

MoDec

MoDec is a motif deconvolution software that finds the motifs and corresponding binding core offsets describing the data in a list of peptides. It is described in the publication (available [here](#)):

Racle, J., et al. Robust prediction of HLA class II epitopes by deep motif deconvolution of immunopeptidomes. *Nat. Biotech.* (2019).

Installation

1. Download MoDec-1.1.zip file and move it to a directory of your choice, where you have writing permissions.
2. Unzip MoDec-1.1.zip package.
3. (Optional) In order to make an html report of the runs from MoDec and to draw the logos from each motifs, you need R and some R-packages:
 1. Rscript needs to be in your path.
 2. You need to install the R-package *ggseqlogo*, based on the fork found at <https://github.com/GfellerLab/ggseqlogo>. For this, from within R, type the command:
`devtools::install_github("gfellerlab/ggseqlogo")`
 3. You also need the R-package *htmlTable* that can be found at <https://cran.r-project.org/web/packages/htmlTable/index.html>. Command to install from within R:
`install.packages("htmlTable")`
 4. You then finally need to open the file *make_logo_report.R*, locate the line starting with `.libPaths` in this file and replace the `"PATH/TO/R/LIBRARIES"` by the path where you installed these R-packages.
4. To test your installation, make sure you are in *MoDec-1.1* directory and run the following command, depending on your operating system:
 - Mac OS: `./MoDec -i test/testData.txt -o test/test_out --Kmax 5 --pepLmin 12 --nruns 2 --specInits --makeReport`
 - Unix: `./MoDec_unix -i test/testData.txt -o test/test_out --Kmax 5 --pepLmin 12 --nruns 2 --specInits --makeReport`
 - Windows: `MoDec.exe -i test/testData.txt -o test/test_out --Kmax 5 --pepLmin 12 --nruns 2 --specInits --makeReport`

Running the software can take few minutes depending on data size and number of motifs. If you didn't do the optional step (3) above, the `--makeReport` option shouldn't be used as MoDec won't be able to build this report.

Your folder *test_out* should contain the same elements than *out_compare* folder (with few differences about the timing and file location indicated within the results files; the motifs found in *test_out/test_out_report.html* and *out_compare/test_out_report.html* should be highly similar but there can be small differences if MoDec is not run on Mac OS due to the

random initial conditions that rely on different random number generators in the various operating systems - note that also the order of the motifs can differ due to these different random number generators).

The *testData.txt* file corresponds to HLA-II peptidomics data from meningioma tissues obtained in our study (sample 3869-GA, with 5987 unique peptides).

5. (Optional) To run MoDec from anywhere on your computer, make an alias of MoDec executable (see above for which one depending on operating system) or add it in your path.

Running

Command

```
MoDec -i input_file -o output_name [additional options]
```

- Depending on your operating system, use Modec, MoDec_unix or MoDec.exe as indicated in the installation instructions.
- Do not use spaces in your file or directory names.
- Do not use other special characters (e.g., *, ?, %, &, ...) in file or directory names.

Required arguments

- Input file (command **-i** or **--input**): File listing all the peptides with one peptide per line. It can also be a fasta file (lines starting with ">" are skipped).
- Output name (command **-o** or **--output**): The name of the output for the current run (including the path where to save it). MoDec will create a folder of this name and save all results in this folder.

Optional arguments

- **--Kmin 1** and/or **--Kmax 6**: To run MoDec searching for motifs between the *Kmin* value (default is 1) and *Kmax* value (default is 6).
- **-L 9** or **--Lmotif 9**: To search for motifs of the given length (default is 9).
- **--pepLmin 12**: To remove from the input the peptides shorter than this length (default is equal to the motif length). **For HLA-II peptidomics, we recommend using a value of 12, to remove potential HLA-I contaminants often present in the data.**
- **-r 20** or **--nruns 20**: To make this number of runs for the deconvolution with each selected number of motifs (default is 20).
- **-b 1** or **--nruns_out 1**: To output the results from this number of the best runs (by default we just output the results from the best run but can also output more). When this is used, the standard *report.html* still contains only the best results while another *..._report_multiRuns.html* file is created showing the results from the various runs.

