# Project title: Subfamily classification and analysis of CAZy family

**Goal: Classify the CAZy family into subfamilies, analyze and visualize them using Bioinformatics tools and methods.**

In this tutorial, CAZy family GH31 was used to explain bioinformatics methods to analyze and visualize them.

**Paper Link:** A subfamily classification to choreograph the diverse activities within glycoside hydrolase family 31. **DOI:** https://doi.org/10.1016/j.jbc.2023.103038

**Required materials are available at https://github.com/sivanr92/Class_project**

**Updated: Xinpeng, 9/19/2023**

**Updated: Siva, 02 Oct 2023, 3.20 AM.**

**Updated: Xinpeng 02 Oct.**

**Methods to identify the subfamily from a family:**

1. Dataset and preprocess
2. Domain annotation using HMM and dbCAN
3. Extraction of modules (based on the annotation)
4. Construction of sequence similarity networks (SSN) using SSNpipe and analyzing SSNs based on characterized IDs from CAZy and EC numbers.
5. Visualization of SSN networks using Cytoscape
6. Phylogenentic Analysis
7. Interpretation and Discussion

**Detailed tutorials and methods are available in the following pages.**

Note: Denotes the program name available at the GitHub.

**Detailed methods section:**

**1. Dataset and preprocess:**
- Dataset can be downloaded from CAZy database
  http://www.cazy.org/
- Select GH classes on CAZy database (Your family under the classification, we will use GH31 as example, please choose your own family)
  http://www.cazy.org/Glycoside-Hydrolases.html
- Go to GH class 31 and download the dataset.
  http://www.cazy.org/GH31.html
- Find the page you can find the dataset
  http://www.cazy.org/IMG/cazy_data/GH31.txt
  Download that using the wget in the terminal as use click and download.
  Command line:

```
$ wget http://www.cazy.org/IMG/cazy_data/GH31.txt
```



Use GH31.txt to get a list of genbank ids of GH31 family.

Use the following command in the terminal make to make the unique list of IDs

```
$ cut -f 4 GH31.txt > list_ids.txt
```

Remove redundant ids from the list

```
$ cat list_ids.txt | sort | uniq > uniq_list_ids.txt
```

24307 unique ids were unique.

To download the sequences from the NCBI. Use the batch service available from NCBI Entrez Direct: E-utilities

Full tutorial Link: https://www.ncbi.nlm.nih.gov/books/NBK179288/

- Register in NCBI using the for using the API-KEY. Helps for fast and terminal download.
- NCBI register using UNL mail id. (Use Institution search)
- Navigate to account settings.
- Find the Generate API-KEY button and generate.

My NCBI    My Bibliography    Account Settings    Site Preferences

MyNCBI Dashboard > NCBI Account Settings

## NCBI Account Settings

### Email

This email is used for delivery of saved searches and recovery of password for your native NCBI account.

| Email | Status | Edit |
|---|---|---|
| snallattinputhurra2@unl.edu | (confirmed) | ✏ |

### NCBI Account

Your username is the email address of the third-party account that you used to set up your NCBI account.

| Username |
|---|
| snallattinputhurra2@unl.edu |

### Linked Accounts

You can log into your NCBI account via these third parties. Contact the third party about any issues related to logging into any of the accounts below.

| Account | Email/ID | Remove |
|---|---|---|
| University of Nebraska-Lincoln | snallattinputhurra2@unl.edu (logged in) | 🗑 |

Add account

Install EDirect using one of the following.

```
$ sh -c "$(curl -fsSL
https://ftp.ncbi.nlm.nih.gov/entrez/entrezdirect/install-
edirect.sh)"
```

**OR**

```
$ sh -c "$(wget -q
https://ftp.ncbi.nlm.nih.gov/entrez/entrezdirect/install-
edirect.sh -O -)"
```

```
$ export PATH=${HOME}/edirect:${PATH}. (After installation, an
automatic path will be shown. Export that)
```

After that export your unique API-KEY, using the following command,

```
$ export NCBI_API_KEY=unique_api_key
```

**Program to download the sequences can be found at** <mark>"batch_download_protein_seqeunces.sh"</mark>

24306 ids were found and downloaded.

