

AMOEBAE documentation

Lael D. Barlow

Version of July 7, 2020

Contents

| | | |
|----------|---|----------|
| 1 | Introduction | 1 |
| 1.1 | What is AMOEBAE? | 1 |
| 1.2 | Why use AMOEBAE? | 1 |
| 1.3 | Key features | 1 |
| 1.4 | A word of caution | 2 |
| 1.5 | User support | 2 |
| 1.6 | Documentation | 2 |
| 1.7 | How to cite AMOEBAE | 2 |
| 1.8 | Acknowledgments | 3 |
| 1.9 | License | 3 |
| 2 | How to start using AMOEBAE | 3 |
| 2.1 | System requirements | 3 |
| 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 | 7 |
| 3.2 | amoebae mkdatadir | 9 |
| 3.3 | amoebae setup_hmmdb | 9 |
| 3.4 | amoebae add_to_dbs | 10 |
| 3.5 | amoebae list_dbs | 10 |
| 3.6 | amoebae add_to_queries | 10 |

| | | |
|------|-------------------------------------|----|
| 3.7 | amoebae list_queries | 11 |
| 3.8 | amoebae get_redun_hits | 11 |
| 3.9 | amoebae setup_fwd_srch | 13 |
| 3.10 | amoebae run_fwd_srch | 13 |
| 3.11 | amoebae sum_fwd_srch | 14 |
| 3.12 | amoebae setup_rev_srch | 15 |
| 3.13 | amoebae run_rev_srch | 16 |
| 3.14 | amoebae sum_rev_srch | 16 |
| 3.15 | amoebae interp_srchs | 17 |
| 3.16 | amoebae find_redun_seqs | 18 |
| 3.17 | amoebae plot | 22 |
| 3.18 | amoebae add_to_models | 22 |
| 3.19 | amoebae list_models | 23 |
| 3.20 | amoebae get_alt_topos | 23 |
| 3.21 | amoebae prune | 24 |
| 3.22 | amoebae auto_prune | 24 |
| 3.23 | amoebae reduce_tree | 25 |
| 3.24 | amoebae constrain_mb | 26 |
| 3.25 | amoebae visualize_tree | 26 |
| 3.26 | amoebae replace_seqs | 27 |
| 3.27 | amoebae csv_to_fasta | 27 |
| 3.28 | amoebae check_depend | 28 |
| 3.29 | amoebae check_imports | 28 |
| 3.30 | amoebae regen_genome_info | 28 |

4 Miscellaneous scripts 29

1 Introduction

1.1 What is AMOEBAE?

Analysis of MOlecular Evolution with BAtch 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?

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 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 done 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 100 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 (searching in 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

1 Birney, 2005) where possible for obtaining better exon predictions.

2 1.4 A word of caution

3 AMOEBAE is not optimized for ease of use, but is meant to be highly configurable. The
4 many options available to AMOEBAE users inevitably provide many opportunities for user
5 errors in specifying search criteria, and user errors in interpreting results detailed in output
6 files. Some prior experience with similarity searching and with running software using the
7 command line are prerequisites for using AMOEBAE, and experience writing scripts in Bash
8 (linux shell) and Python would be highly advantageous. Also, you may need to carefully
9 define the scope of your analysis depending on what additional steps you may find necessary
10 beyond those that may be performed using AMOEBAE (you may find that the maximum
11 30 queries and 100 genomes suggested above may in fact be unmanageable). Moreover,
12 AMOEBAE is still under active development, so some features may not yet be thoroughly
13 tested.

14 1.5 User support

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

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

19 1.6 Documentation

20 This document provides an overview of AMOEBAE and describes the functionality of the
21 various commands/scripts. For a tutorial which includes a working example of a similarity
22 search analysis run using AMOEBAE, see the Jupyter Notebook: `amoebae/notebooks/sim-`
23 `ilarity_search_tutorial_2.ipynb`. For code documentation, please see the html file(s), which
24 can be opened with your web browser: `amoebae/documentation/code_documentation/`
25 `html/index.html`.

