

---

# Translating L-peptides into non-canonical linear and macrocyclic peptides

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Introduction

Protein-protein interactions (PPIs) are critical to several biological functions, from regulating metabolic activities to several disease-causing pathways [1]. With limited success in using small molecules and large biologics to disrupt clinically relevant PPIs, peptides continue to garner further therapeutic interest [2, 3]. As of 2021, around 80 peptides were approved for clinical use worldwide, apart from more than 160 in various stages of clinical trials [4].

Peptides constitute of amino acid monomers, arranged usually in a linear fashion, but also found as macrocycles and other non-linear topologies [5, 6]. The chemical diversity arising from both the amino acid composition and their spatial arrangement [7], synthetic accessibility [8, 9], and potential for cell penetration [10, 11], have made them increasingly sought after modalities in drug discovery.

High-throughput screening for peptides that bind to desired target proteins is often carried out by using libraries made by genetically-encoded or chemical synthesis approaches, such as phage display [12, 13], mRNA display [14, 15], or affinity selection-mass spectrometry [16, 17]. A key challenge in genetically-encoded libraries for screening lies in the limitation of this approach to generate L-peptides, since cells can produce peptides with naturally occurring amino acids. Recent approaches have resulted in libraries with non-canonical amino acids, however, a majority of drug discovery efforts still uses L-peptide libraries [18, 19].

While L-peptide-based libraries are the workhorse of high throughput screening, peptides with non-canonical or unnatural amino acids have been shown to have similar or higher activity and proteolytically more stable [20, 21]. Thus, it is highly desirable to obtain non-canonical variants for L-peptides. Unfortunately, flipping the stereochemistry for each amino acid, and obtaining a D-peptide analogue of the L-peptide, does not work for most cases, necessitating further investigation into developing chemically similar non-canonical variants.

In this work, we developed a method to translate L-peptides into non-canonical linear and macrocyclic peptides. We used a genetic algorithm with a hierarchical dual-objective function, mutating the seed peptide sequence to be chemically similar to the reference L-peptide, while ensuring that the binding affinity remains same or is higher than the seed sequence. Using DRD2 kinase inhibitor as the target protein, we obtained both linear and macrocyclic peptides with non-canonical amino acids with higher activity than observed by random sampling of L-peptides.

## 2 Results and Discussion

### 2.1 L-peptides data set generation

To simulate a data set of L-peptides, as the one that can be obtained from macrophage display, we randomly sampled individual L-amino acids, dipeptides, and tripeptides (Figure 1A). We docked these peptides against DRD2 kinase inhibitor to calculate their respective binding affinities characterized

