326 support values. Numerous statistics and log files are also provided so that users can  
327 track the number of reads that pass each major bioinformatic step. Given the current  
328 use of MetaWorks by large-scale national initiatives such as STREAM and Ecobiomics,  
329 we foresee additional developments and enhancements. Future improvements include  
330 the development of additional HMM models for pseudogene filtering, updated and  
331 additional classifiers for taxonomic assignment, and support for processing larger jobs  
332 both on HPCs and in a cloud environment. We welcome suggestions and potential  
333 collaborative work to further advance this pipeline for the scientific community.

334

### 335 **Availability and requirements**

336 MetaWorks runs at the command-line in a Conda environment on a Linux

337 operating system. Miniconda can be downloaded from

338 [https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86\\_64.sh](https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh). If

339 pseudogene filtering will be performed ORFfinder can be downloaded from the NCBI ftp

340 site at <ftp://ftp.ncbi.nlm.nih.gov/genomes/TOOLS/ORFfinder/linux-i64/ORFfinder.gz>.

341 MetaWorks workflows, quick-start example, and a tutorial for new users is available

342 from GitHub at <https://github.com/terrimporter/MetaWorks>. Additional trained classifiers

343 that can be used by MetaWorks for taxonomic assignment are available from GitHub at  
344 <https://github.com/terrimporter>.

345

### 346 **Funding**

347 MH received funding from Genome Canada and Ontario Genomics for the Sequencing  
348 the Rivers for Environmental Assessment and Monitoring (STREAM) project. TMP  
349 received funding from the Government of Canada through the Genomics Research and  
350 Development Initiative (GRDI), Metagenomics-based ecosystem biomonitoring  
351 (Ecobiomics) project.

352

### 353 **Authors' contributions**

354 TP and MH conceived of the work. TP designed and wrote the pipeline. All the authors  
355 helped to write the manuscript, approve the submitted version, and agree to be  
356 personally accountable for the author's own contributions and to ensure that questions  
357 related to the accuracy or integrity of any part of the work, even ones in which the  
358 author was not personally involved, are appropriately investigated, resolved, and the  
359 resolution documented in the literature.

360

### 361 **Acknowledgements**

362 We would like to thank Josip Rudar, Katie M. McGee, Chloe V. Robinson, Victoria C.  
363 Maitland, and Michael T.G. Wright from the Hajibabaei lab for helpful discussions and  
364 testing the pipeline with various datasets.

## 365 References

- 366 1. Pace NR. A Molecular View of Microbial Diversity and the Biosphere. *Science*. 1997;276:  
367 734–740. doi:10.1126/science.276.5313.734
- 368 2. Hajibabaei M. The golden age of DNA metasystematics. *Trends in genetics*. 2012;28: 535–  
369 537.
- 370 3. Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E. Towards next-generation  
371 biodiversity assessment using DNA metabarcoding. *Molecular ecology*. 2012;21: 2045–  
372 2050.
- 373 4. Gilbert JA, Jansson JK, Knight R. The Earth Microbiome project: successes and aspirations.  
374 *BMC biology*. 2014;12: 69.
- 375 5. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid Assignment of  
376 rRNA Sequences into the New Bacterial Taxonomy. *Applied and Environmental*  
377 *Microbiology*. 2007;73: 5261–5267. doi:10.1128/AEM.00062-07
- 378 6. Schüßler A. Glomales SSUrRNA gene diversity. *New Phytologist*. 1999;144: 205–207.  
379 doi:10.1046/j.1469-8137.1999.00526.x
- 380 7. Hebert PDN, Cywinska A, Ball SL, deWaard JR. Biological identifications through DNA  
381 barcodes. *Proceedings of the Royal Society B: Biological Sciences*. 2003;270: 313–321.  
382 doi:10.1098/rspb.2002.2218
- 383 8. James TY, Kauff F, Schoch CL, Matheny PB, Hofstetter V, Cox CJ, et al. Reconstructing the  
384 early evolution of Fungi using a six-gene phylogeny. *Nature*. 2006;443: 818–822.  
385 doi:10.1038/nature05110
- 386 9. Hibbett DS, Binder M, Bischoff JF, Blackwell M, Cannon PF, Eriksson OE, et al. A higher-  
387 level phylogenetic classification of the Fungi. *Mycological Research*. 2007;111: 509–547.  
388 doi:10.1016/j.mycres.2007.03.004
- 389 10. 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*  
391 *Academy of Sciences*. 2009;106: 12794–12797.
- 392 11. Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA, et al. Nuclear  
393 ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for  
394 Fungi. *Proceedings of the National Academy of Sciences*. 2012;109: 6241–6246.  
395 doi:10.1073/pnas.1117018109
- 396 12. Zimmermann J, Abarca N, Enk N, Skibbe O, Kusber W-H, Jahn R. Taxonomic Reference  
397 Libraries for Environmental Barcoding: A Best Practice Example from Diatom Research.  
398 Schierwater B, editor. *PLoS ONE*. 2014;9: e108793. doi:10.1371/journal.pone.0108793

