

scanRBP

What is scanRBP?

Installation

Quick Start

Motif Database Resources

Potential additional PWM datasets

Potential CLIP datasets

Gene Annotation

What is scanRBP?

scanRBP loads RNA-protein binding motif PWM and computes the log-odds scores for all the loaded RBPs across a given genomic sequence + draws a heatmap of the scores.

The scores can be described as follows ([biopython docs](#)):

Unset

Here we can see positive values for symbols more frequent in the motif than in the background and negative for symbols more frequent in the background. 0.0 means that it's equally likely to see a symbol in the background and in the motif.

Using the background distribution and PWM with pseudo-counts added, it's easy to compute the log-odds ratios, telling us what are the log odds of a particular symbol to be coming from a motif against the background.

For more information, see the [biopython docs](#).

Installation

The easiest way to install **scanRBP** is to simply run:

Unset

```
pip install scanRBP
```

Note that on some systems, **pip** is installing the executable scripts under `~/.local/bin`. However this folder is not in the PATH which will result in "command not found" if you try to run "scanRBP" on the command line. To fix this, please execute `"export PATH=\"$PATH:~/.local/bin"` (and add this to your `.profile`). Another suggestion is to install inside a virtual environment (using `virtualenv`).

If you would like instead to **install the latest developmental version** from this repository:

```
Unset
# clone scanRBP GitHub repository
git clone https://github.com/grexor/scanRBP.git

# build and install
./build.sh
```

Quick Start

scanRBP quick start:

```
Unset
scanRBP --help
```

Usage for single sequence: scanRBP sequence output [options]
* sequence on the command line, generates output.png/pdf + output.tab

Usage for processing FASTA file: scanRBP filename.fasta [options]
* one heatmap/matrix will be generated per sequence
* output name of the files will be sequence ids provided in the fasta file

Options:

-annotate	Annotate each heatmap cell with the number
-xlabels	Display sequence (x-labels), default False
-only_protein TARDBP	Only for specified protein
-all_protein TARDBP	All possible motifs for provided protein
-figsize "(10,20)"	Change matplotlib/seaborn figure size
-heatmap title	Create heatmap (png+pdf) with title
-output_folder folder	Output folder (default: current workdir)
-nonzero	Set all values<0 to 0 (default: disabled)

Examples:

```
Unset
# random sequence: produce binding scores and a heatmap
# output: example1_PWM.tab
# (log-odds vectors for all proteins for the given command line sequence)
# output: example1.png/pdf
# (heatmap image with clustering of protein binding vectors)
```

```
./scanRBP
AAAGCGGCGACTTATTATATCCCATATATTATATCTTCTTCTTATATATAAACAGAGATAGATGTGTGTGGTGG
example1 -heatmap example1

# input can be a multi-sequence fasta file
./scanRBP data.fasta
```

Motif Database Resources

Currently, **scanRBP** is using two motif database resources:

- **mCross** motif database
 - Feng H, Bao S et al., [Modeling RNA-Binding Protein Specificity In Vivo by Precisely Registering Protein-RNA Crosslink Sites](#), Molecular Cell, 2019
- **CISBP-RNA** motif database
 - Ray D, Kazan H, Cook KB, Weirauch MT, Najafabadi HS, et al. [A compendium of RNA-binding motifs for decoding gene regulation](#), Nature, 2013

Potential additional PWM datasets

- **RCRUNCH**: A compendium of RNA-binding motifs for decoding gene regulation
 - [Data file](#)

Potential CLIP datasets

- [Encode list of eCLIP experiments](#)

Gene Annotation

- [NCBI gene metadata \(names, aliases\)](#)