

HELP MANUAL

1.0 What does this Programme do?

This programme was scripted to select a family of protein sequences from a user-defined subset of the taxonomic tree, which would subsequently be processed, and used to plot and determine the level of protein sequence conservation, across species within the user-defined taxonomic group. Following which, these plotted sequences are then scanned with motifs from the PROSITE database, to identify if any known motifs (domains) are associated with these sequences.

2.0 How does this Programme achieve this?

2.1 Performs an esearch (NCBI) based on user input.

- Taxon ID + Protein Name
- Taxon ID + Protein Name + Gene Name

2.2 Displays statistical information on the protein lengths of the proteins found in the aforementioned esearch.

- Option to filter protein sequences >X std deviation above the mean
- Option to filter protein sequences >X std deviation below the mean

2.3 Downloads the FASTA files of the proteins, from the updated dataset (following removal of sequences >x std deviation from the mean).

2.4 Filters for redundancy within species.

- Option to filter protein sequences with a 95% threshold
- Option to filter protein sequences with a 100% threshold

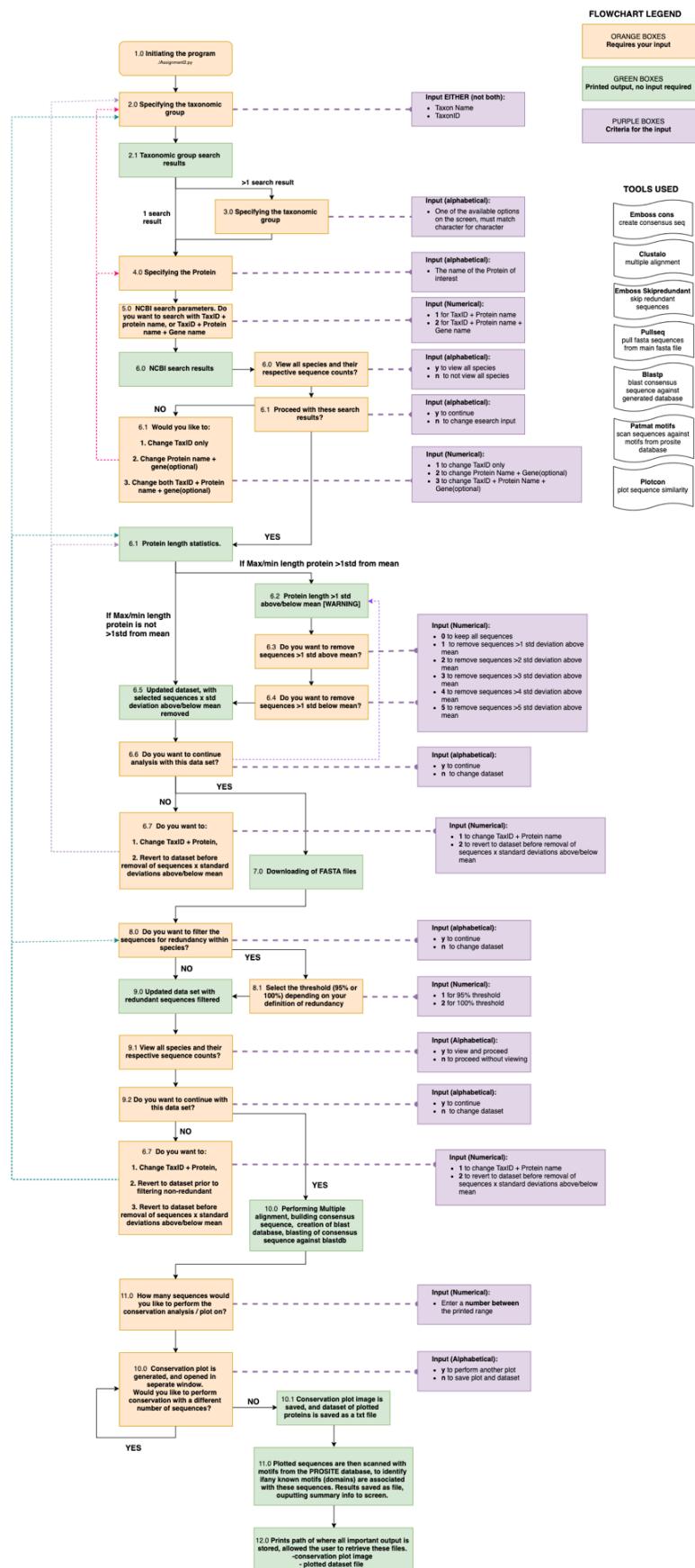
2.5 Performs multiple alignment, builds consensus sequence, creates a blast database from the filtered Fasta file of protein sequences, blasts the consensus sequence against the blast database that was made, to generate a list of proteins from most to least similar.