26 1.7 How to cite AMOEBAE

27 Please cite the GitHub webpage <https://github.com/laelbarlow/amoebae> (or alternative
28 permanent repositories if relevant). Also, the first publication to make use of a version of
29 AMOEBAE was an analysis of Adaptor Protein subunits in embryophytes by Larson *et al.*
30 (2019).

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

3 1.8 Acknowledgments

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

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

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

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

14 1.9 License

15 Copyright 2018 Lael D. Barlow

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

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

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

23 2 How to start using AMOEBAE

24 2.1 System requirements

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

2.2 Dependencies

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 with instructions for how to install it when you type git into the command line. If you have a newer version of MacOS it may prompt you to install developer tools, which may take 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.

```

1      • Print the paths included in the $PATH variable in the container.
2      >>> tr ':' '\n' <<< "$PATH"
3
4      • Check the location of the amoebae executable being run from within the container.
5      >>> command -v amoebae
6
7      • Check that the amoebae executable script can be run (print the help message).
8      >>> amoebae -h
9
10     • Check that all modules can be imported in all python files in the AMOEBAE
11     code.
12     >>> amoebae check_imports
13
14     • Check that key dependencies such as BLASTP can be accessed (they are installed
15     in the Singularity container).
16     >>> amoebae check_depend
17
18 3. Again, running AMOEBAE commands on high-performance computing clusters will
19  require you to write custom job submission scripts. Please refer to documentation
20  provided by your system administrator(s) regarding details specific to your cluster,
21  including the job scheduler used. Also, refer to the Singularity documentation for
22  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

```

23 usage: amoebae <command> [<args>]
24
25 Commands for setting up data structure:
26     mkdatadir      Make a directory with subdirectories and CSV files for
27                     storing sequence data, etc.
28
29 Commands for similarity searching:
30     setup_hmmdb    Construct an HMM database (with hmmpress).
31     add_to_dbs      Format and add a file to a formatted directory.
32     list_dbs        Print a list of all usable database files in the database
33                     directory as defined in the settings file.
34     add_to_queries  Add a query file to a formatted directory.

```

| | | |
|----|--|--|
| 1 | list_queries | Print a list of all usable query files in the query |
| 2 | | directory as defined in the settings file. |
| 3 | get_redun_hits | Run searches with queries to find redundant hits in |
| 4 | | databases (for interpreting results). |
| 5 | setup_fwd_srch | Make directory in which to perform forward searches. |
| 6 | run_fwd_srch | Perform searches with given queries into given dbs. |
| 7 | sum_fwd_srch | Append information about forward searches to csv summary |
| 8 | | file (this is used to organize reverse searches). |
| 9 | setup_rev_srch | Make a directory in which to perform reverse searches. |
| 10 | run_rev_srch | Perform searches with given forward search hits into given db. |
| 11 | sum_rev_srch | Append information about reverse searches to csv summary |
| 12 | | file. |
| 13 | interp_srchs | Interpret search results based on summary. |
| 14 | find_redun_seqs | Identify sequences likely encoded on redundant loci |
| 15 | | predicted for the same species. |
| 16 | plot | Plot search results. |
| 17 | | |
| 18 | Commands for phylogenetic analysis using a reference tree: | |
| 19 | add_to_models | Add an alignment, tree, substitution model, names of |
| 20 | | clade-defining sequences to a directory with other models. |
| 21 | list_models | Print a list of all usable model/reference tree names in |
| 22 | | the models directory as defined in the settings file. |
| 23 | get_alt_topos | Take a tree and make copies with every alternative |
| 24 | | topology for the branches connecting the clades of |
| 25 | | interest. |
| 26 | | |
| 27 | Commands for phylogenetic analysis without a reference tree: | |
| 28 | prune | Identify sequences in a tree, and remove them from a |
| 29 | | given alignment for further phylogenetic analysis. |
| 30 | auto_prune | Automatically identify sequences in a tree, and remove |
| 31 | | them from a given alignment for further phylogenetic |
| 32 | | analysis. |
| 33 | reduce_tree | Remove terminal nodes from a given tree if there are |
| 34 | | not any sequences with the same name in a given multiple |
| 35 | | sequence alignment file. |
| 36 | constrain_mb | Add constraint commands to MrBayes input file based on a |
| 37 | | given tree topology. |
| 38 | visualize_tree | Parse phylogenetic analysis output files for a single |
| 39 | | alignment in a given directory, and write human-readable |
| 40 | | tree figures to PDF files. |
| 41 | replace_seqs | Replace sequences in an alignment with their top hits in a |
| 42 | | given fasta file (useful if genomes or taxon selection has |
| 43 | | been updated). |
| 44 | | |
| 45 | Miscellaneous commands: | |
| 46 | csv_to_fasta | Generate a fasta file from sequences detailed in a |
| 47 | | spreadsheet of similarity search results. |
| 48 | check_depend | Check that all the dependencies are properly installed and |
| 49 | | useable. |