## 2. Domain annotation using HMM and dbCAN.
### (i) download dbCAN-fam-HMMs.txt, hmmscan-parser.sh

Download following files from https://bcb.unl.edu/dbCAN2/download/Databases/dbCAN-old@UGA/

- hmmscan-parser.sh
- dbCAN-fam-HMMs.txt



### (ii) download HMMER 3.0 package [hmmer.org] and install it properly

http://hmmer.org/download.html, File: http://eddylab.org/software/hmmer/hmmer-3.4.tar.gz

HMMER is available in HCC.

- Can be loaded with the following commands.

  ```
  $ ml hmmer/3.3
  ```

- Check the installation by running:

```
$ hmmbuild -h
```

```
[sivanr@login1.swan gh31_project]$ hmmbuild -h
# hmmbuild :: profile HMM construction from multiple sequence alignments
# HMMER 3.4 (Aug 2023); http://hmmer.org/
# Copyright (C) 2023 Howard Hughes Medical Institute.
# Freely distributed under the BSD open source license.
# - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Usage: hmmbuild [-options] <hmmfile_out> <msafile>

Basic options:
  -h     : show brief help on version and usage
  -n <s> : name the HMM <s>
  -o <f> : direct summary output to file <f>, not stdout
  -O <f> : resave annotated, possibly modified MSA to file <f>

Options for selecting alphabet rather than guessing it:
  --amino : input alignment is protein sequence data
  --dna   : input alignment is DNA sequence data
  --rna   : input alignment is RNA sequence data

Alternative model construction strategies:
  --fast          : assign cols w/ >= symfrac residues as consensus  [default]
  --hand          : manual construction (requires reference annotation)
  --symfrac <x>   : sets sym fraction controlling --fast construction  [0.5]
  --fragthresh <x> : if L <= x*alen, tag sequence as a fragment  [0.5]

Alternative relative sequence weighting strategies:
  --wpb     : Henikoff position-based weights  [default]
  --wgsc    : Gerstein/Sonnhammer/Chothia tree weights
  --wblosum : Henikoff simple filter weights
  --wnone   : don't do any relative weighting; set all to 1
  --wgiven  : use weights as given in MSA file
  --wid <x> : for --wblosum: set identity cutoff  [0.62]  (0<=x<=1)
```

**(iii) format HMM db: hmmpress dbCAN-fam-HMMs.txt**

> **$** `hmmpress dbCAN-fam-HMMs.txt`

```
[sivanr@login1.swan gh31_project]$ hmmpress dbCAN-fam-HMMs.txt
Working...    done.
Pressed and indexed 783 HMMs (783 names and 9 accessions).
Models pressed into binary file:   dbCAN-fam-HMMs.txt.h3m
SSI index for binary model file:   dbCAN-fam-HMMs.txt.h3i
Profiles (MSV part) pressed into:  dbCAN-fam-HMMs.txt.h3f
Profiles (remainder) pressed into: dbCAN-fam-HMMs.txt.h3p
[sivanr@login1.swan gh31_project]$
```

**(iv) run: hmmscan --domtblout yourfile.out.dm dbCAN-fam-HMMs.txt yourfile > yourfile.out**

**$** `hmmscan --domtblout gh31_hmmscan.out.dm dbCAN-fam-HMMs.txt batch_500_protein_sequences.fasta > gh31_hmmscan.out`

Use the code below to run hmmscan on HCC

Program available at "hmmscan_run.slurm"

Students can use batch and guest partition.

```
#!/bin/bash
#SBATCH --time=24:00:00
#SBATCH --mem=100gb
#SBATCH --job-name=hmmscan
#SBATCH --error=/Your/path/gh31_project/job.%J.err
#SBATCH --output=/ Your/path/gh31_project/job.%J.out
#SBATCH --partition=batch,guest
ml hmmer/3.3
hmmscan --domtblout gh31_hmmscan.out.dm dbCAN-fam-HMMs.txt
batch_500_protein_sequences.fasta > gh31_hmmscan.out
```