- 399 13. Sato Y, Miya M, Fukunaga T, Sado T, Iwasaki W. MitoFish and MiFish Pipeline: A  
400 Mitochondrial Genome Database of Fish with an Analysis Pipeline for Environmental DNA  
401 Metabarcoding. Kumar S, editor. *Molecular Biology and Evolution*. 2018;35: 1553–1555.  
402 doi:10.1093/molbev/msy074
- 403 14. Ahmed M, Back MA, Prior T, Karssen G, Lawson R, Adams I, et al. Metabarcoding of soil  
404 nematodes: the importance of taxonomic coverage and availability of reference sequences  
405 in choosing suitable marker(s). *MBMG*. 2019;3: e36408. doi:10.3897/mbmg.3.36408
- 406 15. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-  
407 resolution sample inference from Illumina amplicon data. *Nature Methods*. 2016;13: 581–  
408 583. doi:10.1038/nmeth.3869
- 409 16. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet C, Al-Ghalith GA, et al. Reproducible,  
410 interactive, scalable, and extensible microbiome data science using QIIME 2. *Nature*  
411 *Biotechnology*. 2019;37: 852–857. doi:10.7287/peerj.preprints.27295v1
- 412 17. Drummond AJ, Newcomb RD, Buckley TR, Xie D, Dopheide A, Potter BC, et al. Evaluating a  
413 multigene environmental DNA approach for biodiversity assessment. *GigaSci*. 2015;4: 46.  
414 doi:10.1186/s13742-015-0086-1
- 415 18. Edge TA, Baird DJ, Bilodeau G, Gagné N, Greer C, Konkin D, et al. The Ecobiomics project:  
416 Advancing metagenomics assessment of soil health and freshwater quality in Canada.  
417 *Science of The Total Environment*. 2020;710: 135906.  
418 doi:10.1016/j.scitotenv.2019.135906
- 419 19. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, et al. Improved  
420 software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi  
421 and other eukaryotes for analysis of environmental sequencing data. Bunce M, editor.  
422 *Methods in Ecology and Evolution*. 2013;4: 914–919. doi:10.1111/2041-210X.12073
- 423 20. Andújar C, Creedy TJ, Arribas P, López H, Salces-Castellano A, Pérez-Delgado A, et al.  
424 NUMT dumping: validated removal of nuclear pseudogenes from mitochondrial  
425 metabarcode data. *Evolutionary Biology*; 2020 Jun. doi:10.1101/2020.06.17.157347
- 426 21. Porter TM, Hajibabaei M. Profile hidden Markov model sequence analysis can help remove  
427 putative pseudogenes from DNA barcoding and metabarcoding datasets. *BMC*  
428 *Bioinformatics*. 2021;22: 256. doi:10.1186/s12859-021-04180-x
- 429 22. Song H, Buhay JE, Whiting MF, Crandall KA. Many species in one: DNA barcoding  
430 overestimates the number of species when nuclear mitochondrial pseudogenes are  
431 coamplified. *PNAS*. 2008;105: 13486–13491. doi:10.1073/pnas.0803076105
- 432 23. Moulton MJ, Song H, Whiting MF. Assessing the effects of primer specificity on eliminating  
433 numt coamplification in DNA barcoding: a case study from Orthoptera (Arthropoda):