```

1      check_imports      Check that all the import statements used in the AMOEBAE
2                          repository run without error.
3      regen_genome_info Write a new genome info spreadsheet file using filenames
4                          from the Genomes directory.
5
6  positional arguments:
7      command      Specify one of the functionalities of amoebae.
8
9  optional arguments:
10     -h, --help  show this help message and exit
11
12  Copyright 2018 Lael D. Barlow Licensed under the Apache License, Version 2.0
13  (the "License"); you may not use this file except in compliance with the
14  License. You may obtain a copy of the License at
15  http://www.apache.org/licenses/LICENSE-2.0 Unless required by applicable law
16  or agreed to in writing, software distributed under the License is distributed
17  on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either
18  express or implied. See the License for the specific language governing
19  permissions and limitations under the License.

```

20 3.2 amoebae mkdatadir

```

21  usage: amoebae [-h] new_dir_path
22
23  Make a directory with subdirectories and CSV files for storing sequence data,
24  etc.
25
26  positional arguments:
27      new_dir_path  Specify the full file path that you want the new directory to
28                    have.
29
30  optional arguments:
31      -h, --help  show this help message and exit

```

32 3.3 amoebae setup_hmmdb

```

33  usage: amoebae [-h] indirpath
34
35  Construct an HMM database (with hmmpress). This is for later sorting of given
36  sequences into categories based on which HMM the score highest against.
37
38  positional arguments:
39      indirpath  Path to directory containing amino acid sequence alignment
40                 file(s) to be constructed into an HMM database using hmmpress
41                 from the HMMer3 software package.
42
43  optional arguments:
44      -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

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.

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]

Print a list of all usable query files in the query directory as defined in
the settings file.

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

3.6 amoebae add_to_queries

```
usage: amoebae [-h] query_file

Add a query file to a formatted directory. This command adds a given sequence
file to the directory with the path that you have specified in the settings.py
file, and appends a corresponding line to the CSV file that you specified
(e.g., '0_query_info.csv') to indicate the query title, etc.

positional arguments:
```

```

1  query_file Path to a sequence file in FASTA format that can be used as a
2             similarity search query file. Or path to a directory containing
3             only files for addition to the queries. Note: By default, the
4             portion of the input filename preceding the first underscore
5             character will be recorded as the "query title", the remaining
6             substring preceding the second underscore character will be
7             recorded as the taxon (e.g., "Hsapiens"), and the rest of the
8             filename preceding the filename extension will be recorded as
9             the sequence ID. So the filename might look like this:
10            "QUERYTITLE_HSAPIENS_SEQUENCEID.fa". However, the relevant
11            information can be revised in the "Queries/0_query_info.csv"
12            file afterward if necessary.
13
14 optional arguments:
15  -h, --help show this help message and exit

