

# GEMSEC Neuropeptide Docs

January 28, 2019

*Chris Pecunies, Aaron, Savvy Gupta*

*Led by: Jacob Rodriguez*

*PI: Mehmet Sarikaya*

## Abstract

The purpose of this documentation is to serve as a reference for the motivation, background, methods, and results of the GEMSEC project to discover the method of signal transduction in human neuropeptides, to ultimately discover candidate peptides with binding affinity to graphene, as well as to illuminate their mechanism of binding.

## Introduction

The overarching aim of this project is to uncover the mechanism of signal transduction in neuropeptides which have been experimentally determined by GEMSEC to bind to graphene. To achieve this goal, we will be using several different measures of peptide similarity to ultimately select three peptides from a set of several thousand neuropeptides to experimentally determine the graphene binding of affinity of, ultimately in the overarching goal of determining the relative efficacy of these techniques highly specified peptide synthesis. To determine candidate peptides to be experimented on, we have utilized three principal methods: cross-sequence entropy, PAM30 sequence similarity, and RRM relative to two known binders, *GrBP5* and *M6*

## Background

While there is extensive literature surrounding the subject of peptide-solid binding and its most important determinants, we have focused on two primary methods: sequence domain functionality clustering and the *Resonant Recognition Model* proposed by Irena Cosic.

The *Resonant Recognition Model*, or RRM, is a method proposed by Irena Cosic for the analysis of peptides and proteins. The model assumes electron-ion interaction potential (EIIP, hereafter) along the backbone of amino acid chains to be the most significant determinant in a peptide or protein's biochemical phenomenological features. For each amino acid, a unique scalar EIIP value can be determined, which allows a amino acid sequence to be converted to a scalar sequence of EIIP values in its most simple form.

## Methods

### RRM

The Resonant Recognition Model is a theoretical approach in determining protein functionality, proposed by Irena Cosic. The Resonant Recognition Model (hereafter, RRM) posits that the primary determinant of biological or physiochemical functionality for a protein or peptide lies in the stochastic delocalized electron potential along its carbon backbone. This premise is motivated by the observed high correlation between the spectra of an amino acid sequence represented numerically according to some unique categorical physical functionality and its biological activity.

The electron-ion interaction potential is determined by the pseudopotential of delocalized electrons along the backbone of the carbon backbone. In this manner, each of the 20 essential amino acids which form the set of neuropeptides to be analyzed are represented by a scalar EIIP value ranging from 0 (Leucine, Ile) to 0.1263 (Asp). These EIIP values are then used to calculate a spectrum unique to each sequence through a Discrete Fourier Transform (DFT) which is calculated, as proposed by Cosic:

$$X(n) = \sum x(m)e^{-j\frac{2nm\pi}{N}} \subset n = 1, 2 \dots N/2 \quad (1)$$

These DFT coefficients are ultimately calculated for each neuropeptide sequence. Cosic proposes cross-spectral function between EIIP spectra calculated in this manner, which allows for the determination of shared "characteristic peaks" which, depending on the magnitude and signal-to-noise ratio of the peak, allows for one to theoretically determine how similar one peptide's biochemical and physiochemical functionality is to another. This cross spectrum is calculated as follows:

$$S(n) = X(n)Y^*(n) \subset n = 1, 2 \dots N/2 \quad (2)$$

This is implemented in our peptide similarity matrix with a manual DFT implementation corresponding to results which cohere to the DFT spectra in the introductory paper and book on the Resonant Recognition Model by Cosic. We initially attempted to implement the DFT using the NumPy and SciPy FFT packages, but struggled to acquire initial results which corresponded to the DFT and cross spectra in the aforementioned paper and book, so the manual implementation was added in with the help of Sid and Aaron.

In the cross spectrum generated by our manual implementation of the DFT outlined by Cosic, the signal-to-noise ratio, calculated by dividing the peak value (whose frequency, should the peak value be distinct enough from its background, corresponds to its *characteristic frequency*) by the mean of all values in the cross spectrum, can be deemed a measure of similarity between the two peptides whose DFT spectra are used to create the cross spectrum. In this manner, we

used both the signal-to-noise ratio of the cross spectrum of each neuropeptide to the two known-binding peptides, as well as the correlation coefficient of each neuropeptide's DFT spectrum to the DFT spectrum of the two known-binding peptides as similarity values in the output similarity matrix for each binding peptide.

### Sequence-Function Domain Matching

From (where?) we used a mapping of the 20 amino acids to seven different categories, specifically: aromatic structure, hydrophobic function, polar structure, proline (itself), glycine (itself), negative charge, positive charge, and Cystine which is excluded (assigned to a separate category from all others, like proline and glycine). We then implemented a function which returned all possible subsequences of each known-binding peptide, and then for each peptide, returned a list of the longest matching subsequences and their indices for each neuropeptide. This was done for explicit amino acid subsequence matches, as well as for subsequence matches indicated by the seven aforementioned bio-functional categories, which, given the limited domain, corresponded to a less discriminatory pattern-matching algorithm for each neuropeptide. Thus, we implemented two columns, one containing a list of amino acid subsequence matches and one for functional encoding matches (as tuples, paired with their index of coincidence), in our output similarity matrix.