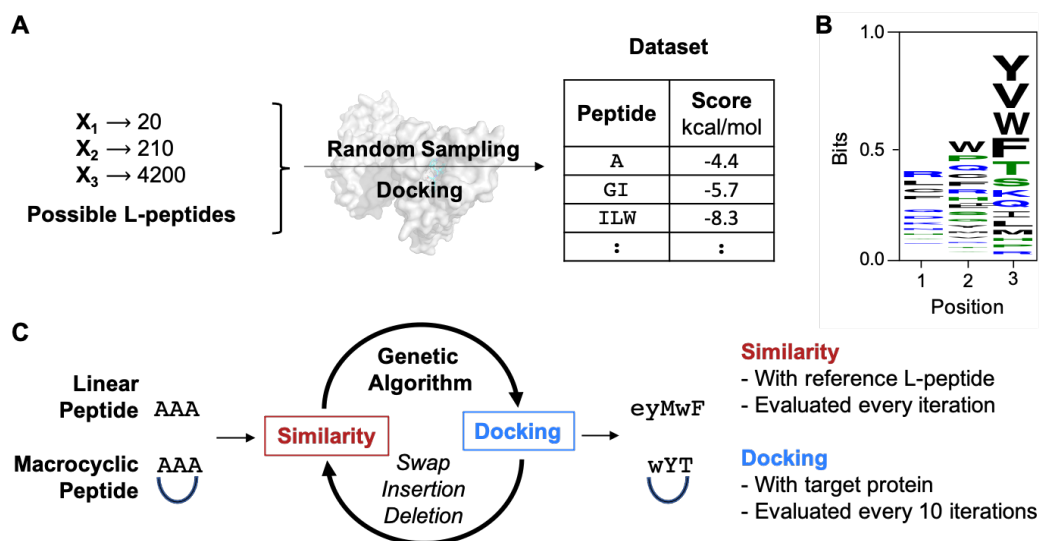


Figure 1: **Outline of the approach for translating L-peptides into non-canonical linear and macrocyclic peptides.** **A.** The data set was generated by random sampling, and docking of L-peptides to DRD2 kinase inhibitor. **B.** Amino acid positional frequency plot, resulting from the weighted multiple-sequence alignment of peptides in the data set was computed for the peptides in the data set. **C.** An overview of the genetic algorithm involving hierarchical similarity and docking for optimization of the seed sequence is shown.

35 by the docking scores. From the data set, we selected the tripeptide YWY, with a docking score of  
 36 -10.3 kcal/mol as the reference peptide (or, the hit compound) for further optimization. We used the  
 37 Python-based DOCKSTRING interface, with AutoDock Vina, for the generation of our initial data  
 38 set, and for later evaluation during the genetic algorithm experiments [22, 23].

39 As a method of evaluating how these peptide binders compared with their small molecule analogues,  
 40 we computed a weighted multiple-sequence alignment of the sequences in the data set [24]. The  
 41 weights were based on the docking scores, with a higher docking score resulting in a directly  
 42 proportional weight, and vice-versa. Visualizing the results using an amino acid positional frequency  
 43 plot, we noted that amino acids with aryl groups, such as Y, tyrosine, and W, tryptophan, were in  
 44 abundance (Figure 1B). The results are aligned to the observation of the interacting phenyl groups in  
 45 small molecules binding to DRD2 kinase [25, 26].

## 46 2.2 Representation of peptide

47 Peptides were converted to molecular structures, and represented using circular fingerprints. Molecu-  
 48 lar structure representations for both linear and macrocyclic peptides, with L- and D-amino acids,  
 49 were obtained by converting amino acid sequences to simplified molecular-input line-entry system  
 50 (SMILES) strings using RDKit [27]. Bit-based and count-based circular fingerprints, with radius 3,  
 51 size 128 and active chirality, were calculated using RDKit [28].

## 52 2.3 Genetic algorithm and optimization

53 Genetic algorithm was set up with a hierarchical objective of (1) achieving high chemical similarity  
 54 to the reference peptide, and (2) increasing binding affinity with the target protein (Figure 1C).  
 55 For mutations, random insertion, deletion, and swapping of amino acids in the seed sequence were  
 56 performed. The ability to mutate in the discrete monomer space, using a pool of L- and D-amino acids,  
 57 helped in increasing the diversity of the peptides, without any worry of synthetic accessibility. In the  
 58 case of macrocyclic peptides, we mutated on the seed peptide string, and then defined a SMILES  
 59 arbitrary target specification (SMARTS) based intra-molecular cyclization reaction to cyclize the  
 60 linear peptides by a head-to-tail amide bond.