**(v) run: sh hmmscan-parser.sh yourfile.out.dm > yourfile.out.dm.ps (if alignment > 80aa, use E-value < 1e-5, otherwise use E-value < 1e-3; covered fraction of HMM > 0.3)**

```
$ sh hmmscan-parser.sh gh31_hmmscan.out.dm >
gh31_hmmscan.out.dm.ps
```



**26390 modules were found.**

**(vi) run: cat yourfile.out.dm.ps | awk '$5<1e-15&&$10>0.35' > yourfile.out.dm.ps.stringent (this allows you to get the same result as what is produced in our dbCAN2 webpage)**

```
$
```

After filtering using the E- Value and coverage cutoff of <1e-15 and >0.35, **25279 modules were found**.

Cols in yourfile.out.dm.ps:
1. Family HMM
2. HMM length
3. Query ID
4. Query length
5. E-value (how similar to the family HMM)
6. HMM start
7. HMM end
8. Query start
9. Query end
10. Coverage
** On our dbCAN2 website, we use E-value < 1e-15 and coverage > 0.35, which is more stringent than the default ones in hmmscan-parser.sh

3. **Extract only GH31 modules from the overview.txt file by matching with the sequences.**

Use a given python program to extract the modules from the dataset.

Refer python Program filter_modules_26sep2023.py

25279 modules were found.

Note: module load in HCC.

```
$ ml biopython
```

**Extract manually characterized IDs from the GH31 page of the CAZy.**

- Cilck on the "Characterized (135)"
- Tabular column page will be shown.
- Copy the Characterized IDs, EC numbers from the page and paste them in Excel and select only the IDs which only have hyperlinks (Characterized by CAZy) and make a list.

Place it in a file with the following format.
IDs EC_Number (example below)



## 4. Construction of sequence similarity networks

**SSNpipe**

**Donwload SSNpipe in the HCC virtual computer by invoking the desktop mode in the interactive application option**

**Open terminal and execute the following command**

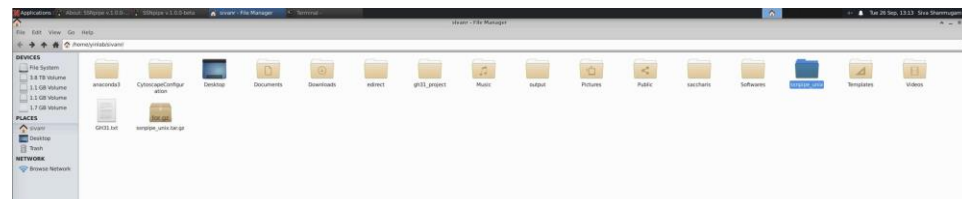**https://github.com/ahvdk/SSNpipe/releases/download/v.1.0-beta/ssnpipe_unix.tar.gz**