2.6 Performs conservation analysis, by plotting x most similar sequences on a conservation plot. Outputting the image, showing the protein sequence conservation across species. Image is saved, and dataset of plotted sequences are saved as a .txt file, so the user can see which protein sequences were plotted.

- Option to select between a minimum of 10 sequences, and maximum of 250 sequences, or maximum number of sequences in the dataset, whichever is lower.

2.7 Identifies if any motifs are associated with any of the plotted sequences, using the PROSITE database. Saves the results as a .txt file.

Biologist Flowchart:



3.0 Example of Program input/output:

Below is a series of screenshots of both the inputs required from the user, and outputs that the user will be expected to see.

- This example was performed with
 - Taxon : Aves
 - Protein : Glucose-6-phosphatase

3.1 Running the script, user input required to specify either the Taxon Name / ID

Inputting Taxon Name	Input required
<pre>bioinfmsc5:~/Assignment2\$./Assignment2.py Please specify the taxonomic group Aves</pre>	Taxon name (alphabetical)
<pre>bioinfmsc5:~/Assignment2\$./Assignment2.py please specify the taxonomic group 8782</pre>	TaxonID (Numerical)

3.2 Taxon search output and user input required to specify Protein of interest

Taxon search output + Protein of interest input	Input required
<pre>===== esearch results are as follows: 1. Aves (birds), class, birds ===== Please specify the family of Protein which you would like to analyse for the Aves taxon Glucose-6-phosphatase</pre> <ul style="list-style-type: none">○ Output of taxon search: Displays the esearch result of the taxonomy database, displaying the results in GenBank format.	Protein name (alphanumeric)

3.3 User input required to specify esearch parameters

Printed instructions and inputting choice of esearch parameters	Input required
<pre>Would you like to: ===== 1. eSearch with: "TaxonID: 8782" "Protein name: glucose-6-phosphatase" WARNING: This will perform an extensive search, but conservation plot would not be as biologically significant 2. eSearch with: "TaxonID: 8782" "Protein name: glucose-6-phosphatase" "Gene name: to be specified" WARNING: This would produce a biologically significant conservation plot, however less species would be covered due to poor maintenance of NCBI gene name database ===== Please type 1 or 2</pre>	1 or 2 (numerical) 1: eSearch w/ Taxon and Protein name [program continues to step 3.4]

<ul style="list-style-type: none"> Printed statement: Informs the user of selection criteria and esearch parameters of each choice, enabling the user to make a sensible decision 	<p>2: research w/ Taxon + Protein name + Gene Name</p>
If input was 2 , program will output the following:	<p>Input required</p> <p>Gene name (alphanumeric)</p> <p>[program continues to step 3.4]</p>