```

16 3.7 amoebae list_queries

```

17 usage: amoebae [-h]
18
19 Print a list of all usable query files in the query directory as defined in
20 the settings file.
21
22 optional arguments:
23  -h, --help show this help message and exit

```

24 3.8 amoebae get_redun_hits

```

25 usage: amoebae [-h] [--csv_file CSV_FILE] [--query_name QUERY_NAME]
26                [--query_list_file QUERY_LIST_FILE] [--db_name DB_NAME]
27                [--db_list_file DB_LIST_FILE] [--query_title QUERY_TITLE]
28                [--outdir OUTDIR]
29                [--blast_report_evalue_cutoff BLAST_REPORT_EVALUATE_CUTOFF]
30                [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
31                [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUATE_CUTOFF]
32                [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
33                [--max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE]
34                [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
35                ]
36                [--predict_redun_hit_selection]
37                srch_dir
38
39 Run searches with queries to find redundant hits in databases (for
40 interpreting results).
41
42 positional arguments:
43  srch_dir Path to directory that will contain output directory
44           as a subdirectory.

```

```

1
2 optional arguments:
3   -h, --help          show this help message and exit
4   --csv_file CSV_FILE  Path to spreadsheet to append summary of result to for
5                        manual annotation. (default: None)
6   --query_name QUERY_NAME
7                        Query filename to use (not full path). (default: None)
8   --query_list_file QUERY_LIST_FILE
9                        Path to file containing a list of query files to use,
10                       if no query_name is specified (or all queries by
11                       default). (default: None)
12   --db_name DB_NAME    Name of database file in the database directory in
13                       which to do searches (not full path). (default: None)
14   --db_list_file DB_LIST_FILE
15                       Path to file containing a list of database files to
16                       use (if no db_name specified). (default: None)
17   --query_title QUERY_TITLE
18                       Name to be assigned to hits in databases that may be
19                       considered redundant with a search query to which the
20                       same title is assigned, otherwise it is taken from the
21                       query info spreadsheet specified in the settings.py
22                       file ('query_info_csv'). (default: None)
23   --outdir OUTDIR      Path to directory to write search results to.
24                       (default: None)
25   --blast_report_evalue_cutoff BLAST_REPORT_EVALUATE_CUTOFF
26                       Maximum E-value for reporting BLAST hits. (default:
27                       0.05)
28   --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
29                       Maximum BLAST target sequences to consider. (default:
30                       500)
31   --hmmer_report_evalue_cutoff HMMER_REPORT_EVALUATE_CUTOFF
32                       Maximum E-value for reporting HMMer hits. (default:
33                       0.05)
34   --hmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF
35                       Minimum sequence score for reporting HMMer hits.
36                       (default: 5)
37   --max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE
38                       Absolute maximum number of hits (BLAST, HMMer, etc) to
39                       summarize in the output spreadsheet. This is important
40                       when working with sequences with WD40 domains, for
41                       example. (default: 50)
42   --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
43                       Number of threads to use for running searches.
44                       (default: 4)
45   --predict_redun_hit_selection
46                       Write a copy of the output spreadsheet with '+' in
47                       rows for hits that may be specific to each query
48                       title, due to not being retrieved as top hits by
49                       queries associated with different query titles.

```



```

1             (default: False)
2
3 Recommendation: For most analyses, use the --query_name option and the
4 --db_name option, and run the get_redun_hits command for each query
5 separately. Otherwise, there will be redundant information in the output
6 spreadsheet(s).

```

3.9 amoebae setup_fwd_srch

```

8 usage: amoebae [-h] [--outdir OUTDIR] srch_dir query_list_file db_list_file
9
10 Make a directory in which to write output files from similarity searches.
11
12 positional arguments:
13   srch_dir             Path to directory that will contain output directory as a
14                       subdirectory.
15   query_list_file      Path to file with list of queries to search with.
16   db_list_file         Path to file with list of databases to search with.
17
18 optional arguments:
19   -h, --help           show this help message and exit
20   --outdir OUTDIR      Path to directory to put search results into (so that this
21                       step can be piped together with other commands). (default:
22                       None)
23
24 Note: Use the bash script to run forward searches on a remote server.