$ wget "https://github.com/ahvdk/SSNpipe/releases/download/v.1.0-beta/ssnpipe_unix.tar.gz"
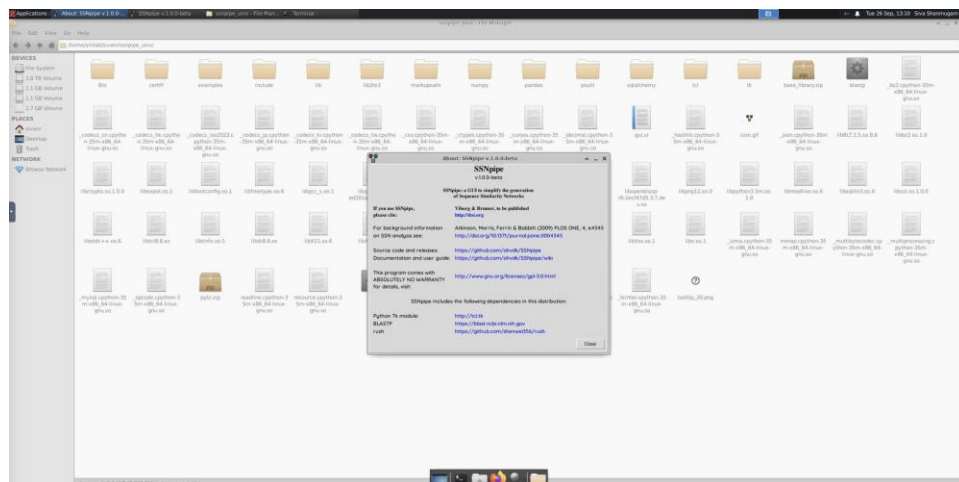
**Extract using following command**
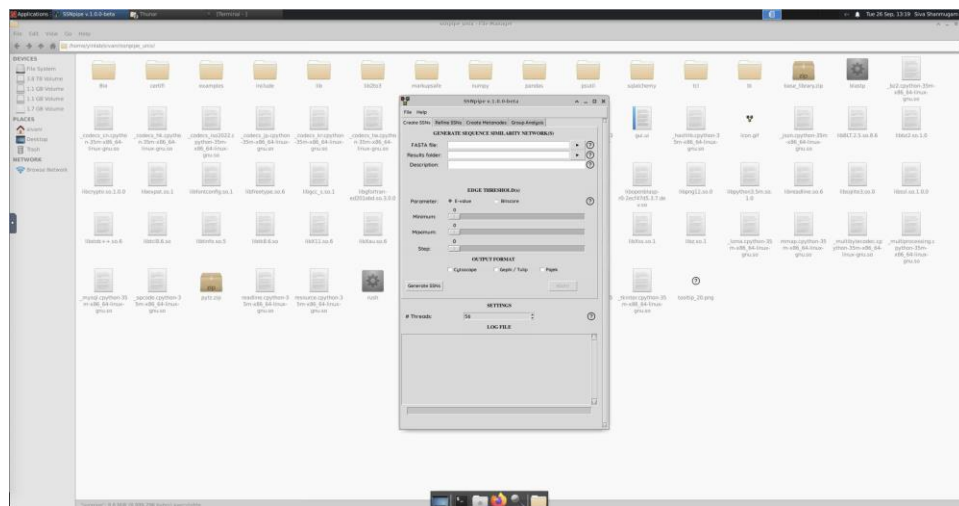
$ tar –xvf ssnpipe_unix.tar.gz