3.4 Esearch (NCBI) results and option for user to view full list of species and their respective sequence counts.

Printed instructions and inputting choice of esearch parameters	Input required
<pre>===== Esearch Results are as follows: Taxon : Aves Protein Family : glucose-6-phosphatase Total sequences: 461 No. of Species : 86 ... The top 3 most represented species are as follows: Species: Wire-tailed manakin Sequences: 17 Species: helmeted guineafowl Sequences: 17 Species: blue-crowned manakin Sequences: 17 ... ----- WARNING: The results above contains redundant sequences. WARNING: To improve the speed of the script, FASTA sequences will only be downloaded later, and redundant sequences removed. WARNING: We shall perform a further check later ----- Would you like to view the full list of species and their respective number of FASTA sequences? y/n</pre> <ul style="list-style-type: none"> Output of esearch w/ taxon + protein name + gene name(optional) First section shows taxon, protein family, and total sequences found, and unique species that these sequences belong to. Second section shows the top 3 most represented sequences Third section is a warning message, informing the user that the above results contain redundant sequences 	<p>y or n (alphabetical)</p> <p>y: displays list of all species and their respective sequence counts.</p> <p>n: nothing displayed [program continues to step 3.5]</p>
if input was y , program will output the following:	<p>Input required</p> <p>-NIL-</p> <p>[Program continues to step 3.5]</p>

3.5 Option for the user to continue with the displayed dataset, or change the esearch input

Continue with the above displayed dataset?	Input required
	<p>y or n (alphabetical)</p>

<pre>Would you like to proceed with the analysis based on the above information? y/n y</pre>	<p>y: [program continues to step 3.6]</p> <p>n: options are printed to the screen</p>
<p>if input was n, program will output the following:</p> <pre>===== This is your current search parameters: 1: [Taxon] Aves [TaxonID] 8782 2: [Protein] glucose-6-phosphatase 3: [Gene name] Unspecified by user ===== Which of the above would you like to change? Enter the digit: 1 : To change Taxon, 2 : To change Protein + Gene(optional), 3 : To change Taxon + Protein + Gene(optional). 2</pre> <ul style="list-style-type: none"> ○ Printed statement: First section shows current research input: Taxon + Protein + Gene Second section informs the user of options, to change Taxon, Protein, Gene. 	<p>Input required</p> <p>1, 2 or 3 (numerical)</p> <p>1: repeat research w/ new Taxon + previous specified Protein [program returns to step 3.1]</p> <p>2: repeat research w/ previously specified Taxon + New Protein + Gene(optional) [program returns to step 3.2]</p> <p>3: repeat research w/ new Taxon + new protein + gene(optional) [program returns to step 3.1]</p>

3.6 Output containing protein length statistics, requiring user input to filter sequences >1 standard deviations above the mean

Protein length statistics and user input to decide whether or not to filter sequences	Input Required
<pre>===== Protein Length Statistics: Minimum Length: 73 Maximum Length: 2554 Mean Length: 315.72 Standard Deviation: 177.19 =====</pre> <ul style="list-style-type: none"> ○ Output of the dataset's Protein statistics minimum protein length maximum protein length mean protein length standard deviation of the dataset 	<p>-NIL-</p> <p>if no sequences were >1 std from the mean [program continues to step 3.7]</p>
If min/max protein length is >1 std from the mean	Input Required
<pre>===== PLEASE READ THIS SECTION: WARNING: Maximum or minimum length sequence is >1 standard deviation away from the mean WARNING: What effect could this have on the output? -If this sequence is significantly longer or shorter than other sequences, it could be an anomaly. Eg.human error while uploading to the NCBI's database -Caution must be exercised by the user when deciding which sequences to remove. Please use the information printed below to make a wise decision =====</pre> <ul style="list-style-type: none"> ○ Warning is printed if min/max length seq is >1 std from the mean. 	<p>-NIL-</p>