```

3.10 amoebae run_fwd_srch

```

26 usage: amoebae [-h] [--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF]
27                [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
28                [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF]
29                [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
30                [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
31                ]
32                fwd_srch_dir
33
34 Perform searches with original queries into subject databases.
35
36 positional arguments:
37   fwd_srch_dir         Path to directory that will contain forward search
38                       output files.
39
40 optional arguments:
41   -h, --help           show this help message and exit
42   --blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF
43                       Maximum E-value for reporting BLAST hits. (default:
44                       0.05)

```

```

1  --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
2      Maximum BLAST target sequences to consider. (default:
3      500)
4  --hmmreport_evalue_cutoff HMMER_REPORT_EVALUATE_CUTOFF
5      Maximum E-value for reporting HMMer hits. (default:
6      0.05)
7  --hmmreport_score_cutoff HMMER_REPORT_SCORE_CUTOFF
8      Minimum sequence score for reporting HMMer hits.
9      (default: 5)
10 --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
11     Number of threads to use for running searches.
12     (default: 4)

```

13 3.11 amoebae sum_fwd_srch

```

14 usage: amoebae [-h] [--max_evalue MAX_EVALUATE]
15               [--max_gap_between_tblastn_hsp MAX_GAP_BETWEEN_TBLASTN_HSPS]
16               [--do_not_use_exonerate]
17               [--exonerate_score_threshold EXONERATE_SCORE_THRESHOLD]
18               [--max_hits_to_sum MAX_HITS_TO_SUM]
19               [--max_length_diff MAX_LENGTH_DIFF]
20               fwd_srch_out csv_file
21

```

22 Append information about forward searches to csv summary file (this is used to
23 organize reverse searches). For TBLASTN searches (protein queries, nucleotide
24 target sequences), HSPs are clustered into groups that are close enough within
25 the target sequence to potentially represent exons from the same coding
26 sequence. The nucleotide subsequences in which these clusters of HSPs are
27 found are then analyzed using exonerate to identify and translate potential
28 exons, in "protein2genome" mode, because exonerate, unlike TBLASTN, attempts
29 to identify exon boundaries, yielding translations that are less likely to
30 include translations of non-coding regions outside exons (which might include
31 apparent stop codons).

32 positional arguments:

```

34   fwd_srch_out      Path to directory where forward search results were
35                     written.
36   csv_file          Path to summary spreadsheet (CSV) file, which may
37                     already contain search summaries, or may not exist
38                     yet.
39

```

40 optional arguments:

```

41   -h, --help        show this help message and exit
42   --max_evalue MAX_EVALUATE
43                     Maximum E-value threshold for reporting forward search
44                     hits. (default: 0.0005)
45   --max_gap_between_tblastn_hsp MAX_GAP_BETWEEN_TBLASTN_HSPS
46                     Maximum number of nucleotide bases between TBLASTN
47                     HSPs to be considered part of the same gene locus.

```

```

1          This is important, because it will be assumed that HSP
2          separated by more than this number of nucleotide bases
3          are not part of the same gene or TBLASTN "hit".
4          (default: 10000)
5  --do_not_use_exonerate
6          Override the default use of exonerate to identify
7          coding sequences and translations, and just use
8          TBLASTN instead. This option is provided because
9          concatenated TBLASTN HSPs may be more inclusive of
10         sequences within the target sequence, and the results
11         of TBLASTN and exonerate may need to be compared.
12         Also, note that HSPs identified by TBLASTN but for
13         which exonerate yields no alignments will be ignored
14         if exonerate is used. (default: False)
15  --exonerate_score_threshold EXONERATE_SCORE_THRESHOLD
16         Set score threshold to be applied when running
17         exonerate on nucleotide sequences identified by
18         TBLASTN. The default for setting of exonerate is 100,
19         but a lower score is set as default here, because
20         otherwise exonerate cannot identify some of the
21         sequences identified by TBLASTN. This option is only
22         relevant if using exonerate. (default: 10)
23  --max_hits_to_sum MAX_HITS_TO_SUM
24         Maximum number of forward search hits to list in the
25         summary spreadsheet. If zero, then reverse searches
26         will be performed for all hits. (default: 0)
27  --max_length_diff MAX_LENGTH_DIFF
28         Maximum number of amino acid residues length
29         difference allowed between the original query and the
30         forward hit sequence. If -1, then a maximum length
31         cutoff will not be applied. (default: -1)