**Move to ssnpipe_unix folder in the explorer**
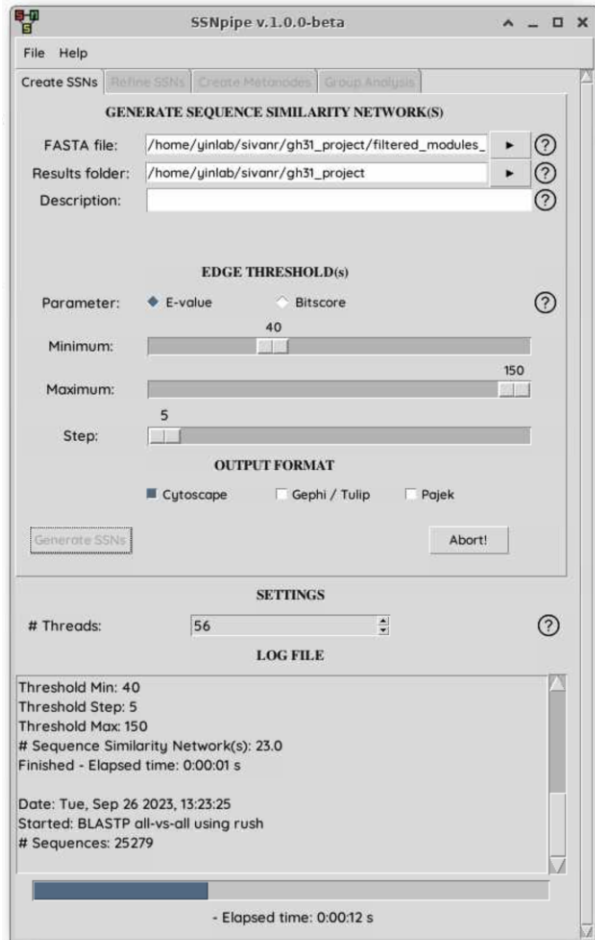


**Then open the application in GUI.**

**You can find the application opened in the top left corner**



Set the input file, output file path, parameters below with the e-value step of 5 from minimum to maximum.

Click on Generate SSNs, the program will run for few hours depending on the database size and computational power.

Note:

Depending on the network files size you may need to create a metanodes (3rd tab in the SSNpipe interface)

Generate group Analysis files from the 4th tab.

Results can be processed with the following python program and EC-Numbers can be mapped.

==Program name: analysis_mapped_1oct2023.py==

Example for the result of "E-value -140".



==Use only one step for the project. Select based on your family. GH31 uses optimized subfamily clusters at E-value of 1E-115.==
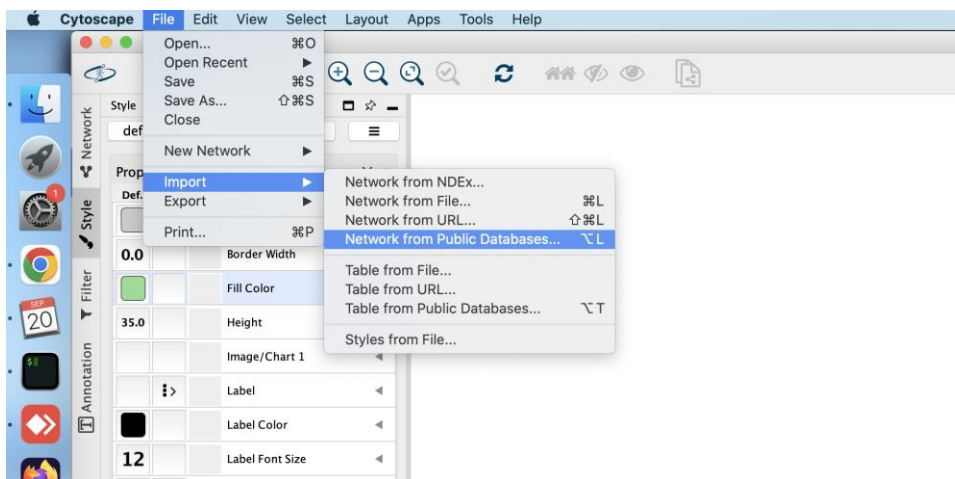
5.  **Visualize the network using the Cytoscape and interpret the network. Create images.**

More details on visualization and analysis are available at: https://cytoscape.org/

Cytoscape can be accessed by HCC Virtual Desktop in the terminal with the following commands.

- Open terminal in the virtual desktop
- Type "ml cytoscape/3.10"
- After loaded type "cytoscape.sh"
- Import file using import option (Ctrl + L)
- Use organic layout for visualization
- Analyze the results of the E-value.

Import the saved network file into the Cytoscape using the



- Use network from File and load them
- Network will be loaded.
- Use a *yFiles Organic layout* from the layout panel. (Install if not available)

