

AMOEBAE documentation

Lael D. Barlow

Version of July 14, 2020

Contents

1	Introduction	1
1.1	What is AMOEBAE?	1
1.2	Why use AMOEBAE?	1
1.3	Key features	2
1.4	A word of caution	2
1.5	User support	2
1.6	Documentation	2
1.7	How to cite AMOEBAE	3
1.8	Acknowledgments	3
1.9	License	3
2	How to start using AMOEBAE	4
2.1	System requirements	4
2.2	Dependencies	4
2.3	Setting up an environment for AMOEBAE using Singularity	4
2.4	Running AMOEBAE using Jupyter notebooks	5
2.5	Running AMOEBAE via the command line	6
3	Command reference	7
3.1	amoebae	8
3.2	amoebae mkdatadir	9
3.3	amoebae setup_hmmdb	10
3.4	amoebae add_to_dbs	10
3.5	amoebae list_dbs	10
3.6	amoebae add_to_queries	11

3.7	amoebae list_queries	11
3.8	amoebae get_redun_hits	12
3.9	amoebae setup_fwd_srch	13
3.10	amoebae run_fwd_srch	14
3.11	amoebae sum_fwd_srch	14
3.12	amoebae setup_rev_srch	16
3.13	amoebae run_rev_srch	17
3.14	amoebae sum_rev_srch	17
3.15	amoebae interp_srchs	18
3.16	amoebae find_redun_seqs	19
3.17	amoebae plot	23
3.18	amoebae add_to_models	23
3.19	amoebae list_models	24
3.20	amoebae get_alt_topos	24
3.21	amoebae prune	25
3.22	amoebae auto_prune	25
3.23	amoebae reduce_tree	26
3.24	amoebae constrain_mb	26
3.25	amoebae visualize_tree	27
3.26	amoebae replace_seqs	27
3.27	amoebae csv_to_fasta	28
3.28	amoebae check_depend	29
3.29	amoebae check_imports	29
3.30	amoebae regen_genome_info	29

4 Miscellaneous scripts 30

1 Introduction

1.1 What is AMOEBAE?

Analysis of MOlecular Evolution with BAth Entry (AMOEBAE) is a bioinformatics software toolkit composed primarily of scripts written in the Python3 language. AMOEBAE scripts use existing Python packages including Biopython (Cock *et al.*, 2009), the Environment for Tree Exploration (ETE3) (Huerta-Cepas *et al.*, 2016), Pandas, and Matplotlib (Hunter, 2007) for setting up, running, and summarizing analyses of molecular evolution using bioinformatics software packages including MUSCLE (Edgar, 2004), BLAST+ (Camacho *et al.*, 2009), HMMer3 (Eddy, 1998), and IQ-TREE (Nguyen *et al.*, 2015). Applications include identifying and classifying predicted peptide sequences according to their evolutionary relationships with homologues. All dependencies are freely available, and AMOEBAE code is open-source (see subsection 1.9) and available on GitHub (<https://github.com/laelbarlow/amoebae>).

1.2 Why use AMOEBAE?

The general problem that AMOEBAE addresses is as follows. Numerous genomes (and transcriptomes) are available for a wide diversity of species of medical, economic, and ecological importance. Yet only a small minority of these species are tractable models for genetic and cell biological experimentation. Effective translation of genetic and cell biological knowledge from model organisms to non-model organisms with sequenced genomes is thus essential to maximize return on investment in scientific research. This translation begins with comparative genomics analyses which aim to compare genes in non-model organisms to characterized genes in model organisms, within the over-arching framework of evolutionary theory. Efficient methods are required to perform such analyses, yet some such methods may not be suited to the scope of particular studies due to their breadth and/or depth.

AMOEBAE is useful for certain mid-scale comparative genomics studies that might otherwise require a much larger investment of repetitive manual/visual manipulation of data. Webservices such as those provided by NCBI (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) (Camacho *et al.*, 2009) and EMBL-EBI (<https://www.ebi.ac.uk/Tools/hmmer/>) provide a means to readily investigate the evolution of one or a few genes via similarity searching, and large-scale analysis pipelines such as OrthoMCL (Li, 2003) and OrthoFinder (Emms and Kelly, 2019) attempt to rapidly perform orthology prediction for all genes among several genomes. AMOEBAE addresses mid-scale analyses which are too cumbersome to be performed via webservices or simple scripts and yet require a level of detail and flexibility not offered by large-scale analysis pipelines. AMOEBAE may be useful for analyzing the distribution of orthologues of up to perhaps 30 genes/proteins among a sampling of no more than approximately 50 eukaryotic genomes. AMOEBAE provides many options which can be tailored to the specific genes/proteins being analyzed, and allow analyses using complex sets of customized criteria to be reproduced more practically.

1.3 Key features

The core functionality of AMOEBAE is to run sequence similarity searches with multiple algorithms, multiple queries, and multiple databases simultaneously and to allow highly customizable implementation of reciprocal-best-hit search strategies. The output includes detailed summaries of results in the form of a spreadsheet and plots.

A particular advantage of AMOEBAE over other tools is the functionality for parsing results of TBLASTN (which searches nucleotide sequences with peptide sequence queries) search results. This allows rapid identification of High-scoring Segment Pair (HSP) clusters at separate gene loci (identified according to user-defined criteria), automatic checking of those loci against information in genome annotation files, and systematic use of Exonerate (Slater and Birney, 2005) where possible for obtaining better exon predictions.

1.4 A word of caution

AMOEBAE is not optimized for ease of use, but is meant to be highly configurable. The many options available to AMOEBAE users inevitably provide many opportunities for user errors in specifying search criteria, and user errors in interpreting results detailed in output files. Some prior experience with similarity searching and with running software using the command line are prerequisites for using AMOEBAE, and experience writing scripts in Bash (linux shell) and Python would be highly advantageous. Also, you may need to carefully define the scope of your analysis depending on what additional steps you may find necessary beyond those that may be performed using AMOEBAE (you may find that the maximum 30 queries and 50 genomes suggested above may in fact be unmanageable). Moreover, AMOEBAE is still under active development, so some features may not yet be thoroughly tested.

1.5 User support

For specific issues with the code, please use the issue tracker on the GitHub webpage here: <https://github.com/laelbarlow/amoebae/issues>.

If you have general questions regarding AMOEBAE, please email the author at lael (at) ualberta (dot) ca.

1.6 Documentation

This document provides an overview of AMOEBAE and describes the functionality of the various commands/scripts. For a tutorial which includes a working example of a similarity search analysis run using AMOEBAE, see the Jupyter Notebook: amoebae/notebooks/similarity_search_tutorial_2.ipynb. For code documentation, please see the html file(s), which

1 can be opened with your web browser: `amoebae/documentation/code_documentation/`
2 `html/index.html`.

3 1.7 How to cite AMOEBAE

4 Please cite the GitHub webpage <https://github.com/laelbarlow/amoebae> (or alternative
5 permanent repositories if relevant). Also, the first publication to make use of a version of
6 AMOEBAE was an analysis of Adaptor Protein subunits in embryophytes by Larson *et al.*
7 (2019).

8 Also, you may wish to cite the software packages which are key dependencies of AMOEBAE,
9 since AMOEBAE would not work without these (see subsection 2.2).

10 1.8 Acknowledgments

11 AMOEBAE was initially developed in the Dacks Laboratory at the University of Alberta, and
12 was supported by National Sciences and Engineering Council of Canada (NSERC) Discovery
13 grants RES0021028, RES0043758, and RES0046091 awarded to Joel B. Dacks, as well as an
14 NSERC Postgraduate Scholarship-Doctoral awarded to Lael D. Barlow.

