

RESEARCH ARTICLE

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

Teresita M. Porter ^{*}, Mehrdad Hajibabaei

Centre for Biodiversity Genomics @ Biodiversity Institute of Ontario & Department of Integrative Biology, University of Guelph, Guelph, ON, Canada

^{*} terrimporter@gmail.com**OPEN ACCESS**

Citation: Porter TM, Hajibabaei M (2022) MetaWorks: A flexible, scalable bioinformatic pipeline for high-throughput multi-marker biodiversity assessments. *PLoS ONE* 17(9): e0274260. <https://doi.org/10.1371/journal.pone.0274260>

Editor: Hideyuki Doi, University of Hyogo, JAPAN

Received: April 22, 2022

Accepted: August 24, 2022

Published: September 29, 2022

Copyright: © 2022 Porter, Hajibabaei. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: MetaWorks runs at the command-line in a Conda environment on linux-64. Miniconda can be downloaded from https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh. If pseudogene filtering will be performed ORFfinder can be downloaded from the NCBI ftp site at <ftp://ftp.ncbi.nlm.nih.gov/genomes/TOOLS/ORFfinder/linux-i64/ORFfinder.gz>. MetaWorks can be downloaded from GitHub at <https://github.com/terrimporter/MetaWorks>. Quickstart examples, detailed pipeline descriptions, a tutorial for new users, and FAQs are available from the MetaWorks website at <https://>

Abstract

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

Introduction

Marker gene sequencing, metabarcoding, or metasytematics are different terms for the same technique that involves extracting DNA from bulk samples such as soil, water, or mixtures of individuals collected from traps. One key strength of this technique is not having to isolate or identify individual specimens. A signature DNA region is then enriched, for example using PCR, to identify biological community composition using bioinformatics [1–3]. In microbial ecology to animal

terrimporter.github.io/MetaWorksSite. Additional trained classifiers that can be used by MetaWorks for taxonomic assignment are available from GitHub at <https://github.com/terrimporter>.

Funding: MH received funding from Genome Canada and Ontario Genomics for the Sequencing the Rivers for Environmental Assessment and Monitoring (STREAM) project. TMP received funding from the Government of Canada through the Genomics Research and Development Initiative (GRDI), Metagenomics-based ecosystem biomonitoring (Ecobiomics) project. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

biodiversity studies, different signature DNA regions are chosen for their ability to identify target taxa. For example, in prokaryotes, the 16S small subunit (SSU) ribosomal RNA (rRNA) region is often used for genus level taxonomic assignments [4,5]. Other popular markers include cytochrome *c* oxidase (COI) for animals; ribulose biphosphate large subunit (rbcL) for plants and diatoms; the internal transcribed spacer (ITS) for fungi and plants; 18S SSU for eukaryotes, arbuscular mycorrhizal fungi, and diatoms; and 12S mitochondrial SSU for fish [6–14].

Existing pipelines such as QIIME2 and DADA2 were initially developed to support the microbial ecology community [15,16]. In comprehensive, multi-trophic, multi-marker studies, there is a need for a pipeline that can handle rRNA genes, spacer regions, as well as protein-coding markers in a single harmonized environment [17,18]. For the ITS region, we needed a pipeline that could remove the conserved flanking rRNA genes as this has been shown to improve taxonomic assignment accuracy [19]. For protein-coding regions, we needed a pipeline that could remove putative pseudogenes [20–23]. We also wanted the ability to generate high quality exact sequence variants (ESVs) for popular metabarcoding markers (not just 16S or ITS) for the additional level of genetic and taxonomic resolution ESVs can provide [24–26]. For taxonomic assignment, we wanted to use a classifier that would provide a measure of confidence for assignments to reduce false-positive assignments [27–29].