If maximum protein length >1 std above the mean	Input Required
<pre>WARNING: ABOVE THE MEAN Maximum length sequence is > 1 standard deviations above the mean Your options are as follows for sequences ABOVE THE MEAN : -Choice- -Action- 0 : [DO NOT REMOVE] any sequences 1 : [REMOVE SEQUENCES] 9 Sequences that are 1 Standard Deviations above the Mean 2 : [REMOVE SEQUENCES] 2 Sequences that are 2 Standard Deviations above the Mean 3 : [REMOVE SEQUENCES] 2 Sequences that are 3 Standard Deviations above the Mean 4 : [REMOVE SEQUENCES] 2 Sequences that are 4 Standard Deviations above the Mean 5 : [REMOVE SEQUENCES] 2 Sequences that are 5 Standard Deviations above the Mean Based on the above information, please input a digit for your choice 2</pre>	0, 1, 2, 3, 4 or 5 (numerical) <p>0: no sequences to be removed</p> <p>1: sequences >1 std deviation above the mean</p> <p>2: sequences >2 std deviation above the mean</p> <p>etc:</p>
<ul style="list-style-type: none"> ○ Output of maximum protein length >x std above the mean: Options are printed onto the screen for the user, at increments of 1 std above the mean, and the respective number of sequences found. Capping at 5 standard deviations above the mean, or when 0 sequences are found. 	
If input was 1, 2, 3, 4 or 5 , the program will output the following:	Input Required
<pre>Species Name Species TaxID Prot Accession Prot Length 0 Saker falcon 345164 XP_005447099.1 2664 1 peregrine falcon 8954 XP_005239713.1 2664 The above sequences will be removed. Would you like to continue? y/n y</pre>	y or n (alphabetical) <p>y: the displayed/selected sequences will be removed from the dataset</p> <p>n: program returns to previous step, asking user to input his choice</p>
<ul style="list-style-type: none"> ○ Output based on user selection of sequences >1 std above the mean to be removed: Based on the user input of previous step, the selected proteins are printed on the screen, for the user to make an educated decision to remove them from the dataset or keep them. 	
If maximum protein length >1 std above the mean	Input Required
<pre>===== WARNING: BELOW THE MEAN Maximum length sequence is > 1 standard deviations below the mean Your options are as follows for sequences BELOW THE MEAN : -Choice- -Action- 0 : [DO NOT REMOVE] any sequences 1 : [REMOVE SEQUENCES] 26 Sequences that are 1 Standard Deviations below the Mean Based on the above information, please input a digit for your choice. Sequences to be removed will then be displayed for confirmation 0</pre>	0, 1, 2, 3, 4 or 5 (numerical) <p>0: no sequences to be removed</p> <p>1: sequences >1 std deviation above the mean</p> <p>2: sequences >2 std deviation above the mean</p> <p>etc:</p>
<ul style="list-style-type: none"> ○ Output of minimum protein length >x std below the mean: Options are printed onto the screen for the user, at increments of 1 std below the mean, and the respective number of sequences found. ○ Capping at 5 standard deviations above the mean, or when 0 sequences are found 	
If input was 1, 2, 3, 4 or 5 , the program will output the following:	Input Required
same as above step, no differences!	