15 We acknowledge the support of the Natural Sciences and Engineering Research Council of
16 Canada (NSERC).

17 Cette recherche a été financée par le Conseil de recherches en sciences naturelles et en génie
18 du Canada (CRSNG).

19 Also, help with testing AMOEBAE has been kindly provided by Raegan T. Larson, Shweta
20 V. Pipalya, Kira More, and Kristína Záhonová.

21 1.9 License

22 Copyright 2018 Lael D. Barlow

23 Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file
24 except in compliance with the License. You may obtain a copy of the License at

25 <http://www.apache.org/licenses/LICENSE-2.0>

26 Unless required by applicable law or agreed to in writing, software distributed under the
27 License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS
28 OF ANY KIND, either express or implied. See the License for the specific language governing
29 permissions and limitations under the License.

2 How to start using AMOEBAE

2.1 System requirements

Please note that the commands shown likely only work on MacOS or Linux operating systems (you may have trouble running AMOEBAE directly on Windows).

2.2 Dependencies

You do not need to install these dependencies yourself. All dependencies are free and open-source, and are automatically installed in a virtual environment for AMOEBAE scripts (see subsection 2.3).

The main dependencies of AMOEBAE include the following:

- Python3.
- Biopython, a Python package for bioinformatics (Cock *et al.*, 2009).
- The Environment for Tree Exploration 3 (ETE3), a Python package for working with phylogenetic trees (Huerta-Cepas *et al.*, 2016).
- Matplotlib, a Python package for generating plots (Hunter, 2007).
- (gffutils).
- NCBI BLAST+, a software package for sequence similarity searching (Camacho *et al.*, 2009).
- HMMer3, a software package for profile sequence similarity searching (Eddy, 1998).
- MUSCLE, for multiple sequence alignment (Edgar, 2004).
- IQ-TREE, for phylogenetic analysis (Nguyen *et al.*, 2015).

2.3 Setting up an environment for AMOEBAE using Singularity

Follow the steps below to set up AMOEBAE on your personal computer, or on a Linux cluster with Singularity (<https://sylabs.io/singularity/>) pre-installed. This setup process should take approximately 5 minutes to complete.

1. If you are setting up AMOEBAE on a high performance computing cluster, then you will probably not be able to install Singularity yourself, or may need to use specific

procedures to load Singularity prior to use. Contact your system administrator(s) if Singularity is not installed, and direct them to this webpage: <https://sylabs.io/guides/3.5/admin-guide/>.

2. If you are setting up AMOEBAE on a personal computer, ensure that you have at least 30GB of empty storage space available (and keep in mind that it is generally recommended that you leave at least 20% of your storage space empty for efficient performance). This is important for running virtual machines.
3. If using a personal computer, ensure that Git is installed on your computer. If you do not already have git installed, then your computer will prompt you to install it when you type git into the command line. If you are using MacOS, the easiest way to install Git is by installing Xcode via the App Store (this will use up a considerable amount of storage space). Documentation for Git is available here: <https://git-scm.com/doc>. You can check which version you have (or whether it is installed at all) by running the command below. Please note: Here ">>>" is used to indicate that the following text in the line is to be entered in you terminal command prompt.

```
>>> git --version
```

4. Clone the AMOEBAE repository using Git. If you simply download the code from GitHub, instead of cloning the repository, then AMOEBAE cannot record specifically what version of the code you use, and will not run properly. Make sure to use the appropriate directory path (the path shown is just an example). Also, replace the path shown below with the path to the directory on your system where you wish to put the main AMOEBAE directory.

```
>>> cd /path/to/directory/where/you/keep/files
>>> git clone https://github.com/laelbarlow/amoebae.git
```

5. Set up AMOEBAE. This performs several steps including checking for whether singularity is installed and attempting to use VirtualBox and Vagrant to run Singularity in a pre-built Ubuntu virtual machine with Singularity installed. This is because Singularity does not run on MacOS (or Windows), and installation of Singularity on Linux is complex, as several dependencies are required. This script downloads a pre-built singularity container, which was built using the singularity.recipe file, and provided on the Singularity Library (https://cloud.sylabs.io/library/_container/5e8ca8fff0f8eb90a8a7b60d).

```
>>> cd amoebae
>>> bash setup.sh
```

2.4 Running AMOEBAE using Jupyter notebooks

1. After setting up AMOEBAE according to the instructions above, the easiest way to start running analyses using AMOEBAE is via the tutorials, which are in the form of Jupyter notebooks (<https://jupyter.org/>). These Jupyter notebooks can be run

using the installation of Jupyter in the Singularity container, and can be accessed using your browser (on a personal computer). To start a Jupyter server, run the bash script as indicated below (assuming your current working directory is the main amoebae directory that you cloned with Git).

```
>>> bash singularity_jupyter.sh
```

2. Copy and past the resulting URL (the one at the bottom of the output) into the address bar of your web browser (either Firefox, Chrome, or Safari will work). This will open Jupyter to the notebooks subdirectory, which contains several tutorial and example notebooks (.ipynb files). These files are the files on your regular (host) filesystem, as the amoebae directory is synced with the Singularity container. Thus changes to files will persist after you shut down the Jupyter server and the Singularity container. Documentation on Jupyter is available here: <https://jupyter-notebook.readthedocs.io/en/stable/>.
3. Click on one of the tutorial files (.ipynb). These Jupyter notebooks include information on how to use them once opened. The first tutorial (amoebae_tutorial_1.ipynb) provides a simple example of similarity searching with BLASTP using a Jupyter notebook. The second tutorial (amoebae_tutorial_2.ipynb) provides an example using most of the similarity searching functionality that AMOEBAE provides.
4. To shut down the Jupyter server, click the logout button in the jupyter browser tab(s), and then go to the terminal window that you used to startup the Jupyter server, and press CTRL-C to kill the Jupyter kernel. This will close the Jupyter notebooks, but the analysis output files will remain, because they are saved to your amoebae/notebooks folder which is on your host machine and accessed from within the container.
5. Working with the Jupyter notebooks interactively in this manner on high-performance computing clusters is likely possible but inconvenient, and procedures will vary. Also, running the tutorial notebooks would require access to the internet from compute nodes (as opposed to login nodes) which may not be supported. Therefore, it is recommended that you run the tutorials on a personal laptop/desktop computer if possible. To run your own notebooks on a cluster, you will need to write a job submission script that will be specific to the cluster, the job scheduler it uses, and your account details. Please refer to documentation provided by your system administrators for this. For an example script that writes a script for running a notebook as a job to a SLURM job scheduler see https://github.com/laelbarlow/amoebae/blob/master/notebooks/write_notebook_slurm_script.sh.

2.5 Running AMOEBAE via the command line

1. The easiest way to access AMOEBAE dependencies via the command line is to use the bash script provided. If you are running AMOEBAE on a personal computer (running singularity in a virtual machine), then, without customizing the code, only one shell session may be opened at once (and these cannot be opened at the same time as the

singularity_jupyter.sh script is running Singularity in a virtual machine). Running the script as indicated below will open a shell session in the Singularity container, with the amoebae directory being the only one accessible. Also, the amoebae executable script is added to the \$PATH in the container, so you can run amoebae commands from any directory.

```
>>> bash singularity_shell.sh
```

2. You may find it useful to explore and test the environment using the following commands.

- Print the paths included in the \$PATH variable in the container.

```
>>> tr ':' '\n' <<< "$PATH"
```

- Check the location of the amoebae executable being run from within the container.

```
>>> command -v amoebae
```

- Check that the amoebae executable script can be run (print the help message).

```
>>> amoebae -h
```

- Check that all modules can be imported in all python files in the AMOEBAE code.

```
>>> amoebae check_imports
```

- Check that key dependencies such as BLASTP can be accessed (they are installed in the Singularity container).

```
>>> amoebae check_depend
```

3. Again, running AMOEBAE commands on high-performance computing clusters will require you to write custom job submission scripts. Please refer to documentation provided by your system administrator(s) regarding details specific to your cluster, including the job scheduler used. Also, refer to the Singularity documentation for formulating Singularity commands (<https://sylabs.io/docs/>).

