

CLANS-Python Tutorial

Version 0.1

Overview

CLANS-Python is a program for visualizing the relationship between proteins based on their all- against-all pairwise sequence similarities. The program implements a version of the Fruchterman-Reingold force directed graph layout algorithm to present the sequence similarities in a 2D or 3D graph.

CLANS-Python is composed of two components:

- A command-line tool (clans_cmd.py), which gets as an input a set of sequences in FASTA format, performs all-against-all BLAST search to obtain a matrix of sequence similarities.
- A GUI-based visualizing tool (clans.py), which gets a matrix of sequence similarities and displays them as a dynamic graph using the Fruchterman-Reingold force-directed layout.

The pairwise similarity scores (E-values in case of BLAST search) are used to calculate the attractive forces between each sequence pair. The better the score (lower E-value), the higher the attractive force. In addition, each sequence repulses every other sequence with a certain force (inversely proportional to their distance in space). Clustering is achieved by iteratively moving sequences according to the force vector resulting from all pairwise interactions (attraction and repulsion).

Installation

Using conda:

Requirements: Conda / Anaconda installed on the target computer. An OS-specific version of Anaconda can be downloaded from: <https://www.anaconda.com/products/individual-d>.

- Download the CLANS-Python package code.
- Create a conda environment using the 'CLANS_python.yml' file:
`conda env create --file CLANS-Python/CLANS_python.yml`
- Activate this environment when running the CLANS-Python program:
`conda activate CLANS_python`

BLAST search:

In order to run a BLAST search using clans_cmd.py, BLAST+ (version 2.6.0. or newer) must be installed on the computer.

Using the command-line tool

The command-line tool can be used to perform a BLAST search in order to create a matrix of sequence similarities and/or to perform a specific number of iterations of the force-directed graph layout calculation (which can be later loaded and displayed in the visualizing tool).

Usage:

```
python clans_cmd.py -infile <fasta_file_path> -saveto <destination_file_path> [options]
```

or

```
python clans_cmd.py -load <network_file_path> -dorounds <number of iterations> -saveto  
<destination_file_path> [options]
```

Mandatory arguments:

-infile <fasta file path>: a FASTA file input for BLAST search
or

- load <network file path>: Load an existing network file in CLANS or tab-delimited formats
- saveto <destination file path>: A destination path for saving the output file (in CLANS format, by default)

Optional arguments:

- h, --help: Show this help message and exit
- input_format <'clans'/'mini-clans'/'delimited'>: The format of the network input file (when using the -load option).
Accepted formats: 'clans' (default) or tab-delimited.
- output_format <'clans'/'mini-clans'/'delimited'>: The format of the network output file (default is 'clans' format).

BLAST-related arguments:

- eval <E-value threshold>: E-value threshold for extracting BLAST HSPs (default=1.0)
- matrix <scoring matrix>: Scoring matrix for BLAST search.
Options: 'BLOSUM62', 'BLOSUM45', 'BLOSUM80', 'PAM30', 'PAM70' (default: BLOSUM62)

Clustering-related argumenets:

- dorounds <number of iterations>: Number of clustering iterations to perform (default=0)
- pval <similarity threshold>: A threshold for the similarity score (default=0.0001).
In case the similarity scores type is 'hsp' (result of a BLAST search), sequences with scores below the threshold are considered connected.
In case the similarity scores type is 'att' (attraction values), the sequences with scores above the threshold are considered connected. The default in this case is 0.1.
- cluster2d: Perform the clustering in 2D instead of 3D (default: cluster in 3D).
- cooling <COOLING>: A multiplier for the 'maxmove' parameter. $0 < \text{Cooling} \leq 1$. By default, set to 1 which causes the graph to keep moving until the user stops it. When cooling<1, maxmove gradually converges to 0 and the graph points stop moving.
- maxmove <MAXMOVE>: The maximum distance a point is allowed to move per iteration (default=0.1).
- att_val <ATT_VAL>: A multiplier factor for the calculation of the attractive force between each two sequences (default=10.0).
- att_exp <ATT_EXP>: An integer number - determines how the attractive force scales with the distance between each two vertices in the graph. Default = 1, attraction increases linearly with the distance.
- rep_val <REP_VAL>: A multiplier factor for the calculation of the repulsive force between each two sequences (default=10.0).
- rep_exp <REP_EXP>: An integer number - determines how the repulsive force scales with the distance between each two vertices in the graph. Default = 1, repulsion decreases linearly with the distance.
- dampening <DAMPENING>: A value between 0 and 1, determines to what extent the movement vector of the last movement affects the current movement. The lower it is, the greater the last movement's influence (default=0.2).
- gravity <GRAVITY>: A minimal force that attracts each sequence towards the origin of the graph and prevents unconnected clusters/sequences from drifting apart indefinitely. It scales linearly with the distance from origin (default=1.0).

Using the GUI-based visualization tool

Open the GUI from the command-line

Usage:

python clans.py [-load <network file path>] [options]

When *clans.py* is run without an input file, the GUI is opened empty and loading a file is possible from the 'File' menu of the GUI.

Optional arguments:

-h, --help: Show this help message and exit

-load <network file path>: Load an existing network file in CLANS or tab-delimited formats.

-format <'clans'/'mini-clans'/'delimited'>: The format of the network input file (default is 'clans' format).

-dorounds <number of iterations>: Number of clustering iterations to perform (default=0)

-pval <similarity threshold>: A threshold for the similarity score (default=0.0001).

In case the similarity scores type is 'hsp' (result of a BLAST search), sequences with scores below the threshold are considered connected.

In case the similarity scores type is 'att' (attraction values), the sequences with scores above the threshold are considered connected. The default in this case is 0.1.

Input / Output file formats:

The visualization tool requires as an input a file with sequences identifiers and pairwise similarity scores. The input file can be provided in one of two file formats: 'CLANS' format or tab-delimited format.

CLANS format

A file in 'CLANS' format (.clans), can be created by *clans_cmd.py* (as an output of the BLAST search), by the CLANS web-utility in the MPI Bioinformatics Toolkit (<https://toolkit.tuebingen.mpg.de/tools/clans>) or saved in a previous session of the visualization tool.

The 'CLANS' file format must contain the following blocks of information:

- The first line must be: *sequences=<number of sequences>*
- The sequences block: the original sequences in FASTA format (the order of the sequences is important and is further used to index the sequences, starting from 0).

```
<seq>
>seq0
MSGRGKQGGKARAKAKTRSSRAGLQFPVGR
>seq1
LAAEVLELAGNAARDNKKTRIIPRHLQLAIRNDEELNKLLSGVT
</seq>
```
- The coordinates block: the positions of the sequences in the 3D space. Every line contains the sequence index and a value for the X, Y, Z coordinates ($-1 \geq X, Y, Z \leq 1$).

```
<pos>
0 0.142 0.281 0.104
1 0.298 0.631 0.913
</pos>
```
- The similarity scores block: the E-values or attraction values for the pairwise sequence similarities.

```
<hsp>
0 1:5.1e-05
0 4:1.1e-02
```

```

0 5:6.8e-04
</hsp>
or
<att>
0 1 0.1
0 4 0.2
0 5 0.3
</att>

```

The file may contain additional blocks:

- Parameters block: a list of all the parameters that were used in the calculation and presentation of the saved session.

```

<param>
rounds_done=264
cluster2d=false
pval=0.2
attfactor=10.0
attvalpow=1
repfactor=10.0
repvalpow=1
cooling=1.0
dampening=0.2
maxmove=0.1
minattract=1.0
</param>

```

- Groups block: a list of groups with their definitions and assigned sequences.

```

<seqgroups>
name=Proteobacteria
size=10
hide=0
color=255;0;0;255
numbers=1;3;5;6;10;13;18;19;20;21;23;27;33;35;42;49;52;
</seqgroups>

```

Minimal-CLANS format (mini-clans)

This format is a shorter version of the standard CLANS format, as it does not include the sequences block (<seq>). It may be useful when dealing with a large amount of sequences, to reduce the file size.

Tab-delimited format

This file type should contain a list of non-redundant sequence-pairs and their similarity scores.

Mandatory columns:

- Sequence1 ID / unique name
- Sequence2 ID / unique name
- Similarity score

Optional column:

- Type of score ('hsp' for P-value and 'att' for attraction value)

The Graphical User Interface (GUI) controls

Once loading a network file (from the command-line or by using File -> Load file from the GUI), the graph presentation of the sequences appears in the central part of the GUI and most of the button-controls become enabled.

Menus

File menu:

- **Load file:** Loads a network file, containing at least sequences names and pairwise similarity scores in one of the two accepted formats (see detailed explanations in the previous section):
 - CLANS format
 - Standard (full) CLANS format – includes the sequences block and is compatible with older versions of the CLANS program.
 - Minimal CLANS – does not contain the sequences in FASTA format.
 - Tab-delimited format.
- **Save to file:** Saves the current graph presentation into a file, either in 'CLANS' format or as tab-delimited file.
 - The 'CLANS' format can keep the data-points current coordinates, the groups definitions and all the running parameters. Thus, it is recommended for further visualization using the CLANS-Python tool.
 - Minimal CLANS format does not contain the sequences and can be used to reduce the file size.
 - The tab-delimited format is recommended only as an API with other tools, such as Cytoscape, as it keeps only the minimal information needed to present a network (pairs of sequence names and similarity scores).
- **Save as image:** saves the currently displayed graph session in PNG format.
- **Quit:** closes the application including all its open windows.

Configure menu:

- **Layout**
 - **Fruchterman-Reingold:** Opens a dialog window, allowing the user to configure the parameters that are used while clustering the sequences using the Fruchterman-Reingold graph layout (a detailed explanation of all the parameters is found in the command-line tool section).
- **Data-points default parameters:** Sets the default size, color and outline-color of the data-points, when no other setting is defined. In case different values are set via the groups definitions (each group can have a different setting), they will take over the default setting.

Tools menu:

- **Group data by:**
 - **Taxonomy:** This feature automatically divides the data into groups according to a taxonomic level of the user's choice. It can be applied when the input is a CLANS/FASTA formatted file and the sequences headers contain the organism names in one of the following formats:
 1. Inside square brackets, as obtained by the NCBI database.
 2. 'OS=organism_name' as obtained by the UniProt database.In other cases, the organism names cannot be extracted.
The taxonomic information is taken from the NCBI Taxonomy database.

Prerequisites: The NCBI taxonomy dump files 'names.dmp' and 'rankedlineage.dmp' should be downloaded from the following ftp site and located in the clans/taxonomy/ folder: https://ftp.ncbi.nlm.nih.gov/pub/taxonomy/new_taxdump/

Usage: The first time this feature is applied on a certain dataset, a taxonomic mapping of the organisms that are found in the input file against the Taxonomy database is performed. This process may take a while. Then, a dialog opens, allowing the user to select the taxonomic level to group according to in addition to other groups-related size and font parameters. Once the taxonomic information for a dataset was collected, the user can easily change the selection of the taxonomic level by opening this dialog again without having to wait (Tools -> Group data by -> Taxonomy). Once a certain taxonomic level was selected, it is added to the 'Group by' combo-box, which switches between different grouping options.

The colors of the groups are generated automatically according to the number of groups within each taxonomic level. However, once the groups are generated, it is possible to edit them manually and change their parameters, such as color.

The sequences, for which taxonomic information could not be extracted, or a specific taxonomic level is missing, are grouped together as 'Not assigned'.

- **Add custom grouping category:** This option allows the user to upload a tab-delimited metadata file (see details below), with one or more grouping-categories, by which the data can be grouped. The new grouping category is added to the 'Group by' combo-box, which switches between different grouping options.

The colors of the groups are generated automatically according to the number of groups within each category. However, once the groups are generated, it is possible to edit them manually and change their parameters, such as color via the 'Edit groups' dialog. If there are sequences, for which no group was assigned, they are grouped together as 'Not assigned'.

Metadata file format: Tab-delimited file, containing at least two columns and a header:

Sequence ID	Category1
0	group1
1	group2

Sequence ID: There are two options to provide a unique sequence ID:

1. The sequence name as it is given in the input file. When the input was provided in CLANS/FASTA format, the sequence ID is the sequence header until the first space character.
2. The sequence serial number according to the order in the FASTA section of the clans file (starting by 0).

Note: Once the grouping by either taxonomy or user-defined parameter action is performed, the 'Group by' combo-box is enabled and it is easily possible to switch between the newly added grouping options and the 'Manual definition' default grouping option (where groups are manually defined by selecting sequences/data-points).

- **Color data by:**

- **Sequence length:** This feature colors the data-points according to the length of the sequences. The color-range can be set by clicking the colored 'change' buttons.
- **User-defined parameter:** This option allows the user to upload a tab-delimited metadata file (see details below), with one or more numeric parameters, according to which the data can be colored. Once a metadata file is loaded, it is possible to set a different color-range for each parameter. An unlimited number of metadata files with different parameters can be uploaded and the information is added to the already existing user-defined parameters.

Metadata file format: Tab-delimited file, containing at least two columns and a header:

Sequence ID	Param1	[Param2, Param3, ...]
0	30	0.22
1	50	0.75

Sequence ID: There are two options to provide a unique sequence ID:

1. The sequence name as it is given in the input file. When the input was provided in CLANS/FASTA format, the sequence ID is the sequence header until the first space character.
2. The sequence serial number according to the order in the FASTA section of the clans file (starting by 0).

Parameter: Int / float number.

Note: Once the coloring by either sequence length or user-defined parameter action is performed, the 'Color by' combo-box is enabled and it is easily possible to switch between the newly added coloring options and the coloring by groups (or default color in case there are no defined groups).

Interaction with the graph area

CLANS-Python has three distinct interaction modes within the graph area:

Rotate/Pan graph mode:

- Rotating the graph is done by holding the left mouse button pressed + dragging the mouse.
- Panning the graph is done by shift + left mouse button + dragging the mouse.
- Moving the selected data-points to a different location in the graph is done by CTRL + left mouse button + dragging the mouse.

Select data-points mode:

Allows to select specific data-points from the graph, by using the mouse:

- To select a specific data-point, locate the mouse cursor on that point and click the left mouse button. Another click on a selected point will deselect it.
- To select all the data-points within a rectangular area, hold the left mouse button pressed + drag the mouse.

The selection can be done in two modes:

- Data-points mode
- Groups mode: available when there is at least one defined group. Clicking on a point selects all the data-points that belong to the same group. When selecting an area by mouse-dragging, all the data-points that belong to groups within that area will be selected.

Move/Edit text mode:

This mode allow to use the mouse in order to change the location of the group names text elements in the graph area or to edit them.

- Left button click on a group name + mouse drag: drags the text to the desired position in the scene.
- Double-click on a group name opens a dialog in which it is possible to edit the text (name, size, color, etc.). Please note that changing the color of the group name will also change the color of the data-points belong to that group.

Zoom in/out is done by scrolling the mouse down/up. It is in possible in all interaction modes.

GUI Controls

Clustering options:

- **Initialize:** generates random coordinates to all the data-points.
- **Start:** starts the iterative clustering process. The number of iterations that were done is displayed below the graph area.
- **Stop:** stops the clustering process. It be resumed from the same point by clicking 'Resume clustering'.
- **Cluster in 3D / 2D** combo-box: determines whether the clustering process is done in 3D or 2D (default is 3D, unless stated differently in the running parameters). Setting the clustering to be done in 2D will change the graph view to 2D as well.
- **P-value / Attraction value threshold:** this threshold determines at which similarity score sequences are considered as connected. When the similarity score is a p-value, all the sequence-pairs with lower score are considered connected. In case the similarity score is an 'attraction value', sequence-pairs with similarity score above this threshold are considered connected.

Interaction mode combo-box:

Switches between the three modes of interaction with the graph (see above). The default is 'Rotate/Pan graph' mode.

View options:

- **3D / 2D view**: switches the graph view (not the clustering) between 3D and 2D view. The 2D view can be useful for producing images, as the connecting lines are displayed "behind" the data-points. Another option which is only available in 2D view is displaying and moving the group names (if any).
- **Full / selected dataset view**: switches between the full dataset presentation (default) to a presentation of the selected subset only. When the selected subset is displayed, any operation that is done is performed on the subset only. For example, it is possible to start the clustering considering the subset data-points only. When switching back to full dataset, the view changes back to the exact full-data presentation that was displayed before. If a clustering was performed on the subset, it is not automatically saved, but can rather be saved to a file while in 'selected subset' mode. The 'Selected subset' viewing mode is only enabled when there are at least 5 selected data-points.
- **Z-index** combo-box: this option is only enabled if there are defined groups and the view is set to 2D. It allows the user to determine whether the Z-indexing of the data-points will be done automatically (the default) or according to the groups order (which can be set in the 'Manage groups' window). Note: Z-index by groups order is only allowed when the clustering process has stopped, as it slows the graphics significantly.

Display options:

- **Connections**: when this button is checked, the connecting lines (edges) are presented in the graph. In order to display the edges "behind" the data-points, switch to a 2D view mode.
- **Selected names**: when this button is checked, the names (headers) of the selected data-points are presented next to them (enabled when there is at least one selected point).
- **Group names**: When this button is checked, a list of the group names is presented on the left side of the graph area.
 - **All / Selected**: By default, when the 'Group names' button is checked, all the group names are presented. If there is at least one selected group, it is possible to display the selected group names only.
 - **Reset names**: Clicking on this button brings the group names back to the top-left part of the scene. It is useful after rotating the graph, zooming-in/out, etc.
- **Color by**: Determines by which parameter the data is colored. This combo-box is enabled once a 'color by' action (Tools -> Color data by:) was performed (either by sequence length or by a user-defined parameter). It is possible to switch between coloring the data by the defined groups (or default color when no group is defined) or by any of the numeric parameters, including sequence length. When each parameter has a different color-range setting, the data is colored according to the selected parameter's setting.

Selection options:

CLANS-Python enables to select a subset of data-points and perform all kind of operations on this subset. The selected data-points are marked in the graph by having a bigger size and a magenta outline color.

- **Selection mode** combo-box: switches between the selection of specific data-points and the selection of whole groups. In data-points selection mode, the points that are clicked or within the selected area, are being selected. In groups selection mode, when a point is clicked (or within the selected area), all the points that belong to the same group are being selected. This combo-box is only enabled in 'Select data-points' interaction mode and in case there is at least one defined group.

- **Select all:** selects all the data-points.
- **Clear:** clears the selection.
- **Select by name:** this option enables to enter a search term and select data-points according to their sequence name or header. Clicking on this button opens a find dialog, in which a search term can be entered. The results are displayed in the 'Search results' window (see detailed explanation in the 'Windows' section).
- **Edit selected sequences:** This button is enabled only in case there is at least one selected data-point. Clicking on it opens a window, presenting the selected sequences IDs and headers. The 'Selected subset' window can remain open and the display gets updated whenever there is a change in the selection subset (see detailed explanation in the 'Windows' section).

Groups options:

In order to help the user to navigate in the clustering map and examine whether the clustering meets other classifications, CLANS-Python allows to define groups of sequences/data-points, which can be displayed in the graph by different colors. There are several ways to define groups:

1. **Manual definition** (default option): the groups can be manually defined by selecting data-points and adding them to a new or existing group or use the 'select by name' option.
2. **Input file:** The groups can be pre-defined in the CLANS input-file, using the <seqgroups> block.
3. **Taxonomy:** Use the Tools -> Group data by -> Taxonomy feature to automatically group the sequences according to the NCBI taxonomic classification. Each taxonomic level is a grouping category containing groups of sequences.
4. **User-defined parameter:** Upload a tab-delimited file with one or more parameter, by which the data is grouped.

Groups Controls:

- **Edit groups:** opens a dialog presenting a list of all the existing groups (in parentheses the number of group members). When selecting a group, the following options are possible:
 - **Edit group:** opens another dialog, in which it is possible to change the group name, the size and the color in which the group's data-points are presented in the graph. At the bottom of this window there is a list of the group's members, in which it is possible to select sequences and remove them from the group. Any change that is done in this dialog is updated in the graph.
 - **Delete group:** Remove this group's definition. The data-points that were assigned to this group are not assigned to any group anymore and get the default presentation (black color).
 - **Move up:** Move this group one step up in the order of the groups list.
 - **Move down:** Move this group one step down in the order of the groups list.
The groups order is only relevant when viewing the graph in 2D and setting the Z-index to be done by the groups order.
- **Add selected to group:** this option is enabled when there is at least one selected data-point. It opens a dialog, allowing the user to add the selected data-points to a new group or to an existing group (if any). If the user chooses to create a new group, a new dialog is opened, allowing the user to define the group name, size and color parameters.
- **Remove selected from group:** removes the selected data-points from their group(s) if they were assigned to any group. They get the default presentation (black color).
- **Group by:** Determines by which grouping category the data is grouped. This combo-box is enabled once there is at least one more grouping category in addition to the default 'Manual definition'. (The additional categories can be added by methods no. 2-4 described above). It is then possible to switch between all the defined grouping categories. When selecting a specific 'group by' category, all the groups-related controls refer to the selected category (for example: the 'edit groups' dialog will present the groups that belong to the selected group-by category).

Windows

Selected subset window

This window presents the IDs and headers of the sequences in the selected subset. It is possible to mark one or multiple lines (sequences) by using the Shift / CTRL buttons or by clicking the left mouse button + dragging the mouse.

The following operations can be performed on the marked sequences:

- **Highlight in graph:** highlights the relevant data-points in the graph with turquoise color.
- **Set as selected subset:** leaves only the marked sequences as the selected subset (removing the other from the subset).
- **Remove from subset:** removes the marked sequences from the selected subset.
- **Find in subset:** opens a find dialog, in which a search term can be entered. If there are sequences in the subset that match the search term, they are marked in blue and the above operations can be done on them.

All the changes to the selected subset are updated in the graph as well.

The window can either remain open or be closed. The sequences display gets updated whenever there is a change to the selection subset.

Search results window:

This window presents the sequences that match the search term entered by the user.

It is possible to mark one or multiple lines (sequences) by using the Shift / CTRL buttons or by clicking the left mouse button + dragging the mouse. Clicking the 'Highlight all' button marks all the lines.

The following operations can be performed on the marked sequences:

- **Add to selected subset:** adds the marked sequences to the selected subset.
- **Set as selected subset:** sets the marked sequences as the selected subset. If there were other sequences in the subset, they are removed.
- **New search:** opens the find dialog again and allows the user to enter a new search term.

All the changes to the selected subset are updated in the graph as well.