

# **CLTree 1.0**

## **User's Manual**

Guanghong Zuo  
ghzuo@ucas.ac.cn

March 24, 2022

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The Installation and Testing</b>	<b>1</b>
2.1	Normal Unix-like Mode . . . . .	1
2.1.1	Preparation . . . . .	1
2.1.2	Compile by CMake . . . . .	1
2.1.3	Testing with Example . . . . .	1
2.2	Run Collapse in Container . . . . .	2
<b>3</b>	<b>Workflow of CLTree</b>	<b>3</b>
3.1	Basic Workflow and Options . . . . .	3
3.2	Lineage Information Preparation . . . . .	4
3.3	Advance Lineage Tools . . . . .	4
<b>4</b>	<b>Command Usage</b>	<b>5</b>
4.1	Basic Usage . . . . .	5
4.2	Default Task . . . . .	6
4.3	Lineage Tasks . . . . .	6
4.4	NCBI Taxonomy Database Tasks . . . . .	7
<b>5</b>	<b>Algorithm</b>	<b>8</b>
5.1	Annotate the Phylogeny Tree . . . . .	8
5.2	Measure Difference by Shannon Entropy . . . . .	8
<b>6</b>	<b>Citing Collapse in a Publication</b>	<b>8</b>
	<b>Reference</b>	<b>8</b>

# 1 Introduction

CLTree is a tool to annotate the phylogenetic tree by lineage and measure their differences in topology by Shannon entropy. This manual is for the Text Version of the CLTree.

## 2 The Installation and Testing

CLTree is distributed by source code. It can be downloaded from Internet (<https://github.com/ghzuo/Collapse>). There are two ways to compile the source codes of CLTree: compiling the source code by CMake; or using the docker image.

### 2.1 Normal Unix-like Mode

The program is implemented in the C++ language. The following build tools and libraries are required.

#### 2.1.1 Preparation

- `cmake`  $\geq 3.0$
- `g++`  $\geq 7.0$  or other compiler supporting C++11 standard
- require library: `libz`, `nlohmann-json`

#### 2.1.2 Compile by CMake

1. unzip the package file and change into it
2. `mkdir build` and change into it
3. `cmake ..` and some options you wanted
4. `make`
5. `make install` (*option*)

#### 2.1.3 Testing with Example

If this is the first time you use the CLTree package, please go to the “example” folder. Please run the `cltree` command to get an annotated phylogenetic tree and monophyly status by:

```
../build/bin/cltree
```

More detail of the command usage can be obtain by ‘-h’ option or read the follow sections.

## 2.2 Run Collapse in Container

The docker containers make the programs can be performed on both Windows and Linux/-MacOS, and transfer the programs easily. To employ the container with Collapse, you should install docker at first. You can download docker free and reference from <https://docs.docker.com/> to how to install it. After installing docker, basic usages for CLTree in the container are shown below:

1. Obtain image: You can build the Collapse docker image based on Dockerfile in the source code by command

```
docker build -t="cltree-img" .
```

Here option ‘-t’ set the image name. After build image, you can delete the dangling images for build by `docker image prune`. This will save much hard disk space. You can also download prebuilt Collapse image from internet by command:

```
docker pull ghzu0/cltree .
```

In this step, an image with Collapse programs will obtained.

2. Start container from image: run the follow command in the Collapse directory, i.e. the directory which include the ‘example’ directory of the Collapse

```
docker run --rm -it -v $PWD/example:/root/data cltree-img
```

In this step, you will enter the Collapse container, and the “example” folder of this project will be find in the “data” folder. Change path to the data folder, and run

```
cltree
```

You will get the result for eight genomes in the `list` file. You can change the path ‘\$PWD/example’ to your own data directory.

3. Exit and stop container: `exit` in docker terminal.
4. Run Collapse in a temporary container by one command without enter the container:

```
cd <example> or <other data folder>
docker run --rm -v $PWD:/root/data/ cltree-img cltree
```

5. More usage for docker can reference <https://docs.docker.com/>.

### 3 Workflow of CLTree

#### 3.1 Basic Workflow and Options

The scheme of program 'cltree' is shown in Figure 1. The main command is "cltree". User selects different tokens (the read block in Figure 1) for different tasks. Before running the program, users must prepare two objects:

- A phylogenetic tree (in newick format, e.g. Tree.nwk).
- Lineage information for the leaves of the phylogenetic tree.