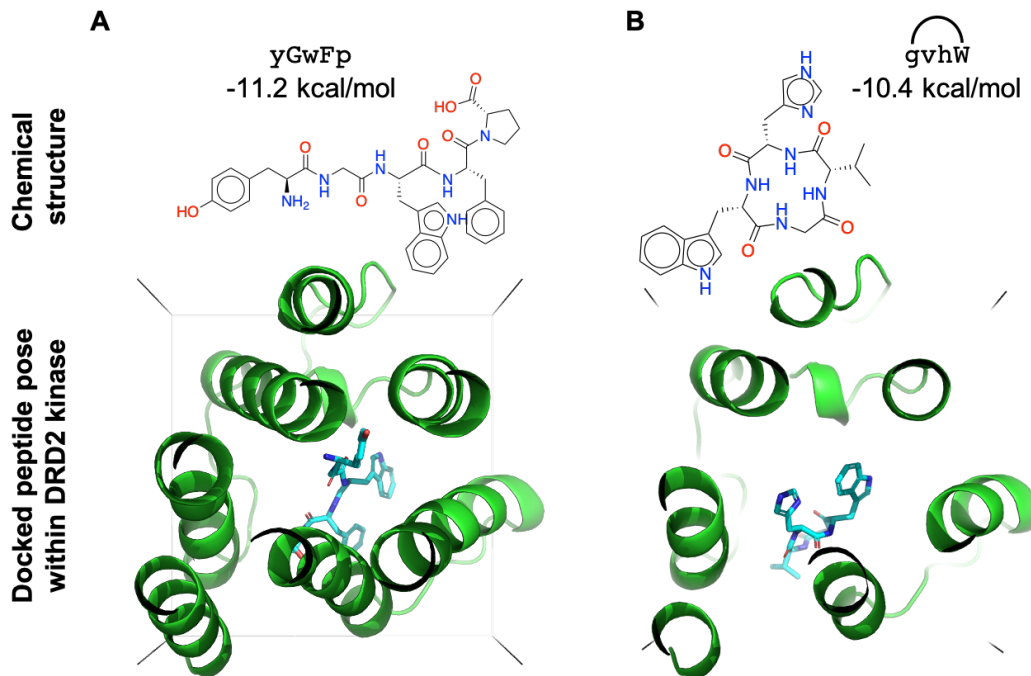


Figure 2: **Optimization for non-canonical linear and macrocyclic peptides.** Sequence information, chemical structure, docking score against DRD2 kinase, and the docked pose of the peptide in the protein are shown for **A.** linear yGwFp, and **B.** macrocyclic gvhw peptides.

61 In our experiments, we used a hierarchical approach of optimizing for similarity, followed by  
 62 optimization for docking scores, motivated by the hypothesis that higher similarity to reference  
 63 peptide will lead to better docking. Briefly, we approached the problem with the idea that if we could  
 64 match the chemical similarity of the reference peptide with the mutated sequence, we could optimize  
 65 further for a better docking score. We often started with simpler peptides, such as AAA and GGG,  
 66 since initial experiments seeding with the reference peptide and other peptides with high docking  
 67 scores did not result in any improvement beyond their initial scores. Additionally, we observed that  
 68 while keeping the computational capacity the same (Dual-core Intel i5, with 8GB memory), chemical  
 69 similarity between two peptides could be computed in microseconds, while docking took 1-3 minutes.  
 70 This observation supported our hypothesis from a computational cost/time standpoint.

## 71 2.4 Optimizing for non-canonical linear and macrocyclic peptides

72 To benchmark different approaches, we set up two distinct experiments, first, optimizing for similarity  
 73 and docking score, and another optimizing for docking score alone. In the first experiment, the  
 74 chemical similarity for the mutated peptide to the reference peptide was computed, and the selection  
 75 for further docking was based on the higher similarity between the seed or mutated peptide to the  
 76 reference. Similar to the similarity evaluation, if the mutated peptide had a higher docking score,  
 77 it was carried over to the next iteration, or else, the seed peptide was mutated again. In the second  
 78 experiment, there was no similarity optimization done. In both cases, the seed peptide was mutated  
 79 for at least 10 times, before being evaluated for docking. The genetic algorithm was run for 50  
 80 iterations, thus, a minimum of 500 mutations of the seed sequence was done for each experiment.