3.7 Outputting the updated dataset, following removal of sequences >1 std from the mean.

Displaying updated dataset, and user input to show all species and their respective sequence counts	Input required
<pre>===== Filtered search results are as follows: Taxon : Aves Protein Family : glucose-6-phosphatase Total sequences: 459 No. of Species : 80 The top 3 most represented species are as follows: Species: blue-crowned manakin Sequences: 17 Species: helmeted guineafowl Sequences: 17 Species: Wire-tailed manakin Sequences: 17 ===== WARNING: The results above contains redundant sequences. WARNING: To improve the speed of the script, FASTA sequences will only be downloaded later, and redundant sequences removed. ===== Would you like to view the full list of species and their respective number of FASTA sequences? y/n n</pre> <ul style="list-style-type: none"> ○ <u>Output of updated dataset following removal of sequences >1 std from mean</u> <p>First section shows taxon, protein family, and total sequences found, and unique species that these sequences belong to. Second section shows the top 3 most represented sequences Third section is a warning message, informing the user that the above results contain redundant sequences</p> <p>User input (y/n) required to decide whether or not to view full list of species and their respective sequence counts. Similar to that in step 3.4.</p>	y or n (alphabetical) y: displays list of all species and their respective sequence counts. n: nothing displayed
Checking if user would like to continue analysis with current dataset	Input required
<pre>Would you like to view the full list of species and their respective number of FASTA sequences? y/n n Would you like to go ahead with the analysis with this data set? y/n n</pre> <ul style="list-style-type: none"> ○ <u>Printed statement:</u> <p>Asks the user if he/she wants to continue analysis with current dataset</p>	y or n (alphabetical) y: [program continues to step 3.8] n: user decides not to continue analysis with dataset
If input was n, program prints the following:	Input required
<pre>Would you like to go ahead with the analysis with this data set? y/n n Which dataset would you like to revert to? 1 : Start again, change TaxID and Protein 2 : Revert to dataset before removal of sequences x standard deviations above/below mean Please input one of the digits above 2</pre> <ul style="list-style-type: none"> ○ <u>Printed statement</u> <p>Options are printed on the screen for the user, giving him 2 options for changing his dataset.</p>	1 or 2 (numerical) 1: restart research w/ new input [program returns to step 3.1] 2: revert to dataset prior to removal of sequences >1 std from mean [program returns to step 3.6]

3.8 Input required from the user, to skip redundancies within the species, or not filter for redundancies and continue with the dataset

Asking user if he would like to skip redundancies within the species	Input required
<pre>Fasta sequences are now being downloaded, please be patient Would you like to filter the sequences for redundancy within species? -NCBI is a redundant database, and therefore redundant sequences could be present -Filtering these sequences will help reduce the bias of the consensus sequence Please input y/n: y</pre>	<p>y or n (alphabetical)</p> <p>y: user decides to skip redundant sequences</p> <p>n: user decides not to skip redundant sequences [program continues to step 3.9]</p>
If previous input was y , program outputs the following:	<p>Input required</p> <p>1 or 2 (numerical)</p> <p>1: skip redundancy within species at 95% threshold</p> <p>2: skip redundancy within species at 100% threshold</p>

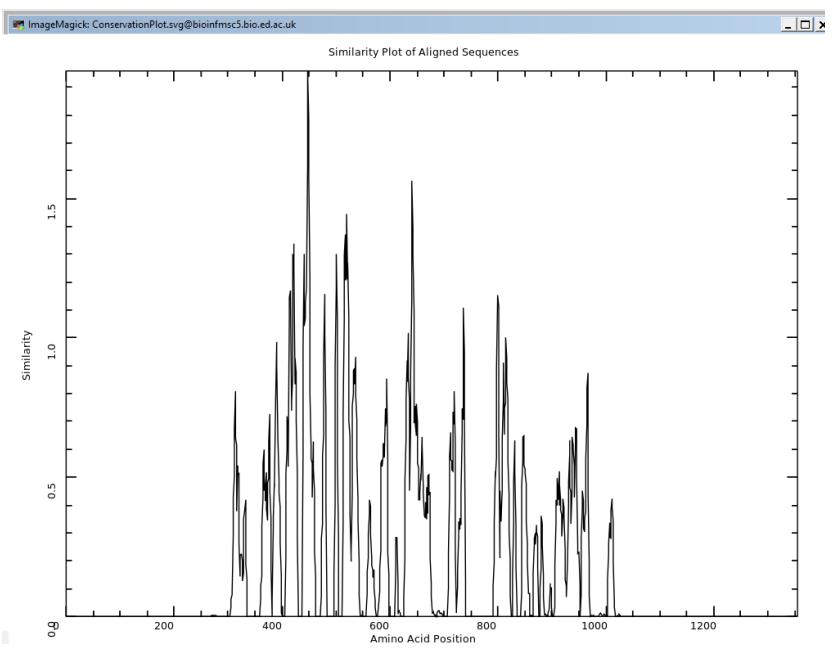
3.9 Printing updated dataset with redundant sequences removed, and checking if user