3 Command reference

Documentation for each AMOEBAE command and the various options may be accessed from the command line via the "-h" options. The following command reference information is the output of running amoebae (and each command) with the "-h" option.

3.1 amoebae

usage: amoebae <command> [<args>]

Commands for setting up data structure:

mkdatadir Make a directory with subdirectories and CSV files for
 storing sequence data, etc.

Commands for similarity searching:

setup_hmddb Construct an HMM database (with hmmpress).
add_to_dbs Format and add a file to a formatted directory.
list_dbs Print a list of all usable database files in the database
 directory as defined in the settings file.
add_to_queries Add a query file to a formatted directory.
list_queries Print a list of all usable query files in the query
 directory as defined in the settings file.
get_redun_hits Run searches with queries to find redundant hits in
 databases (for interpreting results).
setup_fwd_srch Make directory in which to perform forward searches.
run_fwd_srch Perform searches with given queries into given dbs.
sum_fwd_srch Append information about forward searches to csv summary
 file (this is used to organize reverse searches).
setup_rev_srch Make a directory in which to perform reverse searches.
run_rev_srch Perform searches with given forward search hits into given db.
sum_rev_srch Append information about reverse searches to csv summary
 file.
interp_srchs Interpret search results based on summary.
find_redun_seqs Identify sequences likely encoded on redundant loci
 predicted for the same species.
plot Plot search results.

Commands for phylogenetic analysis using a reference tree:

add_to_models Add an alignment, tree, substitution model, names of
 clade-defining sequences to a directory with other models.
list_models Print a list of all usable model/reference tree names in
 the models directory as defined in the settings file.
get_alt_topos Take a tree and make copies with every alternative
 topology for the branches connecting the clades of
 interest.

Commands for phylogenetic analysis without a reference tree:

prune Identify sequences in a tree, and remove them from a
 given alignment for further phylogenetic analysis.
auto_prune Automatically identify sequences in a tree, and remove
 them from a given alignment for further phylogenetic
 analysis.
reduce_tree Remove terminal nodes from a given tree if there are
 not any sequences with the same name in a given multiple
 sequence alignment file.

```

1      constrain_mb      Add constraint commands to MrBayes input file based on a
2                          given tree topology.
3      visualize_tree    Parse phylogenetic analysis output files for a single
4                          alignment in a given directory, and write human-readable
5                          tree figures to PDF files.
6      replace_seqs      Replace sequences in an alignment with their top hits in a
7                          given fasta file (useful if genomes or taxon selection has
8                          been updated).
9
10     Miscellaneous commands:
11         csv_to_fasta    Generate a fasta file from sequences detailed in a
12                          spreadsheet of similarity search results.
13         check_depend    Check that all the dependencies are properly installed and
14                          useable.
15         check_imports   Check that all the import statements used in the AMOEBAE
16                          repository run without error.
17         regen_genome_info Write a new genome info spreadsheet file using filenames
18                          from the Genomes directory.
19
20     positional arguments:
21         command          Specify one of the functionalities of amoebae.
22
23     optional arguments:
24         -h, --help      show this help message and exit
25
26     Copyright 2018 Lael D. Barlow Licensed under the Apache License, Version 2.0
27     (the "License"); you may not use this file except in compliance with the
28     License. You may obtain a copy of the License at
29     http://www.apache.org/licenses/LICENSE-2.0 Unless required by applicable law
30     or agreed to in writing, software distributed under the License is distributed
31     on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either
32     express or implied. See the License for the specific language governing
33     permissions and limitations under the License.

```

3.2 amoebae mkdatadir

```

35     usage: amoebae [-h] new_dir_path
36
37     Make a directory with subdirectories and CSV files for storing sequence data,
38     etc.
39
40     positional arguments:
41         new_dir_path      Specify the full file path that you want the new directory to
42                          have.
43
44     optional arguments:
45         -h, --help      show this help message and exit

```

3.3 amoebae setup_hmmdb

```
usage: amoebae [-h] indirpath

Construct an HMM database (with hmmpress). This is for later sorting of given
sequences into categories based on which HMM the score highest against.

positional arguments:
  indirpath  Path to directory containing amino acid sequence alignment
              file(s) to be constructed into an HMM database using hmmpress
              from the HMMer3 software package.

optional arguments:
  -h, --help  show this help message and exit
```

3.4 amoebae add_to_dbs

```
usage: amoebae [-h] [--split_char SPLIT_CHAR] [--split_pos SPLIT_POS]
               [--skip_header_reformat] [--auto_extract_accs]
               new_file main_data_dir

Format and add a file to a formatted directory.

positional arguments:
  new_file              Can be a fasta file (prot or nucl) or HMM databases,
                        generated using the hmmpress program in the HMMer
                        software package. Or a GFF3 annotation file.
  main_data_dir         Path to main data directory (with Genomes, Queries,
                        and Models subdirectories).

optional arguments:
  -h, --help            show this help message and exit
  --split_char SPLIT_CHAR
                        Character to split the header string on for extracting
                        the accession. (default: )
  --split_pos SPLIT_POS
                        Position that the accession will be in after
                        splitting. (default: 0)
  --skip_header_reformat
                        Skip reformatting of header lines in input fasta file.
                        (default: False)
  --auto_extract_accs   Automatically identify accessions/IDs in sequence
                        headers (overrides split_char and split_pos options
                        above). (default: False)
```

3.5 amoebae list_dbs

```
usage: amoebae [-h] main_data_dir
```

```

1
2 Print a list of all usable query files in the query directory as defined in a
3 given AMOEBAE data directory.
4
5 positional arguments:
6   main_data_dir  Path to main data directory (with Genomes, Queries, and
7                   Models subdirectories).
8
9 optional arguments:
10  -h, --help      show this help message and exit

```

11 3.6 amoebae add_to_queries

```

12 usage: amoebae [-h] query_file main_data_dir
13
14 Add a query file to a formatted directory. This command adds a given sequence
15 file to the directory with the path that you have specified in the settings.py
16 file, and appends a corresponding line to the CSV file that you specified
17 (e.g., '0_query_info.csv') to indicate the query title, etc.
18
19 positional arguments:
20   query_file      Path to a sequence file in FASTA format that can be used as a
21                   similarity search query file. Or path to a directory
22                   containing only files for addition to the queries. Note: By
23                   default, the portion of the input filename preceding the
24                   first underscore character will be recorded as the "query
25                   title", the remaining substring preceding the second
26                   underscore character will be recorded as the taxon (e.g.,
27                   "Hsapiens"), and the rest of the filename preceding the
28                   filename extension will be recorded as the sequence ID. So
29                   the filename might look like this:
30                   "QUERYTITLE_HSAPIENS_SEQUENCEID.fa". However, the relevant
31                   information can be revised in the "Queries/0_query_info.csv"
32                   file afterward if necessary.
33   main_data_dir   Path to main data directory (with Genomes, Queries, and
34                   Models subdirectories).
35
36 optional arguments:
37  -h, --help      show this help message and exit

```

38 3.7 amoebae list_queries

```

39 usage: amoebae [-h] main_data_dir
40
41 Print a list of all usable query files in the query directory as defined in a
42 given AMOEBAE data directory.
43
44 positional arguments:

```

```
1  main_data_dir Path to main data directory (with Genomes, Queries, and
2                  Models subdirectories).
```

```
3
4  optional arguments:
```

```
5  -h, --help      show this help message and exit
```