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

Implementation and workflow

Implementation

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

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

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

Marker	Target taxa	Classifier availability	Number of included sequences	Number of included taxa at all ranks (species)	Source data
COI	Eukaryotes	https://github.com/terrimporter/CO1Classifier	1,221,528	154,351 (114,687)	BOLD [34], INSDC [35]
rbcl	Diatoms	https://github.com/terrimporter/rbcLdiatomClassifier	3,504	1,432 (1,023)	Diat.barcode [36]
rbcl	Land plants	https://github.com/terrimporter/rbcL_landPlant_Classifier	148,258	61,398 (50,778)	INSDC [35]
rbcl	Eukaryotes	https://github.com/terrimporter/rbcLClassifier	164,454	65,742 (53,344)	INSDC [35]
12S	Fish	https://github.com/terrimporter/12SfishClassifier	2,853	4,751 (2,833)	MitoFish [37]
12S	Vertebrates	https://github.com/terrimporter/12SvertebrateClassifier	10,654	15,007 (9,564)	INSDC [35] and MitoFish [37]
SSU (18S)	Diatoms	https://github.com/terrimporter/SSUdiatomClassifier	2,962	1,198 (828)	Diat.barcode [36]
SSU (16S)	Vertebrates	https://github.com/terrimporter/16SvertebrateClassifier	72,195	21,282 (15,155)	INSDC [35]
SSU (18S)	Eukaryotes	https://github.com/terrimporter/18SClassifier	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)	https://github.com/terrimporter/UNITE_ITSClassifier	1,393,203	376,167 (352,588)	UNITE [40]
ITS	Plants	https://github.com/terrimporter/PLANIITS_ITSClassifier	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]

<https://doi.org/10.1371/journal.pone.0274260.t001>

Workflow

The pipeline begins with demultiplexed Illumina paired-end reads as this is the format most often provided by sequencing centres to their clients. Several workflows are available as Snake-make pipelines such as taxonomic assignment of ESVs (Fig 1), clustering of ESVs into OTUs, or for processing single reads. For each of these workflows described below, parameter settings for each bioinformatic step can be customized in the config.yaml file. The user also needs to provide a file of primer sequences so we provide a template for the adapters.fasta file as well as a small set of raw Illumina sequences for the COI amplicon that can be used to test the installation. The online documentation provides a tutorial example using the provided COI test data. The tutorial also walks users through the steps necessary to set up their environment to run the pipeline for the first time, assuming the user has never worked with Conda or Snakemake before.

Exact sequence variants

The ESV workflow will run the pipeline shown in Fig 1. In the config_ESV.yaml file, users indicate the path to the directory that contains the demultiplexed Illumina paired-end reads, specify the unique part of filenames to distinguish between samples and reads, and specify the name of the directory that will contain the outfiles. Default settings for each program are

