

AMOEBAE documentation

Lael D. Barlow

Version of June 19, 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	3
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	9
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	12
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	21
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	25
3.25	amoebae visualize_tree	26
3.26	amoebae replace_seqs	26
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	28

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) provide a means to investigate the evolution of one or a few genes via similarity searching, and automated pipelines such as orthoMCL (Li, 2003) attempt to rapidly perform orthology prediction for all genes in several genomes. AMOEBAE addresses mid-scale analyses which are too cumbersome to be done via webservices and yet require a level of detail and flexibility not offered by automated 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. However, 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 (30 queries and 100 genomes may in fact be unmanageable). 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 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

1 against information in genome annotation files, and systematic use of Exonerate (Slater and
2 Birney, 2005) where possible for obtaining better exon predictions.

3 1.4 A word of caution

4 AMOEBAE is not optimized for ease of use, but is meant to be highly configurable. The
5 many options available to AMOEBAE users inevitably provide many opportunities for errors
6 in specifying search criteria, and errors in interpreting output files. Some prior experience
7 with similarity searching and with running software using the command line is essential
8 for using AMOEBAE, and experience writing scripts in Bash and Python would be highly
9 advantageous. Moreover, AMOEBAE is still under active development, so some features may
10 not yet be thoroughly tested.

11 1.5 User support

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

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

16 1.6 Documentation

17 This document provides an overview of AMOEBAE and describes the functionality of the var-
18 ious commands/scripts. For a tutorial which includes a working example of a similarity search
19 analysis run using AMOEBAE, see the Jupyter Notebook: `amoebae/notebooks/similar-`
20 `ity_search_tutorial_2.ipynb`. For code documentation, please see the html file(s), which can
21 be opened with your web browser: `amoebae/doc/code_documentation/html/index.html`.

22 1.7 How to cite AMOEBAE

23 Please cite the GitHub webpage <https://github.com/laelbarlow/amoebae> (or alternative
24 permanent repositories if relevant). Also, the first publication to make use of a version of
25 AMOEBAE was an analysis of Adaptor Protein subunits in embryophytes by Larson *et al.*
26 (2019).

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

1.8 Acknowledgments

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

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

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

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

1.9 License

Copyright 2018 Lael D. Barlow

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

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

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

All dependencies are free and open-source, and can be automatically installed in a virtual environment (see subsection 2.3).

1 These are the main dependencies of AMOEBAE:

- 2 • Python3 (the Anaconda distribution works well).
- 3 • Biopython, a Python package for bioinformatics (Cock *et al.*, 2009).
- 4 • The Environment for Tree Exploration 3 (ETE3), a Python package for working with
5 phylogenetic trees (Huerta-Cepas *et al.*, 2016).
- 6 • Matplotlib, a Python package for generating plots (Hunter, 2007).
- 7 • (gffutils).
- 8 • NCBI BLAST+, a software package for sequence similarity searching (Camacho *et al.*,
9 2009).
- 10 • HMMer3, a software package for profile sequence similarity searching (Eddy, 1998).
- 11 • MUSCLE, for multiple sequence alignment (Edgar, 2004).
- 12 • IQ-TREE, for phylogenetic analysis (Nguyen *et al.*, 2015).

13 2.3 Setting up an environment for AMOEBAE using Singularity

14 Follow the steps below to set up AMOEBAE on your personal computer. This setup process
15 should take approximately 20 minutes to complete. Additional instructions for setting up
16 AMOEBAE on a remote server will soon be added as well.

- 17 1. Ensure that Git is installed on your computer. If you do not already have git installed,
18 then your computer will prompt you with instructions for how to install it when you
19 type git into the command line. If you have a newer version of macOS it may prompt
20 you to install developer tools, which may take up a considerable amount of storage
21 space. Documentation for Git is available here: <https://git-scm.com/doc>. You can
22 check which version you have (or whether it is installed at all) by running the command
23 below. Please note: Here ">>>" is used to indicate that the following text in the line
24 is to be entered in you terminal command prompt.

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

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

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


3. 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
```

4. If you are setting up AMOEBAE on a high performance computing cluster, then you will not be able to install Singularity yourself, and may need to use specific procedures to load Singularity prior to use.

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 singularity in a virtual machine (*e.g.*, on MacOS), then 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).

```
1 >>> amoebae check_depend
```

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

7 3 Command reference

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

11 3.1 amoebae

```
12 usage: amoebae <command> [<args>]
```

```
13  
14 Commands for setting up data structure:
```

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

```
17  
18 Commands for similarity searching:
```

```
19     setup_hmddb     Construct an HMM database (with hmmpress).  
20     add_to_dbs      Format and add a file to a formatted directory.  
21     list_dbs        Print a list of all usable database files in the database  
22                   directory as defined in the settings file.  
23     add_to_queries  Add a query file to a formatted directory.  
24     list_queries    Print a list of all usable query files in the query  
25                   directory as defined in the settings file.  
26     get_redun_hits  Run searches with queries to find redundant hits in  
27                   databases (for interpreting results).  
28     setup_fwd_srch  Make directory in which to perform forward searches.  
29     run_fwd_srch    Perform searches with given queries into given dbs.  
30     sum_fwd_srch    Append information about forward searches to csv summary  
31                   file (this is used to organize reverse searches).  
32     setup_rev_srch  Make a directory in which to perform reverse searches.  
33     run_rev_srch    Perform searches with given forward search hits into given db.  
34     sum_rev_srch    Append information about reverse searches to csv summary  
35                   file.  
36     interp_srchs    Interpret search results based on summary.  
37     find_redun_seqs Identify sequences likely encoded on redundant loci  
38                   predicted for the same species.  
39     plot            Plot search results.
```

```

1
2 Commands for phylogenetic analysis using a reference tree:
3     add_to_models      Add an alignment, tree, substitution model, names of
4                         clade-defining sequences to a directory with other models.
5     list_models        Print a list of all usable model/reference tree names in
6                         the models directory as defined in the settings file.
7     get_alt_topos      Take a tree and make copies with every alternative
8                         topology for the branches connecting the clades of
9                         interest.
10
11 Commands for phylogenetic analysis without a reference tree:
12     prune              Identify sequences in a tree, and remove them from a
13                         given alignment for further phylogenetic analysis.
14     auto_prune          Automatically identify sequences in a tree, and remove
15                         them from a given alignment for further phylogenetic
16                         analysis.
17     reduce_tree        Remove terminal nodes from a given tree if there are
18                         not any sequences with the same name in a given multiple
19                         sequence alignment file.
20     constrain_mb        Add constraint commands to MrBayes input file based on a
21                         given tree topology.
22     visualize_tree      Parse phylogenetic analysis output files for a single
23                         alignment in a given directory, and write human-readable
24                         tree figures to PDF files.
25     replace_seqs        Replace sequences in an alignment with their top hits in a
26                         given fasta file (useful if genomes or taxon selection has
27                         been updated).
28
29 Miscellaneous commands:
30     csv_to_fasta        Generate a fasta file from sequences detailed in a
31                         spreadsheet of similarity search results.
32     check_depend        Check that all the dependencies are properly installed and
33                         useable.
34     check_imports       Check that all the import statements used in the AMOEBAE
35                         repository run without error.
36     regen_genome_info    Write a new genome info spreadsheet file using filenames
37                         from the Genomes directory.
38
39 positional arguments:
40     command             Specify one of the functionalities of amoebae.
41
42 optional arguments:
43     -h, --help          show this help message and exit
44
45 Copyright 2018 Lael D. Barlow Licensed under the Apache License, Version 2.0
46 (the "License"); you may not use this file except in compliance with the
47 License. You may obtain a copy of the License at
48 http://www.apache.org/licenses/LICENSE-2.0 Unless required by applicable law
49 or agreed to in writing, software distributed under the License is distributed

```

1 on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either
2 express or implied. See the License for the specific language governing
3 permissions and limitations under the License.

4 3.2 amoebae mkdatadir

5 usage: amoebae [-h] new_dir_path
6
7 Make a directory with subdirectories and CSV files for storing sequence data,
8 etc.
9
10 positional arguments:
11 new_dir_path Specify the full file path that you want the new directory to
12 have.
13
14 optional arguments:
15 -h, --help show this help message and exit

16 3.3 amoebae setup_hmmdb

17 usage: amoebae [-h] indirpath
18
19 Construct an HMM database (with hmmpress). This is for later sorting of given
20 sequences into categories based on which HMM the score highest against.
21
22 positional arguments:
23 indirpath Path to directory containing amino acid sequence alignment
24 file(s) to be constructed into an HMM database using hmmpress
25 from the HMMer3 software package.
26
27 optional arguments:
28 -h, --help show this help message and exit

29 3.4 amoebae add_to_dbs

30 usage: amoebae [-h] [--split_char SPLIT_CHAR] [--split_pos SPLIT_POS]
31 [--skip_header_reformat] [--auto_extract_accs]
32 new_file
33
34 Format and add a file to a formatted directory.
35
36 positional arguments:
37 new_file Can be a fasta file (prot or nucl) or HMM databases,
38 generated using the hmmpress program in the HMMer
39 software package. Or a GFF3 annotation file.
40
41 optional arguments:

```

1  -h, --help          show this help message and exit
2  --split_char SPLIT_CHAR
3                      Character to split the header string on for extracting
4                      the accession. (default: )
5  --split_pos SPLIT_POS
6                      Position that the accession will be in after
7                      splitting. (default: 0)
8  --skip_header_reformat
9                      Skip reformatting of header lines in input fasta file.
10                     (default: False)
11  --auto_extract_accs Automatically identify accessions/IDs in sequence
12                     headers (overrides split_char and split_pos options
13                     above). (default: False)

```

14 3.5 amoebae list_db

```

15 usage: amoebae [-h]
16
17 Print a list of all usable query files in the query directory as defined in
18 the settings file.
19
20 optional arguments:
21  -h, --help  show this help message and exit

```

22 3.6 amoebae add_to_queries

```

23 usage: amoebae [-h] query_file
24
25 Add a query file to a formatted directory. This command adds a given sequence
26 file to the directory with the path that you have specified in the settings.py
27 file, and appends a corresponding line to the CSV file that you specified
28 (e.g., '0_query_info.csv') to indicate the query title, etc.
29
30 positional arguments:
31  query_file  Path to a sequence file in FASTA format that can be used as a
32              similarity search query file. Or path to a directory containing
33              only files for addition to the queries. Note: By default, the
34              portion of the input filename preceding the first underscore
35              character will be recorded as the "query title", the remaining
36              substring preceding the second underscore character will be
37              recorded as the taxon (e.g., "Hsapiens"), and the rest of the
38              filename preceding the filename extension will be recorded as
39              the sequence ID. So the filename might look like this:
40              "QUERYTITLE_HSAPIENS_SEQUENCEID.fa". However, the relevant
41              information can be revised in the "Queries/0_query_info.csv"
42              file afterward if necessary.
43
44 optional arguments:

```

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

2 3.7 amoebae list_queries

3 usage: amoebae [-h]

4

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

7

8 optional arguments:

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

10 3.8 amoebae get_redun_hits

11 usage: amoebae [-h] [--csv_file CSV_FILE] [--query_name QUERY_NAME]
12 [--query_list_file QUERY_LIST_FILE] [--db_name DB_NAME]
13 [--db_list_file DB_LIST_FILE] [--query_title QUERY_TITLE]
14 [--outdir OUTDIR]
15 [--blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF]
16 [--blast_max_target_seqs BLAST_MAX_TARGET_SEQS]
17 [--hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF]
18 [--hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF]
19 [--max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE]
20 [--num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING]
21]
22 [--predict_redun_hit_selection]
23 srch_dir

24

25 Run searches with queries to find redundant hits in databases (for
26 interpreting results).

27

28 positional arguments:

29 srch_dir Path to directory that will contain output directory
30 as a subdirectory.

31

32 optional arguments:

33 -h, --help show this help message and exit
34 --csv_file CSV_FILE Path to spreadsheet to append summary of result to for
35 manual annotation. (default: None)
36 --query_name QUERY_NAME Query filename to use (not full path). (default: None)
37 if no query_name is specified (or all queries by
38 --query_list_file QUERY_LIST_FILE Path to file containing a list of query files to use,
39 if no query_name is specified (or all queries by
40 default). (default: None)
41 --db_name DB_NAME Name of database file in the database directory in
42 which to do searches (not full path). (default: None)
43 --db_list_file DB_LIST_FILE

44

```

1          Path to file containing a list of database files to
2          use (if no db_name specified). (default: None)
3  --query_title QUERY_TITLE
4          Name to be assigned to hits in databases that may be
5          considered redundant with a search query to which the
6          same title is assigned, otherwise it is taken from the
7          query info spreadsheet specified in the settings.py
8          file ('query_info_csv'). (default: None)
9  --outdir OUTDIR
10         Path to directory to write search results to.
11         (default: None)
12  --blast_report_evalue_cutoff BLAST_REPORT_EVALUE_CUTOFF
13         Maximum E-value for reporting BLAST hits. (default:
14         0.05)
15  --blast_max_target_seqs BLAST_MAX_TARGET_SEQS
16         Maximum BLAST target sequences to consider. (default:
17         500)
18  --hmmmer_report_evalue_cutoff HMMER_REPORT_EVALUE_CUTOFF
19         Maximum E-value for reporting HMMer hits. (default:
20         0.05)
21  --hmmmer_report_score_cutoff HMMER_REPORT_SCORE_CUTOFF
22         Minimum sequence score for reporting HMMer hits.
23         (default: 5)
24  --max_number_of_hits_to_summarize MAX_NUMBER_OF_HITS_TO_SUMMARIZE
25         Absolute maximum number of hits (BLAST, HMMer, etc) to
26         summarize in the output spreadsheet. This is important
27         when working with sequences with WD40 domains, for
28         example. (default: 50)
29  --num_threads_similarity_searching NUM_THREADS_SIMILARITY_SEARCHING
30         Number of threads to use for running searches.
31         (default: 4)
32  --predict_redun_hit_selection
33         Write a copy of the output spreadsheet with '+' in
34         rows for hits that may be specific to each query
35         title, due to not being retrieved as top hits by
36         queries associated with different query titles.
37         (default: False)
38
39  Recommendation: For most analyses, use the --query_name option and the
40  --db_name option, and run the get_redun_hits command for each query
41  separately. Otherwise, there will be redundant information in the output
42  spreadsheet(s).

```

42 3.9 amoebae setup_fwd_srch

```

43 usage: amoebae [-h] [--outdir OUTDIR] srch_dir query_list_file db_list_file
44
45 Make a directory in which to write output files from similarity searches.
46
47 positional arguments:

```



```

1  srch_dir          Path to directory that will contain output directory as a
2                    subdirectory.
3  query_list_file   Path to file with list of queries to search with.
4  db_list_file      Path to file with list of databases to search with.
5
6  optional arguments:
7  -h, --help        show this help message and exit
8  --outdir OUTDIR   Path to directory to put search results into (so that this
9                    step can be piped together with other commands). (default:
10                     None)
11
12  Note: Use the bash script to run forward searches on a remote server.

```

13 3.10 amoebae run_fwd_srch

```

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

```

3.11 amoebae sum_fwd_srch

```
usage: amoebae [-h] [--max_evalue MAX_EVALUE]
               [--max_gap_between_tblastn_hspes MAX_GAP_BETWEEN_TBLASTN_HSPS]
               [--do_not_use_exonerate]
               [--exonerate_score_threshold EXONERATE_SCORE_THRESHOLD]
               [--max_hits_to_sum MAX_HITS_TO_SUM]
               [--max_length_diff MAX_LENGTH_DIFF]
               fwd_srch_out csv_file
```

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

positional arguments:

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

optional arguments:

-h, --help	show this help message and exit
--max_evalue MAX_EVALUE	Maximum E-value threshold for reporting forward search hits. (default: 0.0005)
--max_gap_between_tblastn_hspes MAX_GAP_BETWEEN_TBLASTN_HSPS	Maximum number of nucleotide bases between TBLASTN HSPs to be considered part of the same gene locus. This is important, because it will be assumed that HSP separated by more than this number of nucleotide bases are not part of the same gene or TBLASTN "hit". (default: 10000)
--do_not_use_exonerate	Override the default use of exonerate to identify coding sequences and translations, and just use TBLASTN instead. This option is provided because concatenated TBLASTN HSPs may be more inclusive of sequences within the target sequence, and the results of TBLASTN and exonerate may need to be compared. Also, note that HSPs identified by TBLASTN but for which exonerate yields no alignments will be ignored

```

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

```

19 3.12 amoebae setup_rev_srch

```

20 usage: amoebae [-h] [--outdir OUTDIR] [--aasubseq] [--nafullseq]
21             srch_dir csv_file databases
22
23 Make directory in which to write results of reverse searches.
24
25 positional arguments:
26  srch_dir              Path to directory that will contain output directory as a
27                        subdirectory.
28  csv_file              Path to summary spreadsheet (CSV) file, which contains a
29                        summary of forward search(es).
30  databases             Database filename (in database directory) or path to file
31                        with list of database filenames. Note that filenames are
32                        needed, not file paths.
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
```

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.

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]
               csv_file rev_srch_out
```

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

positional arguments:

csv_file Path to summary spreadsheet (CSV) file, which may already contain reverse search summaries.

rev_srch_out Path to directory where reverse search results were

```

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

```

3.15 amoebae interp_srchs

```

32 usage: amoebae [-h] [--fwd_only] [--fwd_evalue_cutoff FWD_EVALUE_CUTOFF]
33               [--rev_evalue_cutoff REV_EVALUE_CUTOFF]
34               [--hmmmer_cutoff HMMER_CUTOFF] [--no_overlapping_hits]
35               [--out_csv_path OUT_CSV_PATH]
36               csv_file
37
38 Interpret search results based on final summary, which provides a basis for
39 further analyses of positive hits.
40
41 positional arguments:
42   csv_file              Path to spreadsheet with forward and reverse search
43                       results.
44
45 optional arguments:
46   -h, --help          show this help message and exit
47   --fwd_only          Interpret forward searches based on score (HMMer)

```

```

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

```

3.16 amoebae find_redun_seqs

```

29 usage: amoebae [-h] [--out_csv_path OUT_CSV_PATH]
30                [--remove_tblastn_hits_at_annotated_loci]
31                [--just_look_for_genes_in_gff3] [--ignore_gff3]
32                [--allow_internal_stops ALLOW_INTERNAL_STOPS]
33                [--min_length MIN_LENGTH]
34                [--min_percent_length MIN_PERCENT_LENGTH]
35                [--min_percent_query_cover MIN_PERCENT_QUERY_COVER]
36                [--overlap_required] [--max_percent_ident MAX_PERCENT_IDENT]
37                [--min_alig_res_in_overlap MIN_ALIG_RES_IN_OVERLAP]
38                [--min_ident_res_in_overlap MIN_IDENT_RES_IN_OVERLAP]
39                [--min_sim_res_in_overlap MIN_SIM_RES_IN_OVERLAP]
40                [--min_ident_span_len MIN_IDENT_SPAN_LEN]
41                [--min_sim_span_len MIN_SIM_SPAN_LEN]
42                [--min_percent_ident_in_overlap MIN_PERCENT_IDENT_IN_OVERLAP]
43                [--min_percent_sim_in_overlap MIN_PERCENT_SIM_IN_OVERLAP]
44                [--min_percent_overlap MIN_PERCENT_OVERLAP]
45                [--plot_hit_exclusion] [--add_ali_col]
46                csv_file
47

```

1 Identify hit sequences likely encoded by the same gene loci in the genome of a
2 given species, or otherwise not representing paralogous genes. Criteria are
3 applied in this order: 1. Peptide hits with the same ID as higher-ranking hits
4 for the same query (query title) are excluded. 2. Nucleotide hits for the same
5 loci as peptide sequence hits are excluded. 3. Sequences with internal stop
6 codons are excluded, as these are potentially pseudogenes. 4. Sequences are
7 excluded if they do not meet several minimum length criteria: Absolute minimum
8 length (in amino acids) and percent query cover. 5. Sequences are excluded if
9 they do not overlap to a specified degree with all included higher-ranking
10 hits for the same query (query title) in sequence data for the same
11 species/genome. This is determined by algorithmically comparing pairs of
12 sequences aligned to a reference alignment of homologues, and several minimum
13 measures of alignment overlap may be specified. 6. Secondary hit sequences are
14 excluded if they do not meet a specified maximum percent identity threshold.
15 Highly identical sequences may result from false segmental duplications in the
16 genome assembly, may represent alleles, etc. Note: Applying these criteria
17 requires a column to be manually added to the input csv file prior to running
18 with the header "Alignment for sequence comparison" and filled with the
19 appropriate alignment name to use (one for each query title, as listed in the
20 "Query title" column). Alternatively, you can run this command with the
21 --add_ali_col option to automatically identify appropriate alignments among
22 your aligned FASTA queries used for running HMMer searches. If no alignment
23 (.afaa) file can be found, then the first single sequence query file (.faa)
24 that appears in the summary CSV file will be used instead.

25
26 positional arguments:

27 csv_file Path to spreadsheet with interpreted search results
28 outputted by the interp_srchs command.

29
30 optional arguments:

31 -h, --help show this help message and exit
32 --out_csv_path OUT_CSV_PATH
33 Optionally specify an output file path, so that this
34 command can be piped together with others. (default:
35 None)
36 --remove_tblastn_hits_at_annotated_loci
37 Ignore tblastn hits that overlap with any previously
38 annotated loci. The rationale for this would be that
39 the corresponding protein sequences should have been
40 retrieved if the tblastn hit were a true positive
41 anyway. If this option is not specified, then
42 sequences will still be excluded if they specifically
43 correspond to the same loci as do higher-ranking hits.
44 (default: False)
45 --just_look_for_genes_in_gff3
46 When looking for records in GFF3 annotation files that
47 overlap with subsequences identified by similarity
48 searching (TBLASTN), ignore records that are not
49 explicitly "gene" (for example, "CDS", "mRNA", and

```

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

```



```

1         the overlap_required option is specified. (default:
2         15)
3 --min_ident_span_len MIN_IDENT_SPAN_LEN
4         Minimum number of aligning residues which are
5         identical that must exist in at least one continuous
6         span for two sequences to be considered as potentially
7         distinct hits (not counting positions where both
8         sequences have gaps). This is only relevant if the
9         overlap_required option is specified. (default: 0)
10 --min_sim_span_len MIN_SIM_SPAN_LEN
11         Minimum number of aligning residues which are similar
12         (or identical) that must exist in at least one
13         continuous span for two sequences to be considered as
14         potentially distinct hits (not counting positions
15         where both sequences have gaps). This is only relevant
16         if the overlap_required option is specified. (default:
17         0)
18 --min_percent_ident_in_overlap MIN_PERCENT_IDENT_IN_OVERLAP
19         Minimum percent identity between the two sequences of
20         interest in the alignment. This is only relevant if the
21         overlap_required option is specified. (default: 0)
22 --min_percent_sim_in_overlap MIN_PERCENT_SIM_IN_OVERLAP
23         Minimum percent similarity (including identity)
24         between the two sequences of interest in the
25         alignment. This is only relevant if the
26         overlap_required option is specified. (default: 0)
27 --min_percent_overlap MIN_PERCENT_OVERLAP
28         Minimum number of aligning residues between the two
29         sequences of interest as a percentage of the length of
30         the second sequence (the last sequence in the
31         alignment), not including gaps, for the two sequences
32         to be considered as potentially distinct hits. This is
33         only relevant if the overlap_required option is
34         specified. (default: 0)
35 --plot_hit_exclusion Plot number of hits excluded by the various criteria
36         applied. (default: False)
37 --add_ali_col Add a column to the csv file listing which alignment
38         file in the queries directory to use for comparing
39         sequences. Aligned FASTA queries are selected that
40         match the query titles of the original queries used to
41         retrieve each of the relevant hits listed in the csv
42         file. No other options need to be specified in this
43         case. (default: False)

```

44 3.17 amoebae plot

```

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

```

```

1
2 Plot results of similarity search and sequence classification analyses. The
3 outputs are PDF files.
4
5 positional arguments:
6     csv_file           Path to a spreadsheet with the relevant results to be
7                         plotted. This can be either a CSV file output of the
8                         sum_rev_srch command or from the find_redun_seqs
9                         command. If the output of the sum_rev_srch command is
10                        used, however, redundant hits will be counted (e.g.,
11                        BLASTP and TBLASTN hits corresponding to the same or
12                        highly identical genomic loci).
13
14 optional arguments:
15     -h, --help         show this help message and exit
16     --csv_file2 CSV_FILE2
17                         Path to a second spreadsheet with relevant results to
18                         be compared to the first and plotted. (default: None)
19     --complex_info COMPLEX_INFO
20                         Path to file that specifies which query titles
21                         represent components of which protein complexes (or
22                         otherwise grouped proteins). (default: None)
23     --row_order ROW_ORDER
24                         Path to file that specifies the order in which data
25                         for each species will be displayed. (default: None)
26     --out_pdf OUT_PDF   Path to output pdf file. (default: None)

```

27 3.18 amoebae add _to _models

```

28 usage: amoebae [-h]
29                model_name alignment tree_topology subs_model type_seqs taxon
30
31 Add a phylogenetic model for relationships between members of a gene family
32 (sequence_data matrix, data type, tree topology, type sequence defining each
33 clade of interest, and substitution model) to a directory for use in
34 classifying sequence (via the 'phylo_class' command).
35
36 positional arguments:
37     model_name         An arbitrary name for the model (which will refer to the
38                         alignment, tree, substitution model, etc. collectively).
39     alignment          A multiple amino acid sequence alignment in nexus format.
40     tree_topology      Text file containing a tree (identified previously using
41                         MrBayes, etc) containing the names of all the sequences in
42                         the alignment, in newick format.
43     subs_model         The name of the substitution model used to recover the
44                         provided topology (chosen with ModelFinder or similar
45                         software).
46     type_seqs          Names of sequences (sequence headers) that are to be used to
47                         define clades of interest. A csv file with seq names in one

```

```

1         column and clade names in the next column.
2     taxon         Taxonomic group represented in the model (e.g., "Eukaryotes",
3                   or "Amorphea").
4
5 optional arguments:
6     -h, --help     show this help message and exit

```

7 3.19 amoebae list_models

```

8 usage: amoebae [-h]
9
10 Print a list of all usable model/reference tree names in the models directory
11 as defined in the settings file.
12
13 optional arguments:
14     -h, --help     show this help message and exit

```

15 3.20 amoebae get_alt_topos

```

16 usage: amoebae [-h] [--polytomy] [--not_polytomy_clades]
17                [--keep_original_backbone] [--iqtree_au_test]
18                model_name out_dir_path
19
20 Take a tree and make copies with every alternative topology for the branches
21 connecting the clades of interest. Output as additional models in the Models
22 directory.
23
24 positional arguments:
25     model_name         Name of model/backbone tree to modify (other info
26                       provided in the model info csv file).
27     out_dir_path       Path to directory in which output directory will be
28                       written.
29
30 optional arguments:
31     -h, --help         show this help message and exit
32     --polytomy         Just make one big polytomy connecting the clades of
33                       interest instead of making alternative bifurcating
34                       trees. (default: False)
35     --not_polytomy_clades
36                       Do not make subtrees/clades of interest polytomies in
37                       output topologies. (default: False)
38     --keep_original_backbone
39                       Keep the original backbone topology instead of
40                       generating a polytomy or alternative resolved
41                       topologies. (default: False)
42     --iqtree_au_test   Test all the relevant alternative topologies against
43                       each other using Approximately Unbiased (AU) test with
44                       IQ-tree. (default: False)

```

3.21 amoebae prune

```
usage: amoebae [-h] [--include_seqs] [--output_file OUTPUT_FILE]
              tree_file alignment name_replace_table

Identify sequences in a tree, and remove them from a given alignment for
further phylogenetic analysis.

positional arguments:
  tree_file              Tree in newick format (coded names, because ETE3
                        cannot parse taxon names with space characters without
                        quotation marks around them).
  alignment              Dataset used to make the tree (nexus alignment)
                        (original alignment with original taxon names either
                        trimmed or untrimmed).
  name_replace_table     File for decoding names in input tree file.

optional arguments:
  -h, --help            show this help message and exit
  --include_seqs         Include only listed sequences/nodes instead of
                        removing them. (default: False)
  --output_file OUTPUT_FILE
                        Path to output file. (default: None)
```

3.22 amoebae auto_prune

```
usage: amoebae [-h]
              [--max_bl_iqr_above_third_quartile MAX_BL_IQR_ABOVE_THIRD_QUARTILE]
              [--remove_redun_seqs REMOVE_REDUN_SEQS]
              [--remove_redun_seqs_threshold REMOVE_REDUN_SEQS_THRESHOLD]
              [--output_file OUTPUT_FILE]
              in_dir

Automatically identify sequences in a tree, and remove them from a given
alignment for further phylogenetic analysis.

positional arguments:
  in_dir                Path to directory that contains the phylogenetic
                        analysis output files (sequence name conversion table
                        file and original nexus alignment file can be in the
                        parent directory to this directory as long as their
                        names are mostly identical.

optional arguments:
  -h, --help            show this help message and exit
  --max_bl_iqr_above_third_quartile MAX_BL_IQR_ABOVE_THIRD_QUARTILE
                        Inclusion threshold for number of interquartile ranges
                        above the third quartile of terminal branch lengths
                        the length of a terminal branch can be before it is
```

```

1             considered an outlier (length is total distance from
2             root node after rooting on midpoint, or the longest
3             terminal branch on either side of the midpoint).
4             (default: 1.5)
5  --remove_redun_seqs REMOVE_REDUN_SEQS
6             Remove taxonomically redundant sequences (longest
7             branch of two sister branches when both are sequences
8             from the same species. (default: True)
9  --remove_redun_seqs_threshold REMOVE_REDUN_SEQS_THRESHOLD
10            Minimum support required to consider one of two sister
11            branches/sequences taxonomically redundant. Note: only
12            used if the remove_redun_seqs option is specified.
13            (default: 0.95)
14  --output_file OUTPUT_FILE
15            Path to output file. (default: None)

```

16 3.23 amoebae reduce_tree

```

17 usage: amoebae [-h] [--output_file OUTPUT_FILE] alignment tree_file
18
19 Remove terminal nodes from a given tree if there are not any sequences with
20 the same name in a given alignment.
21
22 positional arguments:
23   alignment            Alignment in nexus format with sequences representing
24                       a subset of those represented in the input tree.
25   tree_file            Tree in newick format.
26
27 optional arguments:
28   -h, --help            show this help message and exit
29   --output_file OUTPUT_FILE
30                       Path to output file. (default: None)

```

31 3.24 amoebae constrain_mb

```

32 usage: amoebae [-h] [--out_alignment OUT_ALIGNMENT] alignment tree
33
34 Add constraint commands to MrBayes input file.
35
36 positional arguments:
37   alignment            Nexus alignment for input to MrBayes (without any
38                       constraint commands).
39   tree                 Tree in newick format with same taxon names as in
40                       alignment. To be used as a topology constraint (all
41                       nodes).
42
43 optional arguments:
44   -h, --help            show this help message and exit

```

```

1  --out_alignment OUT_ALIGNMENT
2          Path to nexus alignment for input to Mrbayes with
3          constraints added. (default: None)

4  3.25 amoebae visualize_tree

5  usage: amoebae [-h] [--root_taxon ROOT_TAXON] [--highlight_paralogues]
6          [--add_clade_names_from_file]
7          input_directory method
8
9  Parse phylogenetic analysis output files in a given directory, and write
10 human-readable tree figures to PDF files.
11
12 positional arguments:
13   input_directory      Path to directory containing input files (must contain
14                        a .table file for decoding taxon names.
15   method               Name of tree searching program used. Either iqtree,
16                        raxml, or mrbayes accepted.
17
18 optional arguments:
19   -h, --help           show this help message and exit
20   --root_taxon ROOT_TAXON
21                        Name of species to root on (e.g.,
22                        "Klebsormidium_nitens").
23   --highlight_paralogues
24                        Highlight clades that contain paralogues found in at
25                        least one other clade in the tree.
26   --add_clade_names_from_file
27                        Use a file in the parent directory with clade names
28                        corresponding to representative sequences to add clade
29                        names to all the taxon names in the output trees.

```

30 **3.26 amoebae replace_seqs**

```

31 usage: amoebae [-h] [--fasta_file FASTA_FILE] alignment
32
33 Replace sequences in an alignment the full-length sequences from the relevant
34 file(s) in the Genomes directory, or with their top hits in a given fasta
35 file. And, align, mask, and trim the identified sequences to the input
36 alignment
37
38 positional arguments:
39   alignment            Path to multiple sequence alignment file in nexus
40                        format (trimmed alignment).
41
42 optional arguments:
43   -h, --help           show this help message and exit
44   --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] [--parologue_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)

--parologue_names Use species name, query title, and parologue 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)

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

--subseq Write subsequences that aligned to forward search query, instead of the full sequences. (default: False)

--all_hits Write all forward hits listed in the input csv file.
(default: False)

--split_by_query_title Write sequences to files according to the query title of the query which retrieved them in a similarity search. (default: False)

--split_by_top_rev_srch_hit SPLIT_BY_TOP_REV_SRCH_HIT

Write sequences to files according to the top hit that they retrieve in a reverse search, for each sequence that meets the reverse search criteria. (Provide the reverse search identifier, eg,

```

1             "rev_srch_20180924122402-1") (default: None)
2  --split_to_query_fastas
3             Write sequences to separate files with filenames that
4             can be easily parsed for loading the the files as
5             queries using the add_to_queries command. (default:
6             False)

```

7 3.28 amoebae check_depend

```

8  usage: amoebae [-h]
9
10 Check that all the dependencies (other than python modules) are properly
11 installed and useable.
12
13 optional arguments:
14  -h, --help  show this help message and exit

```

15 3.29 amoebae check_imports

```

16 usage: amoebae [-h]
17
18 Check that all the import statements used in the AMOEBAE repository run
19 without error.
20
21 optional arguments:
22  -h, --help  show this help message and exit

```

23 3.30 amoebae regen_genome_info

```

24 usage: amoebae [-h] data_dir_path
25
26 Write a new genome info spreadsheet (0_genome_info.csv) file using filenames
27 from the Genomes directory.
28
29 positional arguments:
30  data_dir_path  Specify the full path to an existing AMOEBAE data directory,
31                 which contains a 'Genomes' subdirectory. The new genome info
32                 file will be added to this subdirectory.
33
34 optional arguments:
35  -h, --help      show this help message and exit

```

36 4 Miscellaneous scripts

37 Several scripts of less general applicability than the amoebae commands described above
38 are included in the AMOEBAE toolkit. See the amoebae/misc_scripts directory (https://github.com/amoebae/amoebae/tree/master/misc_scripts)

1 [//github.com/laelbarlow/amoebae/tree/master/misc_scripts](https://github.com/laelbarlow/amoebae/tree/master/misc_scripts)). Most scripts have in-
2 formation regarding usage in the files themselves. More detailed information regarding some
3 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.
- 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.