6 3.8 amoebae get_redun_hits

```
7 usage: amoebae [-h] [--query_name QUERY_NAME]
8                [--query_list_file QUERY_LIST_FILE] [--db_name DB_NAME]
9                [--db_list_file DB_LIST_FILE] [--query_title QUERY_TITLE]
10               [--blast_report_evalue_cutoff BLAST_REPORT_EVALUATE_CUTOFF]
11               [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
12               [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUATE_CUTOFF]
13               [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
14               [--max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE]
15               [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
16               ]
17               [--predict_redun_hit_selection] [--csv_file CSV_FILE]
18               out_dir_path main_data_dir
```

```
19
20 Run searches with queries to find redundant hits in databases (for
21 interpreting results).
```

```
22
23 positional arguments:
```

```
24  out_dir_path      Path to directory to write search results to.
25  main_data_dir     Path to main data directory (with Genomes, Queries,
26                    and Models subdirectories).
```

```
27
28 optional arguments:
```

```
29  -h, --help      show this help message and exit
30  --query_name QUERY_NAME
31                  Query filename to use (not full path). (default: None)
32  --query_list_file QUERY_LIST_FILE
33                  Path to file containing a list of query files to use,
34                  if no query_name is specified (or all queries by
35                  default). (default: None)
36  --db_name DB_NAME
37                  Name of database file in the database directory in
38                  which to do searches (not full path). (default: None)
39  --db_list_file DB_LIST_FILE
40                  Path to file containing a list of database files to
41                  use (if no db_name specified). (default: None)
42  --query_title QUERY_TITLE
43                  Name to be assigned to hits in databases that may be
44                  considered redundant with a search query to which the
45                  same title is assigned, otherwise it is taken from the
46                  query info spreadsheet specified in the settings.py
47                  file ('query_info_csv'). (default: None)
48  --blast_report_evalue_cutoff BLAST_REPORT_EVALUATE_CUTOFF
```



```

1             Maximum E-value for reporting BLAST hits. (default:
2             0.05)
3 --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
4             Maximum BLAST target sequences to consider. (default:
5             500)
6 --hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF
7             Maximum E-value for reporting HMMer hits. (default:
8             0.05)
9 --hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF
10            Minimum sequence score for reporting HMMer hits.
11            (default: 5)
12 --max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE
13            Absolute maximum number of hits (BLAST, HMMer, etc) to
14            summarize in the output spreadsheet. This is important
15            when working with sequences with WD40 domains, for
16            example. (default: 50)
17 --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
18            Number of threads to use for running searches.
19            (default: 4)
20 --predict_redun_hit_selection
21            Write a copy of the output spreadsheet with '+' in
22            rows for hits that may be specific to each query
23            title, due to not being retrieved as top hits by
24            queries associated with different query titles.
25            (default: False)
26 --csv_file CSV_FILE Path to spreadsheet to append summary of result to for
27            manual annotation. (default: None)
28
29 Recommendation: For most analyses, use the --query_name option and the
30 --db_name option, and run the get_redun_hits command for each query
31 separately. Otherwise, there will be redundant information in the output
32 spreadsheet(s).

```

3.9 amoebae setup_fwd_srch

```

34 usage: amoebae [-h] [--outdir OUTDIR] srch_dir query_list_file db_list_file
35
36 Make a directory in which to write output files from similarity searches.
37
38 positional arguments:
39   srch_dir             Path to directory that will contain output directory as a
40                       subdirectory.
41   query_list_file      Path to file with list of queries to search with.
42   db_list_file         Path to file with list of databases to search with.
43
44 optional arguments:
45   -h, --help           show this help message and exit
46   --outdir OUTDIR      Path to directory to put search results into (so that this
47                       step can be piped together with other commands). (default:

```

```

1             None)
2
3 Note: Use the bash script to run forward searches on a remote server.

```

3.10 amoebae run_fwd_srch

```

5 usage: amoebae [-h] [--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF]
6               [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
7               [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF]
8               [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
9               [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
10            ]
11            fwd_srch_dir main_data_dir
12
13 Perform searches with original queries into subject databases.
14
15 positional arguments:
16   fwd_srch_dir          Path to directory that will contain forward search
17                        output files.
18   main_data_dir         Path to main data directory (with Genomes, Queries,
19                        and Models subdirectories).
20
21 optional arguments:
22   -h, --help            show this help message and exit
23   --blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF
24                        Maximum E-value for reporting BLAST hits. (default:
25                        0.05)
26   --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
27                        Maximum BLAST target sequences to consider. (default:
28                        500)
29   --hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF
30                        Maximum E-value for reporting HMMer hits. (default:
31                        0.05)
32   --hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF
33                        Minimum sequence score for reporting HMMer hits.
34                        (default: 5)
35   --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
36                        Number of threads to use for running searches.
37                        (default: 4)

```

3.11 amoebae sum_fwd_srch

```

39 usage: amoebae [-h] [--csv_file CSV_FILE] [--max_evalue MAX_EVALUE]
40               [--max_gap_between_tblastn_hsps MAX_GAP_BETWEEN_TBLASTN_HSPS]
41               [--do_not_use_exonerate]
42               [--exonerate_score_threshold EXONERATE_SCORE_THRESHOLD]
43               [--max_hits_to_sum MAX_HITS_TO_SUM]
44               [--max_length_diff MAX_LENGTH_DIFF]

```

```

1         fwd_srch_out csv_out_path main_data_dir
2
3 Append information about forward searches to csv summary file (this is used to
4 organize reverse searches). For TBLASTN searches (protein queries, nucleotide
5 target sequences), HSPs are clustered into groups that are close enough within
6 the target sequence to potentially represent exons from the same coding
7 sequence. The nucleotide subsequences in which these clusters of HSPs are
8 found are then analyzed using exonerate to identify and translate potential
9 exons, in "protein2genome" mode, because exonerate, unlike TBLASTN, attempts
10 to identify exon boundaries, yielding translations that are less likely to
11 include translations of non-coding regions outside exons (which might include
12 apparent stop codons).
13
14 positional arguments:
15     fwd_srch_out          Path to directory where forward search results were
16                           written.
17     csv_out_path          Path to output summary spreadsheet (CSV) file.
18     main_data_dir         Path to main data directory (with Genomes, Queries,
19                           and Models subdirectories).
20
21 optional arguments:
22     -h, --help            show this help message and exit
23     --csv_file CSV_FILE   Path to summary spreadsheet (CSV) file, which already
24                           contains search summaries. If such a file is
25                           specified, then the output CSV file will contain the
26                           columns from this CSV file with additional columns
27                           summarizing additional forward search results.
28                           (default: None)
29     --max_evalue MAX_EVALUE
30                           Maximum E-value threshold for reporting forward search
31                           hits. (default: 0.0005)
32     --max_gap_between_tblastn_hsp MAX_GAP_BETWEEN_TBLASTN_HSPS
33                           Maximum number of nucleotide bases between TBLASTN
34                           HSPs to be considered part of the same gene locus.
35                           This is important, because it will be assumed that HSP
36                           separated by more than this number of nucleotide bases
37                           are not part of the same gene or TBLASTN "hit".
38                           (default: 10000)
39     --do_not_use_exonerate
40                           Override the default use of exonerate to identify
41                           coding sequences and translations, and just use
42                           TBLASTN instead. This option is provided because
43                           concatenated TBLASTN HSPs may be more inclusive of
44                           sequences within the target sequence, and the results
45                           of TBLASTN and exonerate may need to be compared.
46                           Also, note that HSPs identified by TBLASTN but for
47                           which exonerate yields no alignments will be ignored
48                           if exonerate is used. (default: False)
49     --exonerate_score_threshold EXONERATE_SCORE_THRESHOLD