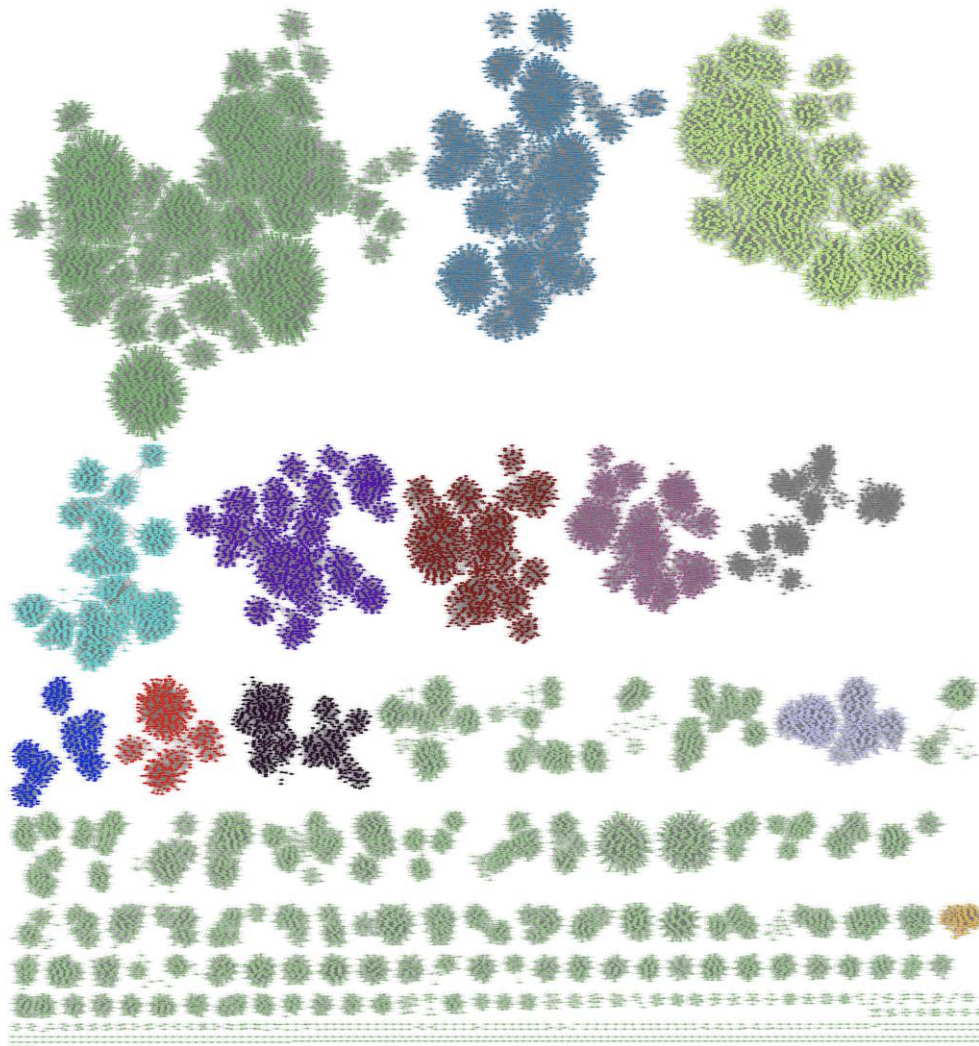- You may notice a change in a layout after the Algorithm is selected.
- Check different nodes and subfamilies, how they are aligned, etc.
- Group them based on the connected nodes and color them for easy visualization.

For example, I have used the network results obtained from the blast results with the cutoff of E-value 1e-140 and their sequence similarity networks (SSN) "NETWORK_cs_ev_140_.txt"

Modified colors based on the connected nodes and subfamilies. (This can be done by selecting the connected nodes in each cluster.)

## 6. Phylogenentic Analysis.

Analyse the results using data obtained from the subfamily random 30 seqs in subfamily where seqs are more than 30 and rest are taken as such (>20 anyways).

**Refer to python program** (phylo_process_random_1oct2023.py)

Use python program to select random 30 sequences each and use them for phylogeny buliding using RAxML.