- **--specInits**: To perform the computations based on multiple types of initial condition biases (e.g. performing some runs with a hydrophobic preference at P1). **This is recommended for HLA-II peptidomics** to increase the chances of finding the best motifs representing the data. When this is used, MoDec will run in fact 5 times more runs than what is given by the option **--nruns** (this is why the result summary files and *report.html* indicate a *n_runs_input* and *n_runs_real* value).
- **--makeReport**: To make an html report containing the logos found in the best run from each number of motifs (need to have R installed and *ggseqlogo* as described in step 3 above).
- **-y 3** or **--logoYmax 3**: Only useful with **--makeReport**. To define the max y value for the plots of the logos (default is 3).
- **--no_flat_mot**: When used, MoDec will not include a flat motif representing contaminant peptides.
- **--out_fullResp**, **--out_w_ks**: To output some additional files giving the full responsibilities of each peptides and the full weight of the motifs.
- **--outAdd**: Normally, MoDec will not run if the folder given in the **--output** argument already existed. When this option is used, MoDec will still run. **This option is to use with care**: previous results might get overwritten as it is using same file names, and, if different parameters are used for the different runs, the report html file will only indicate one of the sets of parameters used, without warning. This option is useful when launching MoDec for different numbers of motifs (**Kmin**, **Kmax** arguments): in such a case the results won't be overwritten as the output has a different name for each number of motifs searched for.
- **--seed 38230811**: To give an initial seed for the random number generator used by MoDec (default value is 38230811). When using the seed value 0, MoDec will try using a fully random seed or use the current time if such fully random device isn't available.
- **--log**: To make that MoDec outputs its logs to a log file created in the result folder instead of outputting it to the terminal.

Results returned and additional information

- MoDec creates a folder containing various files and subfolders with the results.
- The files *Summary/res_K.._out.txt* gives various information, like the parameters used in the current run, the final log-likelihood or a comparison between the motifs (PWM values) found in the various run (of this number of motifs - this comparison is found in the last lines, values correspond to Kullback-Leibler divergences between each motif from the best set of motifs and the set of motifs from each other run).
- The PWMs correspond to the $\theta^{\{k\}}_{\{l,i\}}$ from equation (1) of our manuscript.
- The peptide responsibilities are returned in the *Responsibilities/bestPepResp_....txt* files. This gives for each peptide its two best responsibility values indicating to which motif and offset this corresponds as well as the corresponding core binding sequence.
With the option **--out_fullResp** it is possible to view the complete matrix of the peptide

responsibilities (values from each peptide towards each motif and towards each binding core offset).

- The optimal number of motifs is difficult to automatically determine and existing measures provide unsatisfactory results. We use the Akaike Information Criterion (AIC) as a guide but we encourage the users to manually look at the motifs and determine the optimal number (AIC information is found in the html report and in *Summary/nMotScores....txt*).
- In the html report, the numbers indicated below each logo correspond to the weighed number of peptides associated to this motif (based both on the peptide weights and on the peptide responsibilities). The number in parenthesis next to it is the fraction of peptides assigned to this motif.
- The PWM and logo of the flat motif is also returned by MoDec. Note however that these have one extra position with respect to the other PWMs and the requested motif length (i.e. it has a length of 10 when the default `--Lmotif 9` was asked): the last position corresponds to the real values used for this flat motif in the deconvolution, and the other *Lmotif* first positions correspond to how this motif would have looked based on the peptides associated to this flat motif (this can be for example useful in order to verify if this motif is really flat or not).

Latest version

Latest version of MoDec is available at <https://github.com/GfellerLab/MoDec>.

License

MoDec can be used freely by academic groups for non-commercial purposes (see the license file). The product is provided free of charge, and, therefore, on an "as is" basis, without warranty of any kind.

FOR-PROFIT USERS: If you plan to use MoDec (version 1.1) or any data provided with the script in any for-profit application, you are required to obtain a separate license. To do so, please contact eauffarth@licr.org at the Ludwig Institute for Cancer Research Ltd.

Contact information

For scientific questions, please contact Julien Racle (julien.racle@unil.ch) or David Gfeller (david.gfeller@unil.ch).

For license-related questions, please contact Ece Auffarth (eauffarth@licr.org).

How to cite

To cite MoDec, please refer to:

Racle, J., et al. Robust prediction of HLA class II epitopes by deep motif deconvolution of immunopeptidomes. *Nat. Biotech.* (2019).