```

```

1           Set score threshold to be applied when running
2           exonerate on nucleotide sequences identified by
3           TBLASTN. The default for setting of exonerate is 100,
4           but a lower score is set as default here, because
5           otherwise exonerate cannot identify some of the
6           sequences identified by TBLASTN. This option is only
7           relevant if using exonerate. (default: 10)
8   --max_hits_to_sum MAX_HITS_TO_SUM
9           Maximum number of forward search hits to list in the
10          summary spreadsheet. If zero, then reverse searches
11          will be performed for all hits. (default: 0)
12   --max_length_diff MAX_LENGTH_DIFF
13          Maximum number of amino acid residues length
14          difference allowed between the original query and the
15          forward hit sequence. If -1, then a maximum length
16          cutoff will not be applied. (default: -1)

```

17 3.12 amoebae setup_rev_srch

```

18 usage: amoebae [-h] [--outdir OUTDIR] [--aasubseq] [--nafullseq]
19           srch_dir csv_file databases main_data_dir
20
21 Make directory in which to write results of reverse searches.
22
23 positional arguments:
24   srch_dir           Path to directory that will contain output directory as a
25                     subdirectory.
26   csv_file           Path to summary spreadsheet (CSV) file, which contains a
27                     summary of forward search(es).
28   databases           Database filename (in database directory) or path to file
29                     with list of database filenames. Note that filenames are
30                     needed, not file paths.
31   main_data_dir       Path to main data directory (with Genomes, Queries, and
32                     Models subdirectories).
33
34 optional arguments:
35   -h, --help         show this help message and exit
36   --outdir OUTDIR    Path to directory to put search results into (so that this
37                     step can be piped together with other commands). (default:
38                     None)
39   --aasubseq         Use only the portion of each (amino acid) forward hit
40                     sequence that aligns to the original query used (top HSP
41                     subject sequence). This is default for nucleotide hits.
42                     (default: False)
43   --nafullseq        Use the full (nucleic acid) forward hit sequence. This is
44                     default for amino acid hits. (default: False)

```

3.13 amoebae run_rev_srch

```
usage: amoebae [-h] [--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF]
               [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
               [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF]
               [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
               [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
               ]
               rev_srch_dir main_data_dir
```

Perform searches with forward search hit sequences as queries into the original query databases.

positional arguments:

rev_srch_dir	Path to directory that will contain output of searches.
main_data_dir	Path to main data directory (with Genomes, Queries, and Models subdirectories).

optional arguments:

-h, --help	show this help message and exit
--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF	Maximum E-value for reporting BLAST hits. (default: 0.05)
--blast_max_target_seqs BLAST_MAX_TARGET_SEQS	Maximum BLAST target sequences to consider. (default: 500)
--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF	Maximum E-value for reporting HMMer hits. (default: 0.05)
--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF	Minimum sequence score for reporting HMMer hits. (default: 5)
--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING	Number of threads to use for running searches. (default: 4)

3.14 amoebae sum_rev_srch

```
usage: amoebae [-h] [--redun_hit_csv REDUN_HIT_CSV]
               [--min_evaldiff MIN_EVALDIFF] [--aasubseq] [--nafullseq]
               [--max_rev_srchs MAX_REV_SRCHS]
               fwd_srch_csv rev_srch_out csv_out_path main_data_dir
```

Append information about reverse searches to csv summary file. Use information from redundant hit csv file to interpret results.

positional arguments:

fwd_srch_csv	Path to summary spreadsheet (CSV) file, which contains
--------------	--

```

1         forward search summaries and also may already contain
2         reverse search summaries.
3     rev_srch_out      Path to directory where reverse search results were
4                       written.
5     csv_out_path      Path output spreadsheet (CSV) file with reverse search
6                       results appended to previous results.
7     main_data_dir     Path to main data directory (with Genomes, Queries,
8                       and Models subdirectories).
9
10    optional arguments:
11    -h, --help          show this help message and exit
12    --redun_hit_csv REDUN_HIT_CSV
13                       Path to spreadsheet (CSV) file, which specifies which
14                       hits are redundant positive hits for a given query
15                       (query title) in a given database. If this is not
16                       provided, then it is assumed that any and all reverse
17                       search hits are equivalent to/redundant with the
18                       original query. (default: None)
19    --min_evaldiff MIN_EVALDIFF
20                       Minimum difference in E-value order of magnitude
21                       between top reverse search hit and first reverse
22                       search hit that is not redundant with the original
23                       query. (default: 5)
24    --aasubseq          Use only the portion of each (amino acid) forward hit
25                       sequence that aligns to the original query used (top
26                       HSP subject sequence). This is default for nucleotide
27                       hits. Must be selected if selected when the
28                       setup_rev_srch command was run. (default: False)
29    --nafullseq         Use the full (nucleic acid) forward hit sequence. This
30                       is default for amino acid hits. Must be selected if
31                       selected when the setup_rev_srch command was run.
32                       (default: False)
33    --max_rev_srchs MAX_REV_SRCHS
34                       Maximum number of forward search hits to perform
35                       reverse searches for per query database. If zero, then
36                       reverse searches will be performed for all hits.
37                       (default: 0)

```

3.15 amoebae interp_srchs

```

39    usage: amoebae [-h] [--fwd_only] [--fwd_evalue_cutoff FWD_EVALUE_CUTOFF]
40                  [--rev_evalue_cutoff REV_EVALUE_CUTOFF]
41                  [--hmmer_cutoff HMMER_CUTOFF] [--no_overlapping_hits]
42                  [--out_csv_path OUT_CSV_PATH]
43                  csv_file
44
45    Interpret search results based on final summary, which provides a basis for
46    further analyses of positive hits.
47

```

```

1 positional arguments:
2   csv_file           Path to spreadsheet with forward and reverse search
3                       results.
4
5 optional arguments:
6   -h, --help         show this help message and exit
7   --fwd_only         Interpret forward searches based on score (HMMer)
8                       cutoff. (default: False)
9   --fwd_evalue_cutoff FWD_EVALUE_CUTOFF
10                      Specify an (more stringent) E-value cutoff for forward
11                      search results. (default: None)
12   --rev_evalue_cutoff REV_EVALUE_CUTOFF
13                      Specify an (more stringent) E-value cutoff for reverse
14                      search results. (default: None)
15   --hmmmer_cutoff HMMER_CUTOFF
16                      Specify a score that hits must exceed to be included.
17                      (default: 20)
18   --no_overlapping_hits
19                      If more than one query (query title) retrieves a given
20                      sequence as a positive hit based on the search
21                      criteria, make the sequence a negative hit for all
22                      queries (query titles), except for the one that
23                      retrieved the sequence with the lowest (strongest)
24                      E-value. Warning: Do not use this option if you are
25                      searching sequences that include genomic sequences
26                      that may include more than one genomic locus per
27                      sequence. False-negative results could occur in this
28                      case, because different queries for non-orthologous
29                      genes could retrieve subsequences in the same subject
30                      sequence. (default: False)
31   --out_csv_path OUT_CSV_PATH
32                      Optionally specify an output file path, so that this
33                      command can be piped together with others. (default:
34                      None)

```

3.16 amoebae find_redun_seqs

```

36 usage: amoebae [-h] [--out_csv_path OUT_CSV_PATH]
37                [--remove_tblastn_hits_at_annotated_loci]
38                [--just_look_for_genes_in_gff3] [--ignore_gff3]
39                [--allow_internal_stops ALLOW_INTERNAL_STOPS]
40                [--min_length MIN_LENGTH]
41                [--min_percent_length MIN_PERCENT_LENGTH]
42                [--min_percent_query_cover MIN_PERCENT_QUERY_COVER]
43                [--overlap_required] [--max_percent_ident MAX_PERCENT_IDENT]
44                [--min_alig_res_in_overlap MIN_ALIG_RES_IN_OVERLAP]
45                [--min_ident_res_in_overlap MIN_IDENT_RES_IN_OVERLAP]
46                [--min_sim_res_in_overlap MIN_SIM_RES_IN_OVERLAP]
47                [--min_ident_span_len MIN_IDENT_SPAN_LEN]