- 434 Insecta): DNA BARCODING. *Molecular Ecology Resources*. 2010;10: 615–627.  
435 doi:10.1111/j.1755-0998.2009.02823.x
- 436 24. Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational  
437 taxonomic units in marker-gene data analysis. *The ISME Journal*. 2017;11: 2639–2643.
- 438 25. Glassman SI, Martiny JB. Ecological patterns are robust to use of exact sequence variants  
439 versus operational taxonomic units. *mSphere*. 2018;3: e00148-18. doi:10.1101/283283
- 440 26. Porter TM, Hajibabaei M. Putting COI Metabarcoding in Context: The Utility of Exact  
441 Sequence Variants (ESVs) in Biodiversity Analysis. *Front Ecol Evol*. 2020;8: 248.  
442 doi:10.3389/fevo.2020.00248
- 443 27. Porter TM, Gibson JF, Shokralla S, Baird DJ, Golding GB, Hajibabaei M. Rapid and accurate  
444 taxonomic classification of insect (class Insecta) cytochrome c oxidase subunit 1 (COI) DNA  
445 barcode sequences using a naïve Bayesian classifier. *Mol Ecol Resour*. 2014;14: 929–942.  
446 doi:10.1111/1755-0998.12240
- 447 28. Porter TM, Hajibabaei M. Automated high throughput animal CO1 metabarcode  
448 classification. *Scientific Reports*. 2018;8: 4226. doi:[https://doi.org/10.1038/s41598-018-](https://doi.org/10.1038/s41598-018-22505-4)  
449 [22505-4](https://doi.org/10.1038/s41598-018-22505-4)
- 450 29. Virgilio M, Backeljau T, Nevado B, De Meyer M. Comparative performances of DNA  
451 barcoding across insect orders. *BMC bioinformatics*. 2010;11: 206.
- 452 30. Adamowicz SJ, Boatwright JS, Chain F, Fisher BL, Hogg ID, Leese F, et al. Trends in DNA  
453 barcoding and metabarcoding. *Genome*. 2019;62: v–viii. doi:10.1139/gen-2019-0054
- 454 31. Koster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine.  
455 *Bioinformatics*. 2012;28: 2520–2522. doi:10.1093/bioinformatics/bts480
- 456 32. Anaconda. Anaconda Software Distribution. 2016. Available: <https://anaconda.com>
- 457 33. Wratten L, Wilm A, Göke J. Reproducible, scalable, and shareable analysis pipelines with  
458 bioinformatics workflow managers. *Nat Methods*. 2021 [cited 24 Sep 2021].  
459 doi:10.1038/s41592-021-01254-9
- 460 34. Ratnasingham S, Hebert PD. BOLD: The Barcode of Life Data System ([http://www.](http://www.barcodinglife.org)  
461 [barcodinglife.org](http://www.barcodinglife.org)). *Molecular ecology notes*. 2007;7: 355–364.
- 462 35. Cochrane G, Karsch-Mizrachi I, Takagi T, Sequence Database Collaboration IN. The  
463 International Nucleotide Sequence Database Collaboration. *Nucleic Acids Research*.  
464 2016;44: D48–D50. doi:10.1093/nar/gkv1323

- 465 36. Rimet F, Gusev E, Kahlert M, Kelly MG, Kulikovskiy M, Maltsev Y, et al. Diat.barcode, an  
466 open-access curated barcode library for diatoms. *Sci Rep.* 2019;9: 15116.  
467 doi:10.1038/s41598-019-51500-6
- 468 37. Iwasaki W, Fukunaga T, Isagozawa R, Yamada K, Maeda Y, Satoh TP, et al. MitoFish and  
469 MitoAnnotator: A Mitochondrial Genome Database of Fish with an Accurate and  
470 Automatic Annotation Pipeline. *Molecular Biology and Evolution.* 2013;30: 2531–2540.  
471 doi:10.1093/molbev/mst141
- 472 38. Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, et al. SILVA: a comprehensive  
473 online resource for quality checked and aligned ribosomal RNA sequence data compatible  
474 with ARB. *Nucleic Acids Research.* 2007;35: 7188–7196. doi:10.1093/nar/gkm864
- 475 39. Deshpande V, Wang Q, Greenfield P, Charleston M, Porras-Alfaro A, Kuske CR, et al. Fungal  
476 identification using a Bayesian classifier and the Warcup training set of internal  
477 transcribed spacer sequences. *Mycologia.* 2016;108: 1–5.
- 478 40. Abarenkov K, Nilsson RH, Larsson K-H, Alexander IJ, Eberhardt U, Erland S, et al. The UNITE  
479 database for molecular identification of fungi – recent updates and future perspectives.  
480 *New Phytologist.* 2010;186: 281–285.
- 481 41. Banchi E, Ametrano CG, Greco S, Stanković D, Muggia L, Pallavicini A. PLANITS: a curated  
482 sequence reference dataset for plant ITS DNA metabarcoding. *Database.* 2020;2020:  
483 baz155. doi:10.1093/database/baz155
- 484 42. Liu K-L, Porras-Alfaro A, Kuske CR, Eichorst SA, Xie G. Accurate, Rapid Taxonomic  
485 Classification of Fungal Large-Subunit rRNA Genes. *Appl Environ Microbiol.* 2012;78: 1523–  
486 1533. doi:10.1128/AEM.06826-11
- 487 43. St. John J. SeqPrep. Downloaded 2016. Available:  
488 <https://github.com/jstjohn/SeqPrep/releases>
- 489 44. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
490 *EMBnet journal.* 2011;17: pp-10.
- 491 45. Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon  
492 sequencing. *bioRxiv.* 2016 [cited 28 Jun 2018]. doi:10.1101/081257
- 493 46. Edgar R. UCHIME2: improved chimera prediction for amplicon sequencing. *bioRxiv.* 2016;  
494 074252.
- 495 47. Nearing JT, Douglas GM, Comeau AM, Langille MGI. Denoising the Denoisers: an  
496 independent evaluation of microbiome sequence error-correction approaches. *PeerJ.*  
497 2018;6: e5364. doi:10.7717/peerj.5364