Printing filtered dataset summary	Input Required
<pre>===== Search Results have been filtered -Results below are non-redundant: Taxon : Aves Protein Family : glucose-6-phosphatase Total sequences: 307 No. of Species : 80 ----- The top 3 most represented species are as follows: Species: Erythrura gouldiae Sequences: 9 Species: Gallus gallus Sequences: 7 Species: Lepidothrix coronata Sequences: 7 ===== Would you like to view the full list of species and their respective number of FASTA sequences? y/n n</pre>	<p>y or n (alphabetical)</p> <p>y: displays list of all species and their respective sequence counts.</p> <p>n: nothing displayed</p>

Asking user if he would like to continue with current dataset	Input required
<pre>Would you like to view the full list of species and their respective number of FASTA sequences? y/n n Would you like to perform the conservation analysis on the above listed data? y/n n</pre> <ul style="list-style-type: none"> ○ Printed statement: Asking the user if he would like to continue the analysis and perform conservation analysis on the current dataset. 	y or n (alphabetical) <p>y: user decides to continue with dataset [program continues to step 4.0]</p> <p>n: user decides not to continue with dataset</p>
if previous input was n, program prints the following:	Input required
<pre>Would you like to perform the conservation analysis on the above listed data? y/n n Which dataset would you like to revert to? 1 : Start again, change TaxID and Protein 2 : Revert to dataset prior to filtering for redundancy 3 : Revert to dataset before removal of sequences x standard deviations above/below mean Please input one of the digits above 2</pre> <ul style="list-style-type: none"> ○ Printed statements: Information and instructions are printed on the screen, to help the user make a decision on which dataset to revert to or if he/she wants to start fresh with new research input. 	1, 2 or 3 (numerical) <p>1: start with new research input [program returns to step 3.1]</p> <p>2: revert dataset to the one prior to filtering for redundancy [program returns to step 3.8]</p> <p>3: revert dataset to the one prior to removal sequences >1std from the mean [program returns to step 3.6]</p>

3.10 Multiple sequence alignment, generation of consensus sequence, creation of blast database, and blasting consensus sequence against the newly made blast database, followed by plotting of conservation plot based on user input

Multiple sequence alignment, generation of consensus seq, creation of blastdb, and blasting consensus against blastdb	Input required
<pre>Performing multiple alignment... Creating consensus Sequence... Create a consensus sequence from a multiple alignment Building a new DB, current time: 11/17/2019 11:41:41 New DB name: /localdisk/home/s1962117/Assignment2_11_26/fastafiles/blast/fastadatabase New DB title: /localdisk/home/s1962117/Assignment2_11_26/fastafiles/filtered.fasta Sequence type: Protein Keep MBits: T Maximum file size: 1000000000B Adding sequences from FASTA; added 307 sequences in 0.0165651 seconds. Blasting the consensus sequence against the non-redundant database... How many sequences would you like to perform the conservation analysis on? Minimum number : 10 Maximum Number : 250 Please input a number between the above listed numbers 250</pre>	Minimum: 10 Maximum: 250 or total no. of sequences, depending which one is smaller. (numerical)