```

```

1      [--min_sim_span_len MIN_SIM_SPAN_LEN]
2      [--min_percent_ident_in_overlap MIN_PERCENT_IDENT_IN_OVERLAP]
3      [--min_percent_sim_in_overlap MIN_PERCENT_SIM_IN_OVERLAP]
4      [--min_percent_overlap MIN_PERCENT_OVERLAP]
5      [--plot_hit_exclusion] [--add_ali_col]
6      csv_file main_data_dir
7
8  Identify hit sequences likely encoded by the same gene loci in the genome of a
9  given species, or otherwise not representing paralogous genes. Criteria are
10 applied in this order: 1. Peptide hits with the same ID as higher-ranking hits
11 for the same query (query title) are excluded. 2. Nucleotide hits for the same
12 loci as peptide sequence hits are excluded. 3. Sequences with internal stop
13 codons are excluded, as these are potentially pseudogenes. 4. Sequences are
14 excluded if they do not meet several minimum length criteria: Absolute minimum
15 length (in amino acids) and percent query cover. 5. Sequences are excluded if
16 they do not overlap to a specified degree with all included higher-ranking
17 hits for the same query (query title) in sequence data for the same
18 species/genome. This is determined by algorithmically comparing pairs of
19 sequences aligned to a reference alignment of homologues, and several minimum
20 measures of alignment overlap may be specified. 6. Secondary hit sequences are
21 excluded if they do not meet a specified maximum percent identity threshold.
22 Highly identical sequences may result from false segmental duplications in the
23 genome assembly, may represent alleles, etc. Note: Applying these criteria
24 requires a column to be manually added to the input csv file prior to running
25 with the header "Alignment for sequence comparison" and filled with the
26 appropriate alignment name to use (one for each query title, as listed in the
27 "Query title" column). Alternatively, you can run this command with the
28 --add_ali_col option to automatically identify appropriate alignments among
29 your aligned FASTA queries used for running HMMer searches. If no alignment
30 (.afaa) file can be found, then the first single sequence query file (.faa)
31 that appears in the summary CSV file will be used instead.
32
33 positional arguments:
34   csv_file              Path to spreadsheet with interpreted search results
35                        outputted by the interp_srchs command.
36   main_data_dir         Path to main data directory (with Genomes, Queries,
37                        and Models subdirectories).
38
39 optional arguments:
40   -h, --help            show this help message and exit
41   --out_csv_path OUT_CSV_PATH
42                        Optionally specify an output file path, so that this
43                        command can be piped together with others. (default:
44                        None)
45   --remove_tblastn_hits_at_annotated_loci
46                        Ignore tblastn hits that overlap with any previously
47                        annotated loci. The rationale for this would be that
48                        the corresponding protein sequences should have been
49                        retrieved if the tblastn hit were a true positive

```


1 anyway. If this option is not specified, then
 2 sequences will still be excluded if they specifically
 3 correspond to the same loci as do higher-ranking hits.
 4 (default: False)
 5 --just_look_for_genes_in_gff3
 6 When looking for records in GFF3 annotation files that
 7 overlap with subsequences identified by similarity
 8 searching (TBLASTN), ignore records that are not
 9 explicitly "gene" (for example, "CDS", "mRNA", and
 10 "exon"). This option should probably not be selected,
 11 because in some GFF3 annotation files do not include
 12 "gene" records, but do include predicted coding
 13 sequences for genes. (default: False)
 14 --ignore_gff3 Disregard any information regarding redundancy of
 15 identified nucleotide sequences with identified
 16 protein sequences that may be found in GFF3 annotation
 17 files. (default: False)
 18 --allow_internal_stops ALLOW_INTERNAL_STOPS
 19 Include sequences that have internal stop codons
 20 (anywhere other than the N-terminal position).
 21 (default: True)
 22 --min_length MIN_LENGTH
 23 Absolute minimum length (in AA) of a hit sequence to
 24 be considered a potential distinct paralogue.
 25 (default: 55)
 26 --min_percent_length MIN_PERCENT_LENGTH
 27 Minimum length (in AA) of a hit sequence as a
 28 percentage of query length for the hit to be
 29 considered a potential distinct paralogue. (default:
 30 15)
 31 --min_percent_query_cover MIN_PERCENT_QUERY_COVER
 32 Minimum number of residues aligning with the original
 33 query as a percentage of the original query sequence
 34 length. (default: 0)
 35 --overlap_required True if hits must overlap with a higher-ranking hit to
 36 be considered potential unique paralogues. (default:
 37 False)
 38 --max_percent_ident MAX_PERCENT_IDENT
 39 Maximum percent identity (among aligning residues) for
 40 evaluating whether two sequences are redundant or not
 41 (secondary hits showing a percent identity with a
 42 higher-ranking hit exceeding this value will be
 43 excluded). (default: 98.0)
 44 --min_alig_res_in_overlap MIN_ALIG_RES_IN_OVERLAP
 45 Minimum number of residues which must align for two
 46 sequences to be considered as potentially distinct
 47 hits. This is only relevant if the overlap_required
 48 option is specified. (default: 20)
 49 --min_ident_res_in_overlap MIN_IDENT_RES_IN_OVERLAP

```

1             Minimum number of aligning residues which must be
2             identical for two sequences to be considered as
3             potentially distinct hits. This is only relevant if
4             the overlap_required option is specified. (default:
5             10)
6  --min_sim_res_in_overlap MIN_SIM_RES_IN_OVERLAP
7             Minimum number of aligning residues which must be
8             similar for two sequences to be considered as
9             potentially distinct hits. This is only relevant if
10            the overlap_required option is specified. (default:
11            15)
12  --min_ident_span_len MIN_IDENT_SPAN_LEN
13            Minimum number of aligning residues which are
14            identical that must exist in at least one continuous
15            span for two sequences to be considered as potentially
16            distinct hits (not counting positions where both
17            sequences have gaps). This is only relevant if the
18            overlap_required option is specified. (default: 0)
19  --min_sim_span_len MIN_SIM_SPAN_LEN
20            Minimum number of aligning residues which are similar
21            (or identical) that must exist in at least one
22            continuous span for two sequences to be considered as
23            potentially distinct hits (not counting positions
24            where both sequences have gaps). This is only relevant
25            if the overlap_required option is specified. (default:
26            0)
27  --min_percent_ident_in_overlap MIN_PERCENT_IDENT_IN_OVERLAP
28            Minimum percent identity between the two sequences of
29            interest in the alignment. This is only relevant if the
30            overlap_required option is specified. (default: 0)
31  --min_percent_sim_in_overlap MIN_PERCENT_SIM_IN_OVERLAP
32            Minimum percent similarity (including identity)
33            between the two sequences of interest in the
34            alignment. This is only relevant if the
35            overlap_required option is specified. (default: 0)
36  --min_percent_overlap MIN_PERCENT_OVERLAP
37            Minimum number of aligning residues between the two
38            sequences of interest as a percentage of the length of
39            the second sequence (the last sequence in the
40            alignment), not including gaps, for the two sequences
41            to be considered as potentially distinct hits. This is
42            only relevant if the overlap_required option is
43            specified. (default: 0)
44  --plot_hit_exclusion Plot number of hits excluded by the various criteria
45            applied. (default: False)
46  --add_ali_col Add a column to the csv file listing which alignment
47            file in the queries directory to use for comparing
48            sequences. Aligned FASTA queries are selected that
49            match the query titles of the original queries used to