```

3.12 amoebae setup_rev_srch

```

33 usage: amoebae [-h] [--outdir OUTDIR] [--aasubseq] [--nafullseq]
34                srch_dir csv_file databases
35
36 Make directory in which to write results of reverse searches.
37
38 positional arguments:
39   srch_dir              Path to directory that will contain output directory as a
40                        subdirectory.
41   csv_file              Path to summary spreadsheet (CSV) file, which contains a
42                        summary of forward search(es).
43   databases             Database filename (in database directory) or path to file
44                        with list of database filenames. Note that filenames are
45                        needed, not file paths.
46
47 optional arguments:

```

```

1  -h, --help          show this help message and exit
2  --outdir OUTDIR    Path to directory to put search results into (so that this
3                     step can be piped together with other commands). (default:
4                     None)
5  --aasubseq         Use only the portion of each (amino acid) forward hit
6                     sequence that aligns to the original query used (top HSP
7                     subject sequence). This is default for nucleotide hits.
8                     (default: False)
9  --nafullseq        Use the full (nucleic acid) forward hit sequence. This is
10                     default for amino acid hits. (default: False)

```

11 3.13 amoebae run_rev_srch

```

12 usage: amoebae [-h] [--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF]
13                [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
14                [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF]
15                [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
16                [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
17                ]
18                rev_srch_dir
19
20 Perform searches with forward search hit sequences as queries into the
21 original query databases.
22
23 positional arguments:
24   rev_srch_dir          Path to directory that will contain output of
25                         searches.
26
27 optional arguments:
28   -h, --help            show this help message and exit
29   --blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF
30                         Maximum E-value for reporting BLAST hits. (default:
31                         0.05)
32   --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
33                         Maximum BLAST target sequences to consider. (default:
34                         500)
35   --hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF
36                         Maximum E-value for reporting HMMer hits. (default:
37                         0.05)
38   --hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF
39                         Minimum sequence score for reporting HMMer hits.
40                         (default: 5)
41   --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
42                         Number of threads to use for running searches.
43                         (default: 4)

```

44 3.14 amoebae sum_rev_srch

```

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

```

3.15 amoebae interp_srchs

```

44  usage: amoebae [-h] [--fwd_only] [--fwd_evalue_cutoff FWD_EVALUATE_CUTOFF]
45                [--rev_evalue_cutoff REV_EVALUATE_CUTOFF]
46                [--hmmmer_cutoff HMMER_CUTOFF] [--no_overlapping_hits]
47                [--out_csv_path OUT_CSV_PATH]

```

```

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

```

40 3.16 amoebae find_redun_seqs

```

41 usage: amoebae [-h] [--out_csv_path OUT_CSV_PATH]
42                [--remove_tblastn_hits_at_annotated_loci]
43                [--just_look_for_genes_in_gff3] [--ignore_gff3]
44                [--allow_internal_stops ALLOW_INTERNAL_STOPS]
45                [--min_length MIN_LENGTH]
46                [--min_percent_length MIN_PERCENT_LENGTH]
47                [--min_percent_query_cover MIN_PERCENT_QUERY_COVER]

```

```

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

```

```

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

```



```

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

```

file in the queries directory to use for comparing sequences. Aligned FASTA queries are selected that match the query titles of the original queries used to retrieve each of the relevant hits listed in the csv file. No other options need to be specified in this case. (default: False)