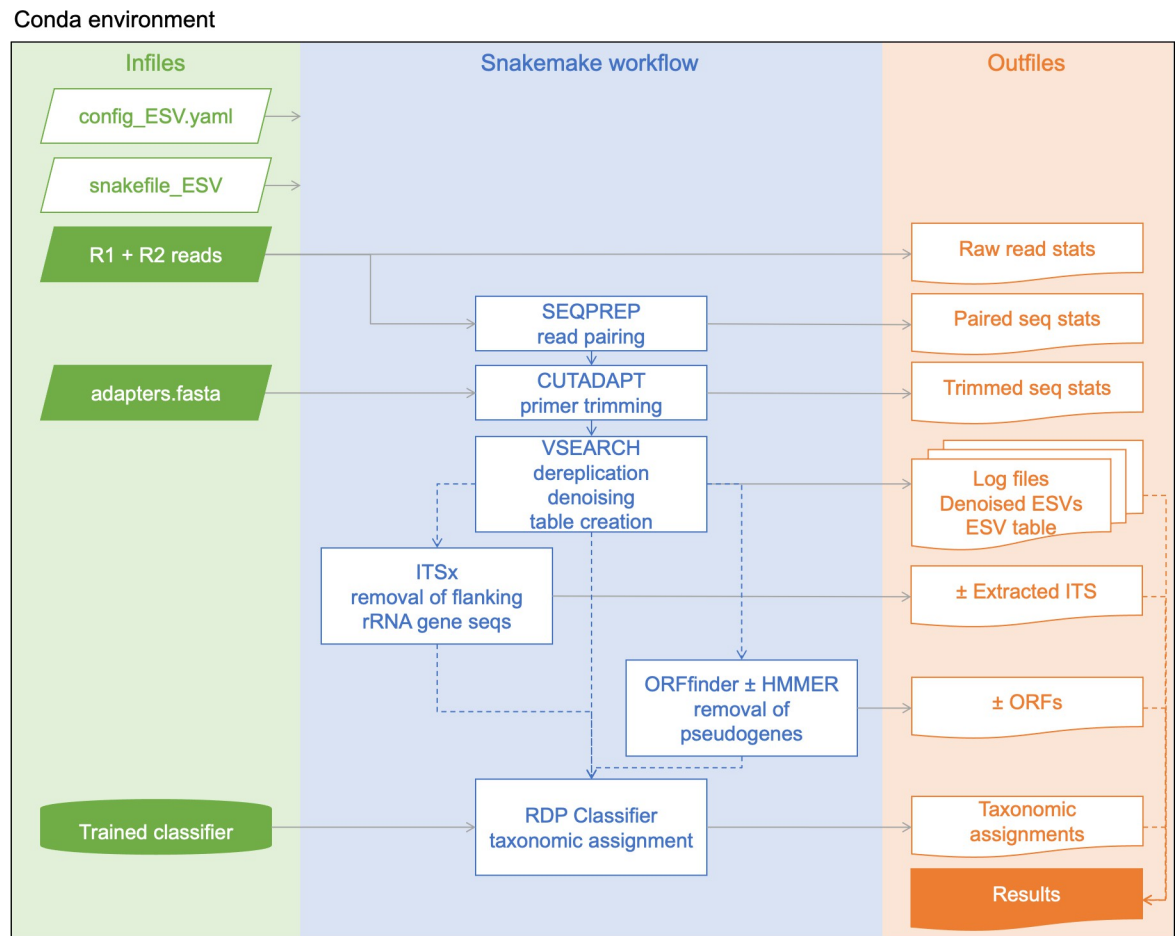


Fig 1. MetaWorks workflow to produce taxonomically assigned exact sequence variants. To aid reproducibility, a Conda environment is provided. Although multiple Snakemake workflows are provided in MetaWorks, here we show the main workflow that generates taxonomically assigned ESVs. Input files are shown in the first panel (green), the ESV workflow is shown in the centre panel (blue), and outfiles are shown in the last panel (orange). The input files in white boxes are required by snakemake to run the appropriate workflow. The input files in green need to be supplied by the user. Note that only custom-trained classifiers such as for COI need to be supplied by the user whereas classifiers built-in to the RDP classifier are used automatically to process prokaryote 16S assignments, for example. The denoising step shown here includes the removal of rare clusters, sequences with putative errors, as well as chimeric sequences. The results are provided in a comma-separated value (CSV) file and shows each ESV per sample with read counts and taxonomic assignments. Abbreviations: Demultiplexed Illumina paired-end reads (R1 + R2), internal transcribed spacer (ITS) region, open reading frame sequences (ORFs).

<https://doi.org/10.1371/journal.pone.0274260.g001>

provided in the `config_ESV.yaml` file but these can be customized by the user. SEQPREP was initially chosen for pairing the forward and reverse reads, because the program comes with the option to output the alignments for visual inspection, an option that most read-pairing programs do not have [43]. For SEQPREP read-pairing, users can specify a Phred score quality cutoff, the minimum overlap between the forward and reverse reads, the maximum fraction of mismatches allowed in the overlap region, and the minimum fraction of matching bases in the overlap region. CUTADAPT was chosen for primer-trimming because it is fast and already widely used in the metabarcoding community for this purpose, so most users will likely already be familiar with how this program works [44]. For CUTADAPT, users need to provide a FASTA-formatted primer sequence file (`adapters.fasta`), they can also specify the minimum sequence length to retain after primer-trimming, a Phred quality score cutoff, the maximum error rate, minimum adapter overlap, and maximum number of ambiguous bases allowed.