We also implemented a primitive scoring algorithm, which assigned a cumulative score of 1 multiplied by each independent subsequence match added to 1.5 multiplied by the length of each subsequence match, which allowed us to assign a scalar magnitude to the subsequence matches found. It is our intention to optimize the weight tuple (instead of using the arbitrary 1 and 1.5 aforementioned weights) used in this scoring method further in the project, and implement more complex scoring based on the weight of each match's functional mapping, some which may carry more information and explanatory power than others.

### String distance

Although not the primary aim of our project, in the aim of reasonable comprehensiveness we included several algorithm implementations of sequence similarity in the hopes of providing possible insight into the validity of other higher-order derived similarity metrics. The specific similarity algorithm outputs which we included are as follows: *Jaro-Winkler*, *Needleman-Wunsch*, *Smith-Waterman*, and *Levenshtein*. Each of the aforementioned sequence alignment metrics is unique enough from one another to justify inclusion without potential redundancy in the similarity matrix, and each implements a scoring method that is principally aimed at producing similarity measures for bioinformatic sequences, especially amino acid sequences. Although these are the four measures we included in our final data output for the neuropeptide sequences, the final `SequenceSimilarity` class provides for many more to be calculated and included

in the output .csv file, all included and implemented through the usage of the *textdistance* Python package [omnium\_textdistance\_2019]

### Substitution Matrix Similarity

The PAM (point accepted mutation) matrices, initially developed by Margaret Dayhoff are commonly used in measuring the similarity of two peptide or protein sequences. The PAM30 matrix, specifically, is a sequence alignment matrix that allows 30 point accepted per 100 amino acids. The PAM30 matrix is a “shallow” sequence alignment matrix, in that it is more appropriate in determining alignment for more “similar” sequences.

## Results

Although we implemented several distinct metrics of similarity in our final two output similarity matrices, the initial and primary motivation for the project laid in exploring Irena Cosic’s Resonant Recognition Model as a means for determining peptide behavior, especially with regards to self-assembly, and as such, we weighted the signal-to-noise ratio and correlation coefficient metrics garnered from RRM results the highest, and all other similarity metrics are used at the behest of validating and interpreting the RRM results. The sequence-function-domain pattern matching algorithm represents a relatively novel similarity metric, but its full validity cannot be verified without more experimental effort (or a project on its own).

Our ultimate goal in the initial phase of this project is to downselect from the full list of neuropeptides to three candidate peptides, as determined primarily through RRM signal-to-noise as a similarity metric. To this end, we used our manual RRM implementation, and were able to sort all possible peptides by their signal-noise ratio and use other metrics as benchmarks to judge candidacy. Below is the similarity matrix (minus distance values and amino acid encodings, for space purposes):

**Table 1:** Top 5 sequences by RRM S/N ratio

Sequences	PAM30	BLOSUM45	RRM_SN	RRM_Corr
INNDCQNFIGNR	-66	-10	5.46	0.69
VPLRPEEDELID	-53	-3	5.42	0.65
LNSLDGAGFGFE	-40	-2	5.42	0.89
RSFGCRFGTCTV	-66	-18	5.39	0.64
SSGGGDGSGMWF	-59	-14	5.27	0.78

**Table 2:** Top 5 sequences and subsequence matches by RRM S/N ratio

Sequences	RRM_SN	sseq_matches	num_matches
INNDQCQNFIGNR	5.46	('I', 0)	('1', 0), ('22', 5)
VPLRPEEDELID	5.42	('D', 7)	('1', 0), ('1', 2), ('5', 7)
LNSLDGAGFGFE	5.42		('1', 0), ('5', 4), ('0', 8)
RSFGCRFGTCTV	5.39		('2', 10)
SSGGGDGSGMWF	5.27		('0', 11)

### Candidate peptides

The five similarity metrics generated for each neuropeptide served as a feature set which could be weighted and used to train a regression model, but without experimental data, the relative weights for each feature were chosen by [METHOD].

After forwarding the candidate peptides for experimentation, we utilized statistical clustering and signal processing methods to predict determinant sequence domains in graphene binding which could then be used with the experimental data once acquired.

### Experimental results

After generating the list of [NUMBER] candidate peptides and receiving binding affinity metrics, we were then able to perform several statistical methods to ascertain the most influential factors in peptide graphene binding.

First, we trained a simple linear regression model, using the gathered experimental results as training data, and [...]