- 498 48. He Y, Caporaso JG, Jiang X-T, Sheng H-F, Huse SM, Rideout JR, et al. Stability of operational  
499 taxonomic units: an important but neglected property for analyzing microbial diversity.  
500 *Microbiome*. 2015;3. doi:10.1186/s40168-015-0081-x
- 501 49. Antich A, Palacin C, Wangensteen OS, Turon X. To denoise or to cluster, that is not the  
502 question: optimizing pipelines for COI metabarcoding and metaphylogeography. *BMC*  
503 *Bioinformatics*. 2021;22: 177. doi:10.1186/s12859-021-04115-6
- 504 50. Robinson CV, Baird DJ, Wright MTG, Porter TM, Hartwig K, Hendriks E, et al. Combining  
505 DNA and people power for healthy rivers: Implementing the STREAM community-based  
506 approach for global freshwater monitoring. *Perspectives in Ecology and Conservation*.  
507 2021;19: 279–285. doi:10.1016/j.pecon.2021.03.001
- 508 51. Smenderovac E, Emilson C, Porter T, Morris D, Hazlett P, Diochon A, et al. Forest soil biotic  
509 communities show few responses to wood ash applications at multiple sites across  
510 Canada. *Sci Rep*. 2022;12: 4171. doi:10.1038/s41598-022-07670-x
- 511 52. Porter TM, Hajibabaei M. Scaling up: A guide to high-throughput genomic approaches for  
512 biodiversity analysis. *Molecular Ecology*. 2018;27: 313–338. doi:10.1111/mec.14478
- 513 53. Haase P, Pauls SU, Schindehütte K, Sundermann A. First audit of macroinvertebrate  
514 samples from an EU Water Framework Directive monitoring program: human error greatly  
515 lowers precision of assessment results. *Journal of the North American Benthological*  
516 *Society*. 2010;29: 1279–1291. doi:10.1899/09-183.1
- 517 54. Robinson CV, Porter TM, Maitland VC, Wright MT, Hajibabaei M. Multi-marker  
518 metabarcoding resolves subtle variations in freshwater condition: Bioindicators, ecological  
519 traits, and trophic interactions. *Ecology*; 2021 Nov. doi:10.1101/2021.11.14.468533
- 520 55. Poelen JH, Simons JD, Mungall CJ. Global biotic interactions: An open infrastructure to  
521 share and analyze species-interaction datasets. *Ecological Informatics*. 2014;24: 148–159.  
522 doi:10.1016/j.ecoinf.2014.08.005
- 523 56. Robinson CV, Porter TM, McGee KM, McCusker M, Wright MTG, Hajibabaei M. Multi-  
524 marker DNA metabarcoding detects suites of environmental gradients from an urban  
525 harbour. *BioRxiv*. 2022; 35.
- 526 57. Buchner D, Leese F. BOLDigger – a Python package to identify and organise sequences  
527 with the Barcode of Life Data systems. *MBMG*. 2020;4: e53535.  
528 doi:10.3897/mbmg.4.53535
- 529 58. Nugent CM, Elliott TA, Ratnasingham S, Adamowicz SJ. coil: an R package for cytochrome C  
530 oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation. *bioRxiv*.  
531 2019; 35. doi:doi: <http://dx.doi.org/10.1101/2019.12.12.865014>

532