VSEARCH was chosen to dereplicate reads (retain unique reads) and remove artefactual sequences using the UNOISE3 and UCHIME3 algorithms [45,46]. We chose the open-source VSEARCH program over alternatives because the program can utilize all the available memory on a system, facilitating the analysis of large datasets on high performance computer systems. We prefer the UNOISE3 method for denoising because it performs up to 1,200 faster and uses less memory than other denoising programs [47]. To map read counts to the newly generated denoised-chimera ESVs to create an ESV x sample table, we use the 'search_exact' method because it is faster and optimized to find exact matches compared with the 'usearch_global' command with the 'id 1.0' parameter, but this is just an intermediate step and further filtering of this table is performed by MetaWorks.

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

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

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

Operational taxonomic units

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

After processing raw reads using the `snakefile_ESV` workflow described in the previous section, users can use the `snakefile_OTU` workflow to cluster ESVs into OTUs. This approach combines the benefits of denoising with clustering using a 97% sequence similarity cutoff using the `snakefile_OTU` workflow [26,49]. This method uses VSEARCH 'cluster_smallmem' method to cluster ESVs using a 97% sequence similarity cutoff. Settings can be adjusted in the `config_OTU.yaml` file such as pointing to the directory that contains the ESVs and choosing a classifier for the OTUs.

Results and discussion

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

Use case 1: As a part of the Canadian STREAM biomonitoring initiative, the MetaWorks pipeline has been used to process macroinvertebrate COI metabarcodes surveyed from stream sites across Canada [50]. One feature of this project is the quick 1–2 month turn-around time from sampling through to the production of watershed biodiversity reports. This is an improvement over reports generated using conventional morphology-based methods that would normally take 6–12 months to produce. The use of a consistent bioinformatics workflow to process metabarcodes has played a key role in the reproducibility, scalability, and throughput to facilitate timely reporting [52]. Generally, samples are processed in batches of 96 per sequencing run then later split into custom reports for stakeholders, processing about 500 samples per year. One feature of these reports are the taxonomic assignments made using the naive Bayesian classifier that provides bootstrap support values. During the data analysis stage, users can use minimum bootstrap support cutoffs to ensure a certain level of expected accuracy (80–99%) and reduce false-positive taxonomic assignments [28]. The cutoffs used are specific to the amplicon, amplicon length, and taxonomic rank of the assignment and assumes the query is represented in the underlying sequence database. This is in contrast with the use of more traditional methods for taxonomic assignment, where taxa are routinely missed during subsampling and taxa detected by primary analysts and auditors may differ by up to 30% [53]. This use-case shows how MetaWorks can be used to create taxon lists for large-scale biodiversity monitoring of streams across Canada.

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

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

As demonstrated in the above examples, MetaWorks supports a wide range of analysis scenarios from metabarcoding data. We envision that MetaWorks will aid broader user communities and fill a need in multi-marker metabarcoding studies that target taxa from multiple different domains of life, to provide a unified processing environment, pipeline, and taxonomic assignment approach for each marker from ribosomal RNA genes, spacers, or protein coding genes. QIIME2 is perhaps the most popular and comprehensive platform for such work, but to date, focuses on processing mainly prokaryote and fungal datasets [16]. To our knowledge, MetaWorks is the only bioinformatic pipeline that can handle rRNA genes but that also integrates special processing steps to handle ITS spacers as well as filter out obvious pseudogenes in protein coding markers such as COI.

There has been a lot of activity with respect to building new bioinformatic tools to handle COI metabarcodes. Recent work, such as the BOLDigger program, makes the BOLD identification engine more suitable for identifying large batches of COI metabarcodes and has both GUI and command-line interfaces for efficient sample processing [57]. A new program, called NUMTdumper, has been developed as a stand-alone program meant to be incorporated into bioinformatic pipelines [20]. NUMTdumper provides a method to screen for NuMTs based on read counts while acknowledging the trade-offs between removing all possible NuMTs while erroneously removing genuine reads. An R package called 'coil' has also recently been developed that will place COI barcode and metabarcode sequences in frame using profile

HMM analysis [58]. MetaWorks aims to extend the COI metabarcoding toolkit that provides a harmonized environment where data from other organismal markers in multi-marker, multi-trophic studies can also be analyzed.