### Signal transduction

From results gathered from model training, statistical clustering, and experimental cross-validation, we were able to identify several sequence patterns and corresponding functions which may prove significant. [...]

## Discussion

### Computational results

Our program, using all aforementioned computational methods, generated two similarity matrices for 1,612 same-length neuropeptides: one matrix for similarity values relative to GrBP5, and one for the wild-type peptide M6. The

**Table 1:** Output schema of .csv similarity tables

seq	numseq	rrmsn	rrmcorr	pam30	numentr	aaentr
IMVSTED	1132335	73	88	23	45	55

seq	numseq	rrmsn	rrmcorr	pam30	numentr	aaentr
GTTYUEI	3224121	11	5	6	14	13

Where the columns correspond as follows:

- **seq**: Original amino acid sequence of the neuropeptide
- **numseq**: The original amino acid sequence converted to numbers corresponding to their biological function as determined by the biological function key-value table<sup>[6]</sup>
- **rrmsn**: The signal-to-noise ratio of the cross-signal of the EIIP frequency spectrum calculated for a given neuropeptide and the sequence of comparison, normalized in the domain  $[0, 100]$ .
- **rrmcorr**: The Pearson correlation coefficient between the EIIP frequency spectrum calculated for a given neuropeptide and the sequence of comparison, mapped to the domain  $[0, 100]$
- **pam30**: The PAM30 similarity score (formula described in methods section); a measure of inverse distance using the PAM30 matrix, mapped to  $[0, 100]$
- **numentr**: The cross-entropy of the numerical sequence represented by numseq
- **aaentr**: The cross-entropy of the amino acid sequence

From this data, we determined that [...]

## Experimental results

[to be completed after experimentaiton]

## Clustering and signal analysis

After gathering the experimental results, we were then able to apply supervised statistical learning techniques to the full set of neuropeptides, as well as signal processing techniques. We first [...]

## How to run

**Prerequisites:** A computer running Windows/Mac OSX/Linux with Python (3+) installed, and the Python libraries pandas, NumPy, and scikit-learn. For visualizations, the library matplotlib should be installed. If these libraries are not already installed, instructions for their installation will be listed below.

1. Download the .zip containing all .py files and sample data set **[!FIX]** and extract to preferred location.
2. Open a terminal and navigate to the directory containing the extracted files

1. On Windows, press the Windows key and type “cmd”, and then press enter. `dir` lists all files and folders in the working directory, while `cd dirName` changes the working directory to the specified folder (in this case, `dirname`). To move up the directory hierarchy, type `cd ../`
  2. On Mac OSX, enter spotlight search with Command + Spacebar and type “Terminal”, then hit enter. Instructions are the same as for windows, but replace `dir` with `ls`.
  3. On Linux, a terminal is likely easily accessible. Commands are the same as for Mac OSX.
3. When in the directory containing the extracted files, type the following command to generate the “similarity tables” for the example dataset: `python main.py example_data.csv`
- This can be run with any .csv sequence file, but it must follow the same schema as the `example_data.csv` table provided, and the sequences must all be of the same length relative to each other as well as to the known binder(s) (non length-matching input sequences will be ignored in output)
  - (*To implement*) The script accepts as a second argument a sequence (or list of sequences), each of which will generate its own similarity table for the input .csv.
    - If no argument is given (as above), it will generate two .csv similarity tables, one for GrBP5 and one for M6. To specify only one .csv output similarity table (for GrBP5), run `python main.py example_data.csv grbp5`.
    - This can be done for an arbitrary number of different sequences, for example: `python main.py example_data.csv IVTSSY UVGEASTT EEVTUSGMII` will output three .csv tables for peptides in `example_data.csv` of lengths corresponding to each sequence specified by the user.
    - Finally, the second argument can itself be a .csv of sequences, following the same schema as the first input .csv. In this way, for a .csv entered as a second argument with 10 rows of sequences will generate 10 separate .csv tables.
4. Similarity tables and visualizations will be generated in a folder titled “output” located in the same directory as the `main.py` file.

## Appendix

### Table 1: AA

Function	AA	Num	EIIP
Aromatic	F	0	0.0946
Aromatic	Y	0	0.0516
Aromatic	W	0	0.0548
Hydrophobic	A	1	0.0373
Hydrophobic	V	1	0.0057
Hydrophobic	I	1	0.0000
Hydrophobic	L	1	0.0000
Hydrophobic	M	1	0.0823
Polar	S	2	0.0829
Polar	T	2	0.0941
Polar	N	2	0.0036
Polar	Q	2	0.0761
Proline	P	3	0.0198
Glycine	G	4	0.0050
Charge (-)	D	5	0.1263
Charge (-)	E	5	0.0058
Charge (+)	K	6	0.0371
Charge (+)	H	6	0.0242
Charge (+)	R	6	0.0959
Excluded	C	7	0.0829

## Bibliography