```

```

1         retrieve each of the relevant hits listed in the csv
2         file. No other options need to be specified in this
3         case. (default: False)

```

3.17 amoebae plot

```

5 usage: amoebae [-h] [--csv_file2 CSV_FILE2] [--complex_info COMPLEX_INFO]
6               [--row_order ROW_ORDER] [--out_pdf OUT_PDF]
7               csv_file
8
9 Plot results of similarity search and sequence classification analyses. The
10 outputs are PDF files.
11
12 positional arguments:
13   csv_file              Path to a spreadsheet with the relevant results to be
14                         plotted. This can be either a CSV file output of the
15                         sum_rev_srch command or from the find_redun_seqs
16                         command. If the output of the sum_rev_srch command is
17                         used, however, redundant hits will be counted (e.g.,
18                         BLASTP and TBLASTN hits corresponding to the same or
19                         highly identical genomic loci).
20
21 optional arguments:
22   -h, --help            show this help message and exit
23   --csv_file2 CSV_FILE2
24                         Path to a second spreadsheet with relevant results to
25                         be compared to the first and plotted. (default: None)
26   --complex_info COMPLEX_INFO
27                         Path to file that specifies which query titles
28                         represent components of which protein complexes (or
29                         otherwise grouped proteins). (default: None)
30   --row_order ROW_ORDER
31                         Path to file that specifies the order in which data
32                         for each species will be displayed. (default: None)
33   --out_pdf OUT_PDF     Path to output pdf file. (default: None)

```

3.18 amoebae add_to_models

```

35 usage: amoebae [-h]
36               model_name alignment tree_topology subs_model type_seqs taxon
37
38 Add a phylogenetic model for relationships between members of a gene family
39 (sequence_data matrix, data type, tree topology, type sequence defining each
40 clade of interest, and substitution model) to a directory for use in
41 classifying sequence (via the 'phylo_class' command).
42
43 positional arguments:
44   model_name            An arbitrary name for the model (which will refer to the

```

```

1         alignment      alignment, tree, substitution model, etc. collectively).
2     alignment          A multiple amino acid sequence alignment in nexus format.
3     tree_topology      Text file containing a tree (identified previously using
4                        MrBayes, etc) containing the names of all the sequences in
5                        the alignment, in newick format.
6     subs_model          The name of the substitution model used to recover the
7                        provided topology (chosen with ModelFinder or similar
8                        software).
9     type_seqs           Names of sequences (sequence headers) that are to be used to
10                        define clades of interest. A csv file with seq names in one
11                        column and clade names in the next column.
12     taxon               Taxonomic group represented in the model (e.g., "Eukaryotes",
13                        or "Amorphea").
14
15 optional arguments:
16     -h, --help          show this help message and exit

```

17 3.19 amoebae list_models

```

18 usage: amoebae [-h]
19
20 Print a list of all usable model/reference tree names in the models directory
21 as defined in the settings file.
22
23 optional arguments:
24     -h, --help          show this help message and exit

```

25 3.20 amoebae get_alt_topos

```

26 usage: amoebae [-h] [--polytomy] [--not_polytomy_clades]
27                [--keep_original_backbone] [--iqtree_au_test]
28                model_name out_dir_path
29
30 Take a tree and make copies with every alternative topology for the branches
31 connecting the clades of interest. Output as additional models in the Models
32 directory.
33
34 positional arguments:
35     model_name           Name of model/backbone tree to modify (other info
36                        provided in the model info csv file).
37     out_dir_path         Path to directory in which output directory will be
38                        written.
39
40 optional arguments:
41     -h, --help          show this help message and exit
42     --polytomy           Just make one big polytomy connecting the clades of
43                        interest instead of making alternative bifurcating
44                        trees. (default: False)

```

```

1  --not_polytomy_clades
2      Do not make subtrees/clades of interest polytomies in
3      output topologies. (default: False)
4  --keep_original_backbone
5      Keep the original backbone topology instead of
6      generating a polytomy or alternative resolved
7      topologies. (default: False)
8  --iqtree_au_test      Test all the relevant alternative topologies against
9      each other using Approximately Unbiased (AU) test with
10     IQ-tree. (default: False)

```

11 3.21 amoebae prune

```

12 usage: amoebae [-h] [--include_seqs] [--output_file OUTPUT_FILE]
13               tree_file alignment name_replace_table
14
15 Identify sequences in a tree, and remove them from a given alignment for
16 further phylogenetic analysis.
17
18 positional arguments:
19   tree_file            Tree in newick format (coded names, because ETE3
20                       cannot parse taxon names with space characters without
21                       quotation marks around them).
22   alignment            Dataset used to make the tree (nexus alignment)
23                       (original alignment with original taxon names either
24                       trimmed or untrimmed).
25   name_replace_table   File for decoding names in input tree file.
26
27 optional arguments:
28   -h, --help           show this help message and exit
29   --include_seqs       Include only listed sequences/nodes instead of
30                       removing them. (default: False)
31   --output_file OUTPUT_FILE
32                       Path to output file. (default: None)

```

33 3.22 amoebae auto_prune

```

34 usage: amoebae [-h]
35               [--max_bl_iqr_above_third_quartile MAX_BL_IQR_ABOVE_THIRD_QUARTILE]
36               [--remove_redun_seqs REMOVE_REDUN_SEQS]
37               [--remove_redun_seqs_threshold REMOVE_REDUN_SEQS_THRESHOLD]
38               [--output_file OUTPUT_FILE]
39               in_dir
40
41 Automatically identify sequences in a tree, and remove them from a given
42 alignment for further phylogenetic analysis.
43
44 positional arguments:

```

```

1  in_dir          Path to directory that contains the phylogenetic
2                  analysis output files (sequence name conversion table
3                  file and original nexus alignment file can be in the
4                  parent directory to this directory as long as their
5                  names are mostly identical.
6
7  optional arguments:
8      -h, --help          show this help message and exit
9      --max_bl_iqr_above_third_quartile MAX_BL_IQR_ABOVE_THIRD_QUARTILE
10                          Inclusion threshold for number of interquartile ranges
11                          above the third quartile of terminal branch lengths
12                          the length of a terminal branch can be before it is
13                          considered an outlier (length is total distance from
14                          root node after rooting on midpoint, or the longest
15                          terminal branch on either side of the midpoint).
16                          (default: 1.5)
17      --remove_redun_seqs REMOVE_REDUN_SEQS
18                          Remove taxonomically redundant sequences (longest
19                          branch of two sister branches when both are sequences
20                          from the same species. (default: True)
21      --remove_redun_seqs_threshold REMOVE_REDUN_SEQS_THRESHOLD
22                          Minimum support required to consider one of two sister
23                          branches/sequences taxonomically redundant. Note: only
24                          used if the remove_redun_seqs option is specified.
25                          (default: 0.95)
26      --output_file OUTPUT_FILE
27                          Path to output file. (default: None)

```

28 3.23 amoebae reduce_tree

```

29 usage: amoebae [-h] [--output_file OUTPUT_FILE] alignment tree_file
30
31 Remove terminal nodes from a given tree if there are not any sequences with
32 the same name in a given alignment.
33
34 positional arguments:
35     alignment          Alignment in nexus format with sequences representing
36                       a subset of those represented in the input tree.
37     tree_file          Tree in newick format.
38
39 optional arguments:
40     -h, --help          show this help message and exit
41     --output_file OUTPUT_FILE
42                       Path to output file. (default: None)

```

43 3.24 amoebae constrain_mb

```

44 usage: amoebae [-h] [--out_alignment OUT_ALIGNMENT] alignment tree