3.17 amoebae plot

```
usage: amoebae [-h] [--csv_file2 CSV_FILE2] [--complex_info COMPLEX_INFO]
               [--row_order ROW_ORDER] [--out_pdf OUT_PDF]
               csv_file
```

Plot results of similarity search and sequence classification analyses. The outputs are PDF files.

positional arguments:

csv_file Path to a spreadsheet with the relevant results to be plotted. This can be either a CSV file output of the sum_rev_srch command or from the find_redun_seqs command. If the output of the sum_rev_srch command is used, however, redundant hits will be counted (e.g., BLASTP and TBLASTN hits corresponding to the same or highly identical genomic loci).

optional arguments:

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

--csv_file2 CSV_FILE2 Path to a second spreadsheet with relevant results to be compared to the first and plotted. (default: None)

--complex_info COMPLEX_INFO Path to file that specifies which query titles represent components of which protein complexes (or otherwise grouped proteins). (default: None)

--row_order ROW_ORDER Path to file that specifies the order in which data for each species will be displayed. (default: None)

--out_pdf OUT_PDF Path to output pdf file. (default: None)

3.18 amoebae add_to_models

```
usage: amoebae [-h]
               model_name alignment tree_topology subs_model type_seqs taxon
```

Add a phylogenetic model for relationships between members of a gene family (sequence_data matrix, data type, tree topology, type sequence defining each clade of interest, and substitution model) to a directory for use in classifying sequence (via the 'phylo_class' command).

```

1
2 positional arguments:
3   model_name      An arbitrary name for the model (which will refer to the
4                   alignment, tree, substitution model, etc. collectively).
5   alignment       A multiple amino acid sequence alignment in nexus format.
6   tree_topology   Text file containing a tree (identified previously using
7                   MrBayes, etc) containing the names of all the sequences in
8                   the alignment, in newick format.
9   subs_model      The name of the substitution model used to recover the
10                   provided topology (chosen with ModelFinder or similar
11                   software).
12   type_seqs       Names of sequences (sequence headers) that are to be used to
13                   define clades of interest. A csv file with seq names in one
14                   column and clade names in the next column.
15   taxon           Taxonomic group represented in the model (e.g., "Eukaryotes",
16                   or "Amorphea").
17
18 optional arguments:
19   -h, --help      show this help message and exit

```

20 3.19 amoebae list_models

```

21 usage: amoebae [-h]
22
23 Print a list of all usable model/reference tree names in the models directory
24 as defined in the settings file.
25
26 optional arguments:
27   -h, --help      show this help message and exit

```

28 3.20 amoebae get_alt_topos

```

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

```

```

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

```

3.21 amoebae prune

```

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

```

3.22 amoebae auto_prune

```

37 usage: amoebae [-h]
38                [--max_bl_iqr_above_third_quartile MAX_BL_IQR_ABOVE_THIRD_QUARTILE]
39                [--remove_redun_seqs REMOVE_REDUN_SEQS]
40                [--remove_redun_seqs_threshold REMOVE_REDUN_SEQS_THRESHOLD]
41                [--output_file OUTPUT_FILE]
42                in_dir
43
44 Automatically identify sequences in a tree, and remove them from a given

```

```

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

```

3.23 amoebae reduce_tree

```

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

```

3.24 amoebae constrain_mb

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

Add constraint commands to MrBayes input file.

positional arguments:
  alignment      Nexus alignment for input to MrBayes (without any
                  constraint commands).
  tree           Tree in newick format with same taxon names as in
                  alignment. To be used as a topology constraint (all
                  nodes).

optional arguments:
  -h, --help      show this help message and exit
  --out_alignment OUT_ALIGNMENT
                  Path to nexus alignment for input to MrBayes with
                  constraints added. (default: None)
```

3.25 amoebae visualize_tree

```
usage: amoebae [-h] [--root_taxon ROOT_TAXON] [--highlight_paralogues]
               [--add_clade_names_from_file]
               input_directory method

Parse phylogenetic analysis output files in a given directory, and write
human-readable tree figures to PDF files.

positional arguments:
  input_directory  Path to directory containing input files (must contain
                   a .table file for decoding taxon names.
  method          Name of tree searching program used. Either iqtree,
                   raxml, or mrbayes accepted.

optional arguments:
  -h, --help      show this help message and exit
  --root_taxon ROOT_TAXON
                   Name of species to root on (e.g.,
                   "Klebsormidium_nitens").
  --highlight_paralogues
                   Highlight clades that contain paralogues found in at
                   least one other clade in the tree.
  --add_clade_names_from_file
                   Use a file in the parent directory with clade names
                   corresponding to representative sequences to add clade
                   names to all the taxon names in the output trees.
```

3.26 amoebae replace_seqs

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

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

positional arguments:

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

optional arguments:

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

--fasta_file FASTA_FILE

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

3.27 amoebae csv_to_fasta

usage: amoebae [-h] [--output_dir OUTPUT_DIR] [--abbrev] [--paralogue_names] [--only_descr] [--subseq] [--all_hits] [--split_by_query_title] [--split_by_top_rev_srch_hit SPLIT_BY_TOP_REV_SRCH_HIT] [--split_to_query_fastas] csv_file

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

positional arguments:

csv_file Path to csv file listing sequences.

optional arguments:

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

--output_dir OUTPUT_DIR

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

--abbrev Add species name instead of sequence description from fasta header. Applicable when the output file is to be used for alignment and phylogenetic analysis. (default: False)

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

```

1  --only_descr          Use the description but not the ID as the new fasta
2                        sequence header. Does not work if the abbrev option is
3                        specified. (default: False)
4  --subseq              Write subsequences that aligned to forward search
5                        query, instead of the full sequences. (default: False)
6  --all_hits            Write all forward hits listed in the input csv file.
7                        (default: False)
8  --split_by_query_title
9                        Write sequences to files according to the query title
10                       of the query which retrieved them in a similarity
11                       search. (default: False)
12  --split_by_top_rev_srch_hit SPLIT_BY_TOP_REV_SRCH_HIT
13                       Write sequences to files according to the top hit that
14                       they retrieve in a reverse search, for each sequence
15                       that meets the reverse search criteria. (Provide the
16                       reverse search identifier, eg,
17                       "rev_srch_20180924122402-1") (default: None)
18  --split_to_query_fastas
19                       Write sequences to separate files with filenames that
20                       can be easily parsed for loading the the files as
21                       queries using the add_to_queries command. (default:
22                       False)

```

23 3.28 amoebae check_depend

```

24 usage: amoebae [-h]
25
26 Check that all the dependencies (other than python modules) are properly
27 installed and useable.
28
29 optional arguments:
30  -h, --help  show this help message and exit

```

31 3.29 amoebae check_imports

```

32 usage: amoebae [-h]
33
34 Check that all the import statements used in the AMOEBAE repository run
35 without error.
36
37 optional arguments:
38  -h, --help  show this help message and exit

```

39 3.30 amoebae regen_genome_info

```

40 usage: amoebae [-h] data_dir_path
41

```



```
1 Write a new genome info spreadsheet (O_genome_info.csv) file using filenames
2 from the Genomes directory.
3
4 positional arguments:
5   data_dir_path  Specify the full path to an existing AMOEBAE data directory,
6                   which contains a 'Genomes' subdirectory. The new genome info
7                   file will be added to this subdirectory.
8
9 optional arguments:
10  -h, --help      show this help message and exit
```

11 4 Miscellaneous scripts

12 Several scripts of less general applicability than the amoebae commands described above
13 are included in the AMOEBAE toolkit. See the amoebae/misc_scripts directory (https://github.com/laelbarlow/amoebae/tree/master/misc_scripts). Most scripts have in-
14 formation regarding usage in the files themselves. More detailed information regarding some
15 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.