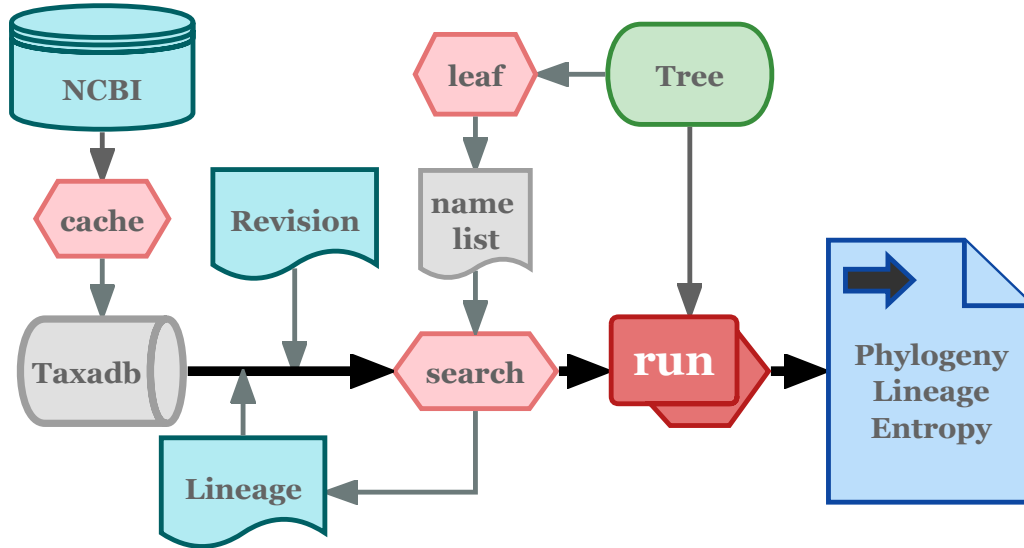


Figure 1: The Basic Workflow of CLTree. The red blocks are the token of cltree, and other blocks indicate the input/output files.

The easiest way for preparing lineage information is utilizing the NCBI Taxonomy database dump file package directly. The file can be downloaded from the NCBI website (<https://ftp.ncbi.nlm.nih.gov/pub/taxonomy/taxdump.tar.gz>). With NCBI taxonomy dump file package and Tree.nwk, users can obtain the result by the command:

```
cltree -i Tree.nwk -d taxdump.tar.gz
```

In default mode, this command will output four files included annotated tree, lineage statistics, and Shannon entropy between phylogeny and lineage (see Table 1)

And to speedup, users can make the cache for the NCBI taxonomy database package by using the command:

Table 1: The default output files of cltree

file name	description
collapsed.entropy	the Shannon entropy for every taxon rank
collapsed-annotated.nwk	a newick tree with every node annotated with the common lineage for all sub branches
collapsed.lineage	lineages for all genomes in CSV format
collapsed.unit	statistics for all taxon ranks in the rank order

```
cltree cache -d taxdump.tar.gz
```

The command will generate a cache file for the database dump, named “taxadb.gz”. It will speedup the program in future tasks.

## 3.2 Lineage Information Preparation

There are two formats for recording the lineage information, i.e. CSV and custom format. The simple way is using the CSV format. It records one lineage information of a genome per line, with the taxon rank at the header (see the lineage.csv in the example folder of source code).

## 3.3 Advance Lineage Tools

There are two columns in this file. The first column is the name of the species, and the lineage is shown in the second column

```
superkingdom D
kingdom K
subkingdom k
superphylum Q
phylum P
subphylum p
superclass L
class C
subclass c
superorder W
order O
suborder o
family F
subfamily f
```

genus	G
subgenus	g
species	S
subspecies	s

- Example of lineage file:

```
Chlamydia_pecorum...RefSeq <D>Bacteria<P>Chlamydiae...<S>Chlamydia_pecorum
Chlamydia_trachomatis...RefSeq <D>Bacteria<P>Chlamydiae...<S>Chlamydia_trachomatis
Chlamydia_pecorum...RefSeq <D>Bacteria<P>Chlamydiae...<S>Chlamydia_pecorum
Chlamydia_trachomatis...RefSeq <D>Bacteria<P>Chlamydiae...<S>Chlamydia_trachomatis
Corynebacterium_jeikeium... <D>Bacteria<P>Actinobacteria...<S>Corynebacterium_jei
Leptospira_interrogans...RefSeq <D>Bacteria<P>Spirochaetes...<S>Leptospira_interro
```

- Example for abbreviation of taxon ranks

#### A Example for batch revision

```
# Phylum B13 Firmicutes
<F>Erysipelotrichaceae<G>Eubacterium <F>Erysipelotrichaceae<G>Erysipelothrix # CVTree
<P>Firmicutes<C>Erysipelotrichia <P>Tenericutes<C>Erysipelotrichia # CVTree
<F>Unclassified<G>Exiguobacterium <F>Bacillaceae<G>Exiguobacterium # CVTree
<F>Planococcaceae<G>Solibacillus <F>Bacillaceae<G>Solibacillus # NCBI taxonomy -> LE
```

## 4 Command Usage

### 4.1 Basic Usage

To obtain a basic usage of the program and tasks, you can use the “help” token, as `cltree help` or `cltree -h`. The output of the help information is shown in below. And users can type `cltree help Task` to get more information about other tasks.