<ul style="list-style-type: none"> ○ Printed statements: Statements will be printed to the screen to inform the user that: multiple alignment is being performed, consensus sequence is being built, blast database is being built, and lastly, the consensus sequence is being blasted against the newly generated blast database. ○ Printed statement for input: Prints the range that the user is able to choose from, how many sequences are to be plotted. The x (user defined) sequences are then selected, these sequences are the most similar sequences to the Consensus sequence 	
output of conservation plot and asking user if he wants to repeat conservation plot on different number of seq.	Input required
<pre>Created /localdisk/home/s1962117/Assignment2_11_26/conservationplot/ConservationPlot.sv Would you like to generate another conservation plot, with a different number of sequences? y/n y</pre> <ul style="list-style-type: none"> ○ Printed statements: ○ Graph is generated, and path of the graph is printed on the screen. Graph is also displayed in a separate window ○ User is prompted if he wants to repeat the conservation plot 	y or n (alphabetical) <p>y: user decides to repeat conservation plot on different no. of seq. [program returns to step 3.10]</p> <p>n: user decides not to repeat conservation plot [program continues to step 3.11]</p>
Conservation plot is displayed in a separate window  <ul style="list-style-type: none"> ○ Example of a conservation plot being displayed in a separate window. This example in particular was for: aves, glucose-6-phosphatase, with 2 sequences >2 std deviations above mean removed, and filtered for redundancy at 100% within species. 250 of the most similar protein sequences were plotted. 	Input required -nil-

3.11 Performing motif analysis on the plotted sequences, Identifies if any motifs are associated with these sequences, using the PROSITE database. Prints the summary of the Motif analysis, and prints path/location of Conservation plot image, plotted sequence dataframe, and motif analysis result file.

Motif analysis summary and paths to conservation plot, plotted sequence dataframe, and full motif analysis output.	Input required
<pre>Motif analysis is complete ===== Motif analysis summary: Total Sequences analyzed/Plotted: 237 Total number of motifs found: 6 Total number of Sequences with Motifs: 64 Motifs that were Found: Motif = AMIDATION For more in-depth analysis of the sequences and their identified motifs, please access the motif analysis output file. Path to this file will be displayed below ===== Location/path of files that were generated by this script are as follows: Conservation Plot Image: /localdisk/home/s1962117/Assignment2_10_24/conservationplot/ConservationPlot.svg Plotted Sequence dataframe: /localdisk/home/s1962117/Assignment2_10_24/conservationplot/PlottedProteinDataFrame.txt Motif analysis output: /localdisk/home/s1962117/Assignment2_10_24/motifsofplotted/Plottedmotifs.txt ===== The script has come to an end, and so has the endless sleepless nights haha bioinfmc5:~/Assignment2\$</pre> <ul style="list-style-type: none"> ○ <u>Output of motif analysis for plotted sequences</u> The first section shows the summary of the motif analysis; Total sequences analysed/plotted, number of motifs found, and total number of motifs found, as well as the names of the motifs that were found. For further in-depth analysis, user is provided the path of the motif analysis output, which he can delve further into. ○ <u>Printed statement:</u> Prints the file locations/path of conservation plot, plotted sequence dataframe, and motif analysis output 	-NIL-

4.0 Important tips so as not to crash the program/script:

- Patience is key, **please do not input anything into the terminal while the script is running tasks in the background.** Prompts for input are extremely self-explanatory and obvious, as long as the user reads the prompts.
- If **esearch result returns with results >1, do not proceed with the analysis.** The script prompts the user to review the dataset/dataframe nearly every time it is altered, giving the user multiple options to revert the dataframe to a prior one, or even perform esearch with different input/parameters.
- If **esearch generates >10,000 results,** the script is designed in such a way for the user to curate and filter sequences, and if the user is not happy with the dataset, to revert the dataset or perform esearch again with different input. Only late into the script are fasta files downloaded. This helps speed up the script, by not having to continuously redownload fasta files and full sequences. If the esearch generates >10,000 sequences, please do not go ahead with the script as I can assure you that you'll be celebrating 5 birthdays before clustalo is done with the multiple alignment.
- **Although error traps are present at every input, please do read the input criteria, as prompts have been made as detailed as they can be.**
- **Output folder is generated in the home space with the hour and minute,** if the user decides to terminate the script immediately after the output folder is generated and runs the script within the same minute, there is a possibility that the file will fail to create as a minute may not have passed yet, thus same folder name (could be bypassed if I had named the folder with the seconds, sorry about that)

Maintenance Manual (competent Python3 code-writer)

5.0 Function Flowchart:

