

CSI: Cleavage Site Investigator

Features

- Run straight from command line
- Compatible with FASTA file format (.fa and .fasta)
- Determine top and bottom strand cleavage events
- Export results to .csv files
- Create visual event distributions as heatmaps and strand linkage plots

Contents

- [Features](#)
- [Contents](#)
- [Installation](#)
- [Usage](#)
 - [Notes](#)
 - [Running CSI \(basic\)](#)
 - [Running CSI \(advanced\)](#)
 - [Generating strand linkage plots directly](#)
 - [Generating heatmap plots directly](#)
- [Outputs](#)
 - [CSI summary file](#)
 - [CSI individual results file](#)

Installation

1. Install Python (tested with Python 3.9.1)
2. Install required libraries ([BioPython](#), [Seaborn](#), [SVGWrite](#) and [TQDM](#))
 - Either using Pip

```
pip install biopython==1.79
pip install tqdm==4.55.1
pip install seaborn==0.11.1
pip install svgwrite==1.4
```

- Or using the provided Anaconda environment file ("csi.yml" in "resources" folder)

```
conda env create -f csi.yml
```

Usage

Notes

- Example files for testing CSI are included in the "data" folder of this repository. These files are:
 - "ex_cassette.fa" - Cassette sequence (must contain one sequence). Example file is for "Splint1TA".
 - "ex_consensus.fa" - Consensus sequence(s) (can contain multiple sequences). Example file is a subset of sequences from "Cas12a_17.fa" sample at [TODO - RDSF link].
 - "ex_reference.fa" - Reference sequence (must contain one sequence). Example file is for "CrisprlasR".
- The above files are used throughout the following code demos.
- Each program (csi.py, heatmap.py and strandlinkageplot.py) can be run entirely from command line. Full argument documentation is accessible using the `-h` (or `--help`) flag (e.g. `python csi.py -h`).

Running CSI (basic)

- The main CSI program is run using csi.py. This will analyse the specified consensus sequences and optionally output event distributions, summary statistics and plots (advanced plotting options available by running heatmap.py and strandlinkageplot.py directly).
- CSI requires a minimum of three arguments, specifying paths to the cassette (`-ca` or `--cassette_path`), reference (`-r` or `--reference_path`) and consensus (`-co` or `--consensus_path`) files.
- The following command is an example

```
python .\src\csi.py -ca .\data\ex_cassette.fa -r .\data\ex_reference.fa -co .\data\ex_consensus.fa
```

- With default parameters (no optional arguments specified) a basic summary will be displayed with the following sections:

Label	Description
"TS position"	Position of the top-strand cleavage event
"BS position"	Position of the bottom-strand cleavage event
"Split seq"	<code>True</code> if the cleavage event spanned the start/end of the reference sequence, <code>False</code> otherwise
"Count"	Number of identified events matching this cleavage event (% of total identified events shown in parenthesis)
"Type"	Type of cleavage event (either "Blunt end", "3' overhang" or "5' overhang")

- An example output is shown below:

RESULTS:

Full sequence frequency:

TS position: 1289
 BS position: 1293
 Split seq: False
 Count: 396/787 (50.3% of events)
 Type: 5' overhang

TS position: 1293
 BS position: 1293
 Split seq: False
 Count: 97/787 (12.3% of events)
 Type: Blunt end

TS position: 1284
 BS position: 1293
 Split seq: False
 Count: 67/787 (8.5% of events)
 Type: 5' overhang

...

Running CSI (advanced)

- CSI offers optional command line parameters to specify execution settings (e.g. the number of bases to fit) as well as additional outputs (e.g. summary CSV files or rendered heatmap plots).

Optional argument	Description	Default
<code>-h, --help</code>	Show help message (lists all required and optional arguments).	NA
<code>-rf, --repeat_filter</code>	Expression defining filter for accepted number of repeats. Uses standard Python math notation, where 'x' is the number of repeats (e.g. 'x>=3' will process all sequences with at least 3 repeats).	NA
<code>-lr, --local_r</code>	When grouping sequences at restriction sites, this is the half width of the local sequences to be extracted. For example, for a sequence 5'...AAT ATT...3', <code>-lr 1</code> would yield "TA", whereas <code>-lr 2</code> would yield "ATAT".	1
<code>-mg, --max_gap</code>	Maximum number of nucleotides between 3' and 5' restriction sites.	10000
<code>-mq, --min_quality</code>	Minimum match quality. Specified in the range 0-1, where 1 is a perfect match.	1.0
<code>-nb, --num_bases</code>	Number of bases to match when comparing sequences (e.g. when searching for cassette ends in a consensus sequence).	20

Optional argument	Description	Default
<code>-pr, --print_results</code>	Prints results to the terminal once a complete file has been processed.	NA
<code>-en, --extra_nt</code>	Number of additional nucleotides to be displayed either side of the cleavage site (when <code>-pr</code> or <code>--print_results</code> is specified).	0
<code>-sp, --show_plots</code>	Display plots showing local sequence distributions as a heatmap and pie-chart.	NA
<code>-wslp, --write_strandlinkageplot</code>	Write strand linkage plot image to SVG file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_strandlinkageplot'. To generate strand linkage plots with greater control over rendering, see Generating strand linkage plots directly	NA
<code>-whsa, --write_heatmap_svg_auto</code>	Write heatmap image (only spanning range of identified event positions) to SVG file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_heatmap'. To generate heatmaps with greater control over rendering, see Generating heatmap plots directly .	NA
<code>-whsf, --write_heatmap_svg_full</code>	Write heatmap image (spanning full range of reference sequence) to SVG file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_heatmap'. To generate heatmaps with greater control over rendering, see Generating heatmap plots directly .	NA
<code>-whca, --write_heatmap_csv_auto</code>	Write heatmap image (only spanning range of identified event positions) to CSV file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_heatmap'. To generate heatmaps with greater control over rendering, see Generating heatmap plots directly .	NA
<code>-whcf, --write_heatmap_csv_full</code>	Write heatmap image (spanning full range of reference sequence) to CSV file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_heatmap'. To generate heatmaps with greater control over rendering, see Generating heatmap plots directly .	NA
<code>-wi, --write_individual</code>	Write individual cleavage results to CSV file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_individual'. For more information on the individual results file format, see CSI individual results file .	NA
<code>-ws, --write_summary</code>	Write summary of results to CSV file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_summary'. For more information on the summary results file format, see CSI summary file .	NA

Optional argument	Description	Default
<code>-wo, --write_output</code>	Write all content displayed in console to a text file. Output file will be stored in consensus file folder with same name as the consensus file, but with the suffix '_output'.	NA
<code>-ad, --append_datetime</code>	Append time and date to all output filenames (prevents accidental file overwriting).	NA
<code>-v, --verbose</code>	Display detailed messages during execution.	NA

Generating strand linkage plots directly

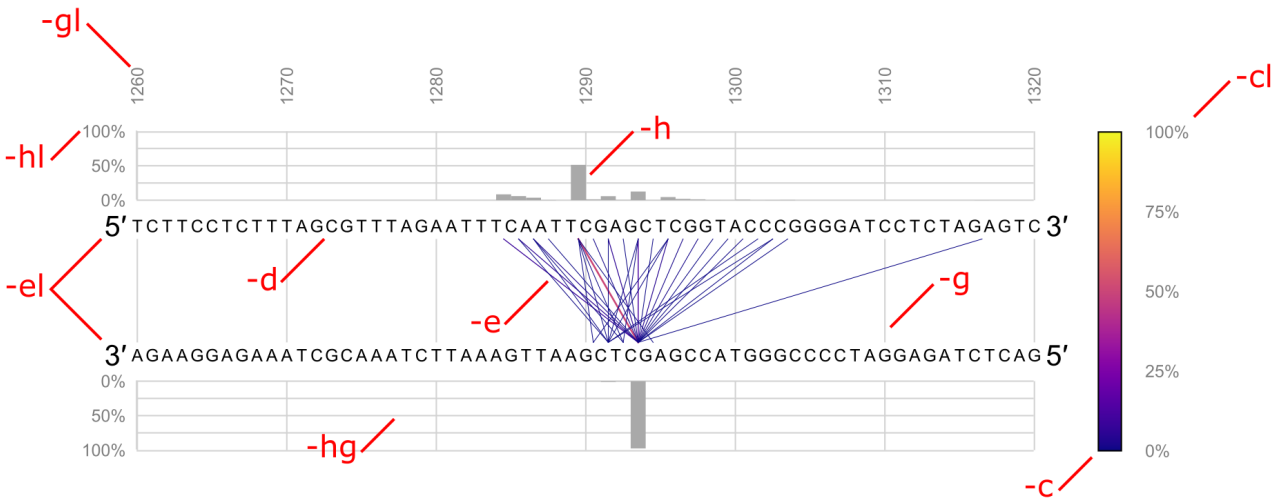
Basics

- Strand linkage plots can be exported to SVG directly from CSI [summary](#) and [individual](#) results files using `strandlinkageplot.py`.
- At a minimum, `strandlinkageplot.py` requires arguments specifying the path to a CSI summary or individual results file (`-d` or `--data_file` argument) and the output SVG path (`-o` or `--out_path` argument).
- For example, the following command will generate a strand linking plot using default parameters:

```
python .\src\strandlinkageplot.py -d .\data\ex_consensus_summary.csv -o .\data\output_strandlinkageplot.svg
```

Advanced control (optional arguments)

- To afford greater control over various aspects of plot rendering, `strandlinkageplot.py` accepts over 50 different command line arguments. Full descriptions for these arguments can be viewed using the `-h` or `--help` flag.
- Optional arguments are grouped by the plot feature they act upon. For example, `-gl_i` controls the grid label increment.
- The following figure and table summarise these regions and arguments:



Root argument	Feature	Instances
<code>-d, --dna</code>	DNA sequence	<code>-d_m (--dna_mode)</code> <code>-d_s (--dna_size)</code> <code>-d_c (--dna_colour)</code> <code>-d_rg (--dna_rel_gap)</code>
<code>-el, --end_label</code>	End label (i.e. 5' and 3')	<code>-el_v (--end_label_vis)</code> <code>-el_s (--end_label_size)</code> <code>-el_c (--end_label_colour)</code> <code>-el_rg (--end_label_rel_gap)</code> <code>-el_p (--end_label_position)</code>

Generating heatmap plots directly

Outputs

CSI summary file

- Summary CSV files contain a pair of information rows (second row containing just bottom-strand sequence) for each unique restriction site identified in the consensus sequence(s).
- The final row of each summary file reports the number of consensus sequences for which cleavage events could not be determined.
- An example summary file is included in the "data" folder ("ex_consensus_summary.csv").
- Summary files include the following columns:

Column	Description
"TYPE"	Type of cleavage event (either "Blunt end", "3' overhang" or "5' overhang").
"COUNT"	Number of identified events matching this cleavage event (% of total identified events shown in parenthesis).
"EVENT_%"	Percentage of all identified events (i.e. doesn't include unmatched sequences) corresponding to this event.
"TOP_POS"	Position of the top-strand cleavage event.
"BOTTOM_POS"	Position of the bottom-strand cleavage event.
"SPLIT_SEQ"	TRUE if the cleavage event spanned the start/end of the reference sequence, FALSE otherwise.
"TOP_LOCAL_SEQ"	Sequence immediately 5' and 3' of the cleavage event on the top strand. The number of nucleotides included either side is determined by the <code>-lr</code> (or <code>--local_r</code>) command line argument.
"BOTTOM_LOCAL_SEQ"	Sequence immediately 5' and 3' of the cleavage event on the bottom strand. The number of nucleotides included either side is determined by the <code>-lr</code> (or <code>--local_r</code>) command line argument.

Column	Description
"SEQUENCE"	Complete top and bottom strand sequences spanning both cleavage sites. The first row corresponds to the top strand and the second to the bottom strand. Cleavage sites on each strand are represented by the " " character.

CSI individual results file

- Individual results files contain a pair of rows (second row containing just bottom-strand sequence) for each consensus sequence processed.
- An example individual results file is included in the "data" folder ("ex_consensus_individual.csv").
- individual results files include the following columns:

Column	Description
"INDEX"	Index of this sequence in the input consensus sequence file. Numbering starts at 1.
"HEADER"	Header text for this sequence. This is any text on the ">" line immediately preceeding the sequence in the FASTA file.
"TYPE"	Type of cleavage event (either "Blunt end", "3' overhang" or "5' overhang").
"TOP_LOCAL_SEQ"	Position of the top-strand cleavage event.
"BOTTOM_POS"	Position of the bottom-strand cleavage event.
"SPLIT_SEQ"	TRUE if the cleavage event spanned the start/end of the reference sequence, FALSE otherwise.
"TOP_LOCAL_SEQ"	Sequence immediately 5' and 3' of the cleavage event on the top strand. The number of nucleotides included either side is determined by the -lr (or --local_r) command line argument.
"BOTTOM_LOCAL_SEQ"	Sequence immediately 5' and 3' of the cleavage event on the bottom strand. The number of nucleotides included either side is determined by the -lr (or --local_r) command line argument.
"SEQUENCE"	Complete top and bottom strand sequences spanning both cleavage sites. The first row corresponds to the top strand and the second to the bottom strand. Cleavage sites on each strand are represented by the " " character.