`cltree Task [options]`

Available Task:

run	Annotate phylogenetic tree with lineage system
cache	Make NCBI database cache from taxdump.tar.gz
query	Query lineage from local NCBI taxonomy database
search	Search lineage from lineage files and NCBI taxonomy database, and revised by the revision file
leaf	Obtain the species name list of phylogenetic tree
rank	Output the default rank names and abbreviations
help	Provide the help information for <Task>
[ -h ]	Display this information

## 4.2 Default Task

`cltree run` perform the main task annotating phylogenetic tree by lineage and measuring their Differences by Shannon entropy. Due to this is the default task of the command, the token `run` can be omitted except for the “-h” option.

```
cltree run
[ -D ./ ]           The work directory, default: ./
[ -i Tree.nwk ]     Input newick tree, default: Tree.nwk
[ -o collapsed ]    Output prefix name: default: collapsed
[ -m <Revision.txt> ] Lineage revision file for batch edit,
                    default: None
[ -l Lineage.txt ]   Lineage file for leafs of tree,
                    default: Lineage.txt or Lineage.csv
[ -d taxadb.gz ]     Taxa database file or directory,
                    default: taxadb.gz or taxdump.tar.gz
[ -R <None> ]        List of rank names and abbreviations,
                    default: use the setting of program
[ -r DKPCOFGS ]      Abbreviation of output taxon rank,
                    default: according to source
[ -O <Outgroup> ]    Set the outgroup for the unroot tree.
                    default: None, rearranged by taxonomy
[ -P ]              Output prediction for undefined leafs
[ -q ]              Run command in quiet mode
[ -h ]              Display this information
```

## 4.3 Lineage Tasks

Users can also obtain the leaf lineages of the phylogenetic tree by other task step by step and review them manually, i.e. using `cltree leaf` obtain the name list of phylogenetic tree, and using `cltree search` obtain the lineages of the species. The function and usage of these two programs are shown below:

- `leaf` – Obtain the name list of phylogenetic tree

```
cltree leaf
[ -i Tree.nwk ]      Input tree file, default: Tree.nwk
[ -o namelist.txt ]  Output name list, default: name.list
```



```
[ -q ]      Run command in quiet mode
[ -h ]      Display this information
```

- **search** – Search the lineage of the genome

```
cltree search
```

```
[ -i namelist.txt ]      Input name list, ':N' after the file name
                          select the N column of the file
                          default: first column of namelist.txt

[ -o Lineage.csv ]       Output lineage file, default: lineage.csv

[ -m <Revision.txt> ]    Lineage revise file for batch edit,
                          default: None

[ -l Lineage.txt ]       Lineage file for leafs of tree,
                          default: Lineage.txt

[ -d taxadb.gz ]         Taxa database file or directory,
                          default: taxadb.gz or taxdump.tar.gz

[ -R <None> ]            List file for rank names and abbreviations,
                          default: use the setting of program

[ -r <DKPCOFGS> ]       Set output taxon rank by abbreviations,
                          default: according to source

[ -q ]                  Run command in quiet mode
[ -h ]                  Display this information
```

## 4.4 NCBI Taxonomy Database Tasks

Other tasks of the program are used to handle the NCBI Taxonomy database.

- **cache** – Package the dump files of NCBI taxonomy database as cache to speed up query lineage from database.

```
cltree cache
```

```
[ -d taxdump.tar.gz ]    NCBI taxon dumpfile directory, default: ta
[ -o taxadb.gz ]         Packaged taxon database, default: taxadb.g
[ -q ]                  Run command in quiet mode
[ -h ]                  Display this information
```

- **rank** – Output an example for the ordered list of the taxon rank names and abbreviations.

```

cltree rank
[ -o ranklist.txt ] Output name list, default: ranklist.txt
[ -q ]             Run command in quiet mode
[ -h ]             Display this information

```

- query – Query the lineage of the genome ONLY from NCBI Taxonomy database dump files or the database cache.

```

cltree query
[ -I <Taxon ID> ]      Query a taxon id
[ -N <Taxon Name> ]    Query a taxon name
[ -i namelist.txt ]     The query list file default: name.list
[ -d taxadb.gz ]        The dump of NCBI taxonomy database
[ -o Lineage.txt ]      Output file, default: Lineage.txt
[ -R <None> ]          List file for rank names and abbreviations,
                        default: use the setting of program
[ -r <DKPCOFGS> ]      Set output taxon rank by abbreviations,
                        default: same to the source
[ -H ]                 Don't output missing items
[ -q ]                 Run command in quiet mode
[ -h ]                 Display this information

```

## 5 Algorithm

### 5.1 Annotate the Phylogeny Tree

### 5.2 Measure Difference by Shannon Entropy

## 6 Citing Collapse in a Publication

Please cite:

1. Guanghong Zuo (2021) Collapse: Annotate Phylogenetic Tree by Lineage and Measure their Differences by Shannon Entropy. in preparation.

## References

Zuo, G. (2021) Collapse: Annotate Phylogenetic Tree by Lineage and Measure their Differences by Shannon Entropy. *Journal*, in submission.