```

```

1
2 Add constraint commands to MrBayes input file.
3
4 positional arguments:
5     alignment      Nexus alignment for input to MrBayes (without any
6                     constraint commands).
7     tree           Tree in newick format with same taxon names as in
8                     alignment. To be used as a topology constraint (all
9                     nodes).
10
11 optional arguments:
12     -h, --help      show this help message and exit
13     --out_alignment OUT_ALIGNMENT
14                     Path to nexus alignment for input to MrBayes with
15                     constraints added. (default: None)

```

16 3.25 amoebae visualize_tree

```

17 usage: amoebae [-h] [--root_taxon ROOT_TAXON] [--highlight_paralogues]
18                [--add_clade_names_from_file]
19                input_directory method
20
21 Parse phylogenetic analysis output files in a given directory, and write
22 human-readable tree figures to PDF files.
23
24 positional arguments:
25     input_directory Path to directory containing input files (must contain
26                     a .table file for decoding taxon names.
27     method          Name of tree searching program used. Either iqtree,
28                     raxml, or mrbayes accepted.
29
30 optional arguments:
31     -h, --help      show this help message and exit
32     --root_taxon ROOT_TAXON
33                     Name of species to root on (e.g.,
34                     "Klebsormidium_nitens").
35     --highlight_paralogues
36                     Highlight clades that contain paralogues found in at
37                     least one other clade in the tree.
38     --add_clade_names_from_file
39                     Use a file in the parent directory with clade names
40                     corresponding to representative sequences to add clade
41                     names to all the taxon names in the output trees.

```

42 3.26 amoebae replace_seqs

```

43 usage: amoebae [-h] [--fasta_file FASTA_FILE] alignment
44

```

1 Replace sequences in an alignment the full-length sequences from the relevant
 2 file(s) in the Genomes directory, or with their top hits in a given fasta
 3 file. And, align, mask, and trim the identified sequences to the input
 4 alignment

5

6 positional arguments:

7 alignment Path to multiple sequence alignment file in nexus
 8 format (trimmed alignment).

9

10 optional arguments:

11 -h, --help show this help message and exit

12 --fasta_file FASTA_FILE

13 Path to file containing sequences with which to
 14 replace sequences in the alignment. If this option is
 15 not specified, then full-length sequences will be
 16 retrieved from files in the Genomes directory.

17 3.27 amoebae csv_to_fasta

18 usage: amoebae [-h] [--output_dir OUTPUT_DIR] [--abbrev] [--parologue_names]
 19 [--only_descr] [--subseq] [--all_hits] [--split_by_query_title]
 20 [--split_by_top_rev_srch_hit SPLIT_BY_TOP_REV_SRCH_HIT]
 21 [--split_to_query_fastas]
 22 csv_file

23

24 Extract sequences described in a spreadsheet output by AMOEBAE, and write to a
 25 file in FASTA format.

26

27 positional arguments:

28 csv_file Path to csv file listing sequences.

29

30 optional arguments:

31 -h, --help show this help message and exit

32 --output_dir OUTPUT_DIR

33 Path for output directory to contain FASTA files.
 34 (default: None)

35 --abbrev Add species name instead of sequence description from
 36 fasta header. Applicable when the output file is to be
 37 used for alignment and phylogenetic analysis.

38 (default: False)

39 --parologue_names Use species name, query title, and parologue number
 40 instead of sequence description from fasta header.
 41 Applicable when the output file is to be used for
 42 alignment and phylogenetic analysis. Does not work if
 43 the abbrev option is specified. (default: False)

44 --only_descr Use the description but not the ID as the new fasta
 45 sequence header. Does not work if the abbrev option is
 46 specified. (default: False)

47 --subseq Write subsequences that aligned to forward search


```

1          query, instead of the full sequences. (default: False)
2  --all_hits      Write all forward hits listed in the input csv file.
3                  (default: False)
4  --split_by_query_title
5                  Write sequences to files according to the query title
6                  of the query which retrieved them in a similarity
7                  search. (default: False)
8  --split_by_top_rev_srch_hit SPLIT_BY_TOP_REV_SRCH_HIT
9                  Write sequences to files according to the top hit that
10                 they retrieve in a reverse search, for each sequence
11                 that meets the reverse search criteria. (Provide the
12                 reverse search identifier, eg,
13                 "rev_srch_20180924122402-1") (default: None)
14  --split_to_query_fastas
15                 Write sequences to separate files with filenames that
16                 can be easily parsed for loading the the files as
17                 queries using the add_to_queries command. (default:
18                 False)

```

19 3.28 amoebae check_depend

```

20 usage: amoebae [-h]
21
22 Check that all the dependencies (other than python modules) are properly
23 installed and useable.
24
25 optional arguments:
26  -h, --help  show this help message and exit

```

27 3.29 amoebae check_imports

```

28 usage: amoebae [-h]
29
30 Check that all the import statements used in the AMOEBAE repository run
31 without error.
32
33 optional arguments:
34  -h, --help  show this help message and exit

```

35 3.30 amoebae regen_genome_info

```

36 usage: amoebae [-h] data_dir_path
37
38 Write a new genome info spreadsheet (0_genome_info.csv) file using filenames
39 from the Genomes directory.
40
41 positional arguments:

```

```
1 data_dir_path Specify the full path to an existing AMOEBAE data directory,  
2               which contains a 'Genomes' subdirectory. The new genome info  
3               file will be added to this subdirectory.  
4  
5 optional arguments:  
6 -h, --help      show this help message and exit
```

7 4 Miscellaneous scripts

8 Several scripts of less general applicablity than the amoebae commands descibed above
9 are included in the AMOEBAE toolkit. See the amoebae/misc_scripts directory ([https:](https://github.com/laelbarlow/amoebae/tree/master/misc_scripts)
10 [//github.com/laelbarlow/amoebae/tree/master/misc_scripts](https://github.com/laelbarlow/amoebae/tree/master/misc_scripts)). Most scripts have in-
11 formation regarding usage in the files themselves. More detailed information regarding some
12 of these scripts may be added to this documentation in the future.

5 References

- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L. (2009). BLAST+: Architecture and applications. *BMC Bioinformatics*, 10(1):421. doi:10.1186/1471-2105-10-421.
- Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., and de Hoon, M.J.L. (2009). Biopython: Freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423. doi:10.1093/bioinformatics/btp163.
- Eddy, S.R. (1998). Profile hidden Markov models. *Bioinformatics*, 14(9):755–763. doi:10.1093/bioinformatics/14.9.755.
- Edgar, R.C. (2004). MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*, 32(5):1792–1797. doi:10.1093/nar/gkh340.
- Emms, D.M. and Kelly, S. (2019). OrthoFinder: Phylogenetic orthology inference for comparative genomics. *Genome Biology*, 20(1):238. doi:10.1186/s13059-019-1832-y.
- Huerta-Cepas, J., Serra, F., and Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Molecular Biology and Evolution*, 33(6):1635–1638. doi:10.1093/molbev/msw046.
- Hunter, J.D. (2007). Matplotlib: A 2D Graphics Environment. *Computing in Science Engineering*, 9(3):90–95. doi:10.1109/MCSE.2007.55.
- Larson, R.T., Dacks, J.B., and Barlow, L.D. (2019). Recent gene duplications dominate evolutionary dynamics of adaptor protein complex subunits in embryophytes. *Traffic*, page tra.12698. doi:10.1111/tra.12698.
- Li, L. (2003). OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes. *Genome Research*, 13(9):2178–2189. doi:10.1101/gr.1224503.
- Nguyen, L.T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology and Evolution*, 32(1):268–274. doi:10.1093/molbev/msu300.
- Slater, G. and Birney, E. (2005). [No title found]. *BMC Bioinformatics*, 6(1):31. doi:10.1186/1471-2105-6-31.