Conclusion

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

Acknowledgments

We would like to thank Josip Rudar, Katie M. McGee, Chloe V. Robinson, Victoria C. Maitland, and Michael T.G. Wright from the Hajibabaei lab for helpful discussions and testing the pipeline with various datasets. We also thank Artin Mashayekhi for assistance with the MetaWorks website.

Author Contributions

Conceptualization: Teresita M. Porter, Mehrdad Hajibabaei.

Formal analysis: Teresita M. Porter.

Funding acquisition: Mehrdad Hajibabaei.

Methodology: Teresita M. Porter.

Resources: Mehrdad Hajibabaei.

Software: Teresita M. Porter.

Supervision: Mehrdad Hajibabaei.

Validation: Teresita M. Porter.

Visualization: Teresita M. Porter.

Writing – original draft: Teresita M. Porter.

Writing – review & editing: Teresita M. Porter, Mehrdad Hajibabaei.

References

1. Pace NR. A Molecular View of Microbial Diversity and the Biosphere. *Science*. 1997; 276: 734–740. <https://doi.org/10.1126/science.276.5313.734> PMID: 9115194

2. Hajibabaei M. The golden age of DNA metasystematics. *Trends in genetics*. 2012; 28: 535–537. <https://doi.org/10.1016/j.tig.2012.08.001> PMID: 22951138
3. Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E. Towards next-generation biodiversity assessment using DNA metabarcoding. *Molecular ecology*. 2012; 21: 2045–2050. <https://doi.org/10.1111/j.1365-294X.2012.05470.x> PMID: 22486824
4. Gilbert JA, Jansson JK, Knight R. The Earth Microbiome project: successes and aspirations. *BMC biology*. 2014; 12: 69. <https://doi.org/10.1186/s12915-014-0069-1> PMID: 25184604
5. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Applied and Environmental Microbiology*. 2007; 73: 5261–5267. <https://doi.org/10.1128/AEM.00062-07> PMID: 17586664
6. Schüßler A. Glomales SSUrRNA gene diversity. *New Phytologist*. 1999; 144: 205–207. <https://doi.org/10.1046/j.1469-8137.1999.00526.x>
7. Hebert PDN, Cywinska A, Ball SL, deWaard JR. Biological identifications through DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences*. 2003; 270: 313–321. <https://doi.org/10.1098/rspb.2002.2218> PMID: 12614582
8. James TY, Kauff F, Schoch CL, Matheny PB, Hofstetter V, Cox CJ, et al. Reconstructing the early evolution of Fungi using a six-gene phylogeny. *Nature*. 2006; 443: 818–822. <https://doi.org/10.1038/nature05110> PMID: 17051209
9. Hibbett DS, Binder M, Bischoff JF, Blackwell M, Cannon PF, Eriksson OE, et al. A higher-level phylogenetic classification of the Fungi. *Mycological Research*. 2007; 111: 509–547. <https://doi.org/10.1016/j.mycres.2007.03.004> PMID: 17572334
10. CBOL Plant Working Group, Hollingsworth PM, Forrest LL, Spouge JL, Hajibabaei M, Ratnasingham S, et al. A DNA barcode for land plants. *Proceedings of the National Academy of Sciences*. 2009; 106: 12794–12797. <https://doi.org/10.1073/pnas.0905845106> PMID: 19666622
11. Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA, et al. Nuclear ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi. *Proceedings of the National Academy of Sciences*. 2012; 109: 6241–6246. <https://doi.org/10.1073/pnas.1117018109> PMID: 22454494
12. Zimmermann J, Abarca N, Enk N, Skibbe O, Kusber W-H, Jahn R. Taxonomic Reference Libraries for Environmental Barcoding: A Best Practice Example from Diatom Research. Schierwater B, editor. *PLoS ONE*. 2014; 9: e108793. <https://doi.org/10.1371/journal.pone.0108793> PMID: 25265556
13. Sato Y, Miya M, Fukunaga T, Sado T, Iwasaki W. MitoFish and MiFish Pipeline: A Mitochondrial Genome Database of Fish with an Analysis Pipeline for Environmental DNA Metabarcoding. Kumar S, editor. *Molecular Biology and Evolution*. 2018; 35: 1553–1555. <https://doi.org/10.1093/molbev/msy074> PMID: 29668970
14. Ahmed M, Back MA, Prior T, Karssen G, Lawson R, Adams I, et al. Metabarcoding of soil nematodes: the importance of taxonomic coverage and availability of reference sequences in choosing suitable marker(s). *MBMG*. 2019; 3: e36408. <https://doi.org/10.3897/mbmg.3.36408>
15. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution sample inference from Illumina amplicon data. *Nature Methods*. 2016; 13: 581–583. <https://doi.org/10.1038/nmeth.3869> PMID: 27214047
16. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet C, Al-Ghalith GA, et al. Reproducible, interactive, scalable, and extensible microbiome data science using QIIME 2. *Nature Biotechnology*. 2019; 37: 852–857. <https://doi.org/10.1038/s41587-019-0209-9> PMID: 31341288
17. Drummond AJ, Newcomb RD, Buckley TR, Xie D, Dopheide A, Potter BC, et al. Evaluating a multigene environmental DNA approach for biodiversity assessment. *GigaSci*. 2015; 4: 46. <https://doi.org/10.1186/s13742-015-0086-1> PMID: 26445670
18. Edge TA, Baird DJ, Bilodeau G, Gagné N, Greer C, Konkin D, et al. The Ecobiomics project: Advancing metagenomics assessment of soil health and freshwater quality in Canada. *Science of The Total Environment*. 2020; 710: 135906. <https://doi.org/10.1016/j.scitotenv.2019.135906> PMID: 31926407
19. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, et al. Improved software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and other eukaryotes for analysis of environmental sequencing data. Bunce M, editor. *Methods in Ecology and Evolution*. 2013; 4: 914–919. <https://doi.org/10.1111/2041-210X.12073>
20. Andújar C, Creedy TJ, Arribas P, López H, Salces-Castellano A, Pérez-Delgado A, et al. NUMT dumping: validated removal of nuclear pseudogenes from mitochondrial metabarcoding data. *Evolutionary Biology*; 2020 Jun. <https://doi.org/10.1101/2020.06.17.157347>

21. Porter TM, Hajibabaei M. Profile hidden Markov model sequence analysis can help remove putative pseudogenes from DNA barcoding and metabarcoding datasets. *BMC Bioinformatics*. 2021; 22: 256. <https://doi.org/10.1186/s12859-021-04180-x> PMID: 34011275
22. Song H, Buhay JE, Whiting MF, Crandall KA. Many species in one: DNA barcoding overestimates the number of species when nuclear mitochondrial pseudogenes are coamplified. *PNAS*. 2008; 105: 13486–13491. <https://doi.org/10.1073/pnas.0803076105> PMID: 18757756
23. Moulton MJ, Song H, Whiting MF. Assessing the effects of primer specificity on eliminating numt coamplification in DNA barcoding: a case study from Orthoptera (Arthropoda: Insecta): DNA BARCODING. *Molecular Ecology Resources*. 2010; 10: 615–627. <https://doi.org/10.1111/j.1755-0998.2009.02823.x> PMID: 21565066
24. Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational taxonomic units in marker-gene data analysis. *The ISME Journal*. 2017; 11: 2639–2643. <https://doi.org/10.1038/ismej.2017.119> PMID: 28731476
25. Glassman SI, Martiny JB. Ecological patterns are robust to use of exact sequence variants versus operational taxonomic units. *mSphere*. 2018; 3: e00148–18. <https://doi.org/10.1101/283283>
26. Porter TM, Hajibabaei M. Putting COI Metabarcoding in Context: The Utility of Exact Sequence Variants (ESVs) in Biodiversity Analysis. *Front Ecol Evol*. 2020; 8: 248. <https://doi.org/10.3389/fevo.2020.00248>
27. Porter TM, Gibson JF, Shokralla S, Baird DJ, Golding GB, Hajibabaei M. Rapid and accurate taxonomic classification of insect (class Insecta) cytochrome c oxidase subunit 1 (COI) DNA barcode sequences using a naïve Bayesian classifier. *Mol Ecol Resour*. 2014; 14: 929–942. <https://doi.org/10.1111/1755-0998.12240>
28. Porter TM, Hajibabaei M. Automated high throughput animal CO1 metabarcoding classification. *Scientific Reports*. 2018; 8: 4226. <https://doi.org/10.1038/s41598-018-22505-4> PMID: 29523803
29. Virgilio M, Bäckeljau T, Nevado B, De Meyer M. Comparative performances of DNA barcoding across insect orders. *BMC bioinformatics*. 2010; 11: 206. <https://doi.org/10.1186/1471-2105-11-206> PMID: 20420717
30. Adamowicz SJ, Boatwright JS, Chain F, Fisher BL, Hogg ID, Leese F, et al. Trends in DNA barcoding and metabarcoding. *Genome*. 2019; 62: v–viii. <https://doi.org/10.1139/gen-2019-0054> PMID: 30998119
31. Koster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics*. 2012; 28: 2520–2522. <https://doi.org/10.1093/bioinformatics/bts480> PMID: 22908215
32. Anaconda. Anaconda Software Distribution. 2016. Available: <https://anaconda.com>.
33. Wratten L, Wilm A, Göke J. Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers. *Nat Methods*. 2021 [cited 24 Sep 2021]. <https://doi.org/10.1038/s41592-021-01254-9> PMID: 34556866
34. Ratnasingham S, Hebert PD. BOLD: The Barcode of Life Data System (<http://www.barcodinglife.org>). *Molecular ecology notes*. 2007; 7: 355–364.
35. Cochrane G, Karsch-Mizrachi I, Takagi T, Sequence Database Collaboration IN. The International Nucleotide Sequence Database Collaboration. *Nucleic Acids Research*. 2016; 44: D48–D50. <https://doi.org/10.1093/nar/gkv1323> PMID: 26657633
36. Rimet F, Gusev E, Kahlert M, Kelly MG, Kulikovskiy M, Maltsev Y, et al. Diat.barcode, an open-access curated barcode library for diatoms. *Sci Rep*. 2019; 9: 15116. <https://doi.org/10.1038/s41598-019-51500-6> PMID: 31641158
37. Iwasaki W, Fukunaga T, Isagozawa R, Yamada K, Maeda Y, Satoh TP, et al. MitoFish and MitoAnnotator: A Mitochondrial Genome Database of Fish with an Accurate and Automatic Annotation Pipeline. *Molecular Biology and Evolution*. 2013; 30: 2531–2540. <https://doi.org/10.1093/molbev/mst141> PMID: 23955518
38. Priesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, et al. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research*. 2007; 35: 7188–7196. <https://doi.org/10.1093/nar/gkm864> PMID: 17947321
39. Deshpande V, Wang Q, Greenfield P, Charleston M, Porras-Alfaro A, Kuske CR, et al. Fungal identification using a Bayesian classifier and the Warcup training set of internal transcribed spacer sequences. *Mycologia*. 2016; 108: 1–5. <https://doi.org/10.3852/14-293> PMID: 26553774
40. Abarenkov K, Nilsson RH, Larsson K-H, Alexander IJ, Eberhardt U, Erland S, et al. The UNITE database for molecular identification of fungi—recent updates and future perspectives. *New Phytologist*. 2010; 186: 281–285. <https://doi.org/10.1111/j.1469-8137.2009.03160.x> PMID: 20409185
41. Banchi E, Ametrano CG, Greco S, Stanković D, Muggia L, Pallavicini A. PLANITS: a curated sequence reference dataset for plant ITS DNA metabarcoding. *Database*. 2020; 2020: baz155. <https://doi.org/10.1093/database/baz155> PMID: 32016319

42. Liu K-L, Porras-Alfaro A, Kuske CR, Eichorst SA, Xie G. Accurate, Rapid Taxonomic Classification of Fungal Large-Subunit rRNA Genes. *Appl Environ Microbiol*. 2012; 78: 1523–1533. <https://doi.org/10.1128/AEM.06826-11> PMID: 22194300
43. St. John J. SeqPrep. Downloaded 2016. Available: <https://github.com/jstjohn/SeqPrep/releases>.
44. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal*. 2011; 17: pp-10.
45. Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing. *bioRxiv*. 2016 [cited 28 Jun 2018]. <https://doi.org/10.1101/081257>
46. Edgar R. UCHIME2: improved chimera prediction for amplicon sequencing. *bioRxiv*. 2016; 074252.
47. Nearing JT, Douglas GM, Comeau AM, Langille MG. Denoising the Denoisers: an independent evaluation of microbiome sequence error-correction approaches. *PeerJ*. 2018; 6: e5364. <https://doi.org/10.7717/peerj.5364> PMID: 30123705
48. He Y, Caporaso JG, Jiang X-T, Sheng H-F, Huse SM, Rideout JR, et al. Stability of operational taxonomic units: an important but neglected property for analyzing microbial diversity. *Microbiome*. 2015;3. <https://doi.org/10.1186/s40168-015-0081-x> PMID: 25995836
49. Antich A, Palacin C, Wangenstein OS, Turon X. To denoise or to cluster, that is not the question: optimizing pipelines for COI metabarcoding and metaphylogeography. *BMC Bioinformatics*. 2021; 22: 177. <https://doi.org/10.1186/s12859-021-04115-6> PMID: 33820526
50. Robinson CV, Baird DJ, Wright MTG, Porter TM, Hartwig K, Hendriks E, et al. Combining DNA and people power for healthy rivers: Implementing the STREAM community-based approach for global freshwater monitoring. *Perspectives in Ecology and Conservation*. 2021; 19: 279–285. <https://doi.org/10.1016/j.pecon.2021.03.001>
51. Smenderovac E, Emilson C, Porter T, Morris D, Hazlett P, Diochon A, et al. Forest soil biotic communities show few responses to wood ash applications at multiple sites across Canada. *Sci Rep*. 2022; 12: 4171. <https://doi.org/10.1038/s41598-022-07670-x> PMID: 35264620
52. Porter TM, Hajibabaei M. Scaling up: A guide to high-throughput genomic approaches for biodiversity analysis. *Molecular Ecology*. 2018; 27: 313–338. <https://doi.org/10.1111/mec.14478> PMID: 29292539
53. Haase P, Pauls SU, Schindehütte K, Sundermann A. First audit of macroinvertebrate samples from an EU Water Framework Directive monitoring program: human error greatly lowers precision of assessment results. *Journal of the North American Benthological Society*. 2010; 29: 1279–1291. <https://doi.org/10.1899/09-183.1>
54. Robinson CV, Porter TM, Maitland VC, Wright MT, Hajibabaei M. Multi-marker metabarcoding resolves subtle variations in freshwater condition: Bioindicators, ecological traits, and trophic interactions. *bioRxiv*. 2021 [cited 16 Nov 2021]. <https://doi.org/10.1101/2021.11.14.468533>
55. Poelen JH, Simons JD, Mungall CJ. Global biotic interactions: An open infrastructure to share and analyze species-interaction datasets. *Ecological Informatics*. 2014; 24: 148–159. <https://doi.org/10.1016/j.ecoinf.2014.08.005>
56. Robinson CV, Porter TM, McGee KM, McCusker M, Wright MTG, Hajibabaei M. Multi-marker DNA metabarcoding detects suites of environmental gradients from an urban harbour. *Sci Rep*. 2022; 12: 10556. <https://doi.org/10.1038/s41598-022-13262-6> PMID: 35732669
57. Buchner D, Leese F. BOLDigger—a Python package to identify and organise sequences with the Barcode of Life Data systems. *MBMG*. 2020; 4: e53535. <https://doi.org/10.3897/mbmg.4.53535>
58. Nugent CM, Elliott TA, Ratnasingham S, Adamowicz SJ. coil: an R package for cytochrome c oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation. *Genome*. 2020; 63: 291–305. <https://doi.org/10.1139/gen-2019-0206> PMID: 32406757