**HCC module load MAFFT**

```
$ ml mafft/7.520

$ mafft --localpair --maxiterate 1000 --genafpair --thread 10
phylo_input_fasttree_140.fasta > phylo_mafft_align_1e-140.fasta
```
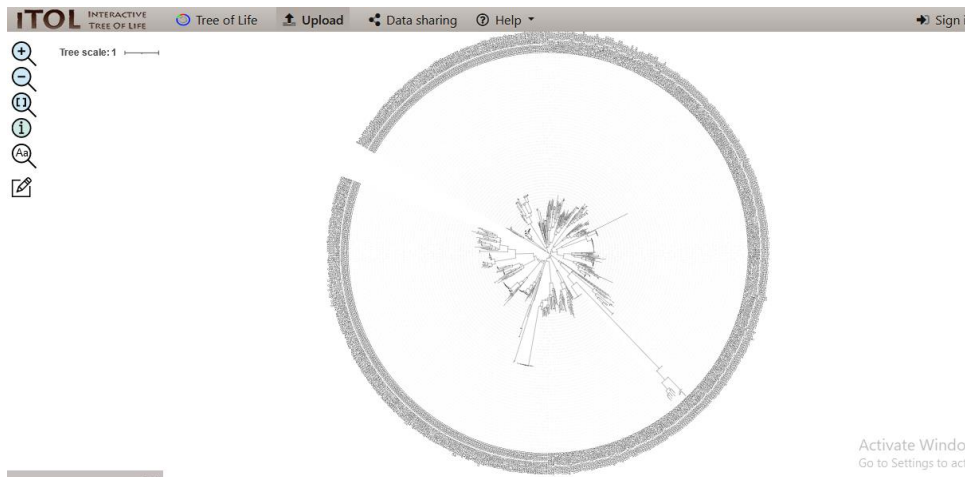
**HCC module load FASTTREE**

```
$ ml fasttree/2.1

$ fasttree -wag -boot 100 -out tree_e-140.nwk
phylo_mafft_align_1e-140.fasta
```

"tree_e-140.nwk" can be visualized using the iTOL available online.

https://itol.embl.de/

Explore the options over the webserver and save the image for the best SSN Evalue

Example for the SSN 1E-140 build tree



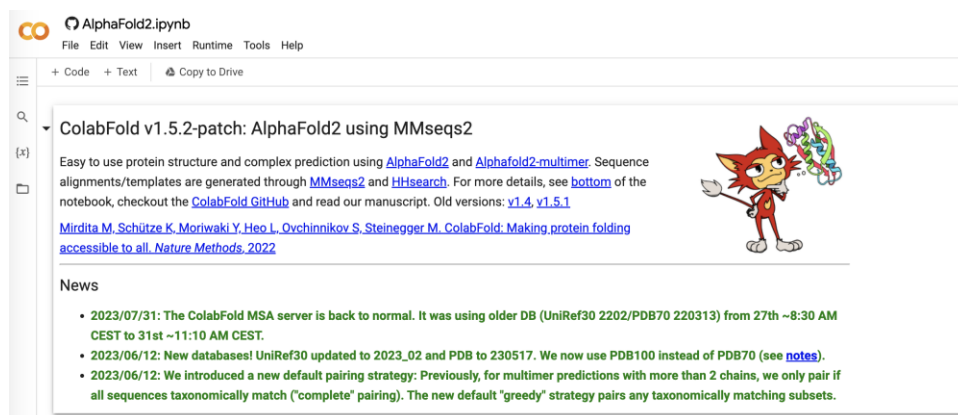Also available at: https://itol.embl.de/tree/13423816579310521696229988

**7. Interpretation and analysis of the subfamilies based on the Phylogeny and subfamilies from the selected E-value.**

**Interpretation steps and protocol:**

(i)    Explain each subfamily in terms of EC code, Taxonomical diversity, enzyme activity and structures.

(ii)   Select any one characterized ID and their EC numbers from each subfamily from the results (Output of "analysis_mapped_1oct2023.py")

(iii)  If PDB ID is available take the structure, or else model them using AlphaFold2.

       a. Use Google Colab code for protein structure modelling, https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb



       b. https://blast.ncbi.nlm.nih.gov/Blast.cgi (helpful to identify the similar sequences/ templates for model building)

(iv)   Identify the common catalytic domain between each subfamily.

(v)    Use PyMol to load and visualize the structures. (Available in HCC virtual Desktop)