81 The chemical similarity was computed by evaluating Tanimoto similarity between the fingerprints of  
 82 the reference and mutated peptides [29]. We noted that similarity between the bit-based fingerprints,  
 83 with active and inactive bits for different molecular fragments, led to similarity score of 1.0 within  
 84 10-20 iterations. Such behavior arises from the bit-based fingerprints capturing individual fragments  
 85 alone, irrespective of their position and number of occurrences [30]. As an alternative, we used the  
 86 count-based fingerprints to compute chemical similarity. However, we noted that the docking scores  
 87 of the resulting peptides at the end of experiments based on different fingerprint types were similar.

88 We optimized the seed sequences for both linear and macrocyclic peptides with non-canonical amino  
89 acids to be similar to the reference peptide YWY with a score of -10.3 kcal/mol, and dock against  
90 DRD2 kinase protein. The best linear peptide was yGwFp with a docking score of -11.2 kcal/mol,  
91 while the best macrocyclic peptide was gvhw with a docking score of -10.4 kcal/mol (Figure 2). In the  
92 optimization run, only the sequence composition could change, while the topology, linear or cyclic,  
93 had to remain constant for all iterations. We pursued the optimization for both linear and macrocyclic  
94 peptides through the similarity-docking and docking-only objective functions, and benchmarking  
95 both bit-based and count-based fingerprints.

96 For a particular experiment involving non-canonical linear peptides, we set the seed to tripeptide aaa  
97 with a docking score of -7.0 kcal/mol, and the seed of the random generator to 0, and noted how the  
98 genetic algorithm proceeded. With the bit-based fingerprints, at 21 iterations, the mutated sequence,  
99 yMwFp, with a docking score of -11.2 kcal/mol, surpassed the docking score of the reference peptide  
100 sequence, and a Tanimoto similarity of 0.73. In further iterations, the sequence was mutated to yGwFp,  
101 although the docking score remained unchanged. With the count-based fingerprints, we noted that the  
102 docking score was surpassed at 3 iterations, with sequence yaNFY having a docking score of -10.7  
103 kcal/mol and a Tanimoto similarity of 0.54. Ultimately, the sequence got mutated to yaVFY with a  
104 docking score of -11.1 kcal/mol. In multiple experiments, with different seed sequences and different  
105 numbers for the random seed generator, we observed similar trends, with the docking score being  
106 surpassed sooner when the similarity was computed using count-based fingerprints. We attribute this  
107 observation to the finer granularity of count-based fingerprints, and the Tanimoto similarity thereof.

108 Using the docking-only objective function, without any reference peptide sequence, the seed sequence,  
109 aaa, could mutate to RFEaa with a docking score of -9.2 kcal/mol, which is worse than that of the  
110 reference sequence. Thus, this experiment underscores the need of a similarity objective function, to  
111 guide the evolution of the sequence towards a high-affinity sequence, and then optimize the docking  
112 even further.

### 113 3 Limitations, Future Work and Conclusion

114 Our work discussed the development of a genetic algorithm-based method to translate linear L-  
115 peptides, commonly obtained through high-throughput library screening, to non-canonical linear and  
116 macrocyclic peptides. With a dual objective function of similarity matching to a high-affinity reference  
117 peptide and maximization of the docking score against the target protein, we were able to obtain  
118 non-canonical peptides with better docking scores in less than 30 iterations. In the current state, the  
119 work is limited in its applicability to other target proteins. We aim to include more proteins and  
120 evaluate our approach in the near future. Additionally, we look forward to working on experimental  
121 data sets to translate L-peptides therein or optimize the non-canonical peptides further. We strongly  
122 believe that this approach can accelerate drug discovery efforts, by enabling the development of more  
123 stable peptides with similar or higher activity, which can ultimately benefit patients.

### 124 References

- 125 [1] Irene MA Nooren and Janet M Thornton. Diversity of protein–protein interactions. *The EMBO*  
126 *journal*, 22(14):3486–3492, 2003.
- 127 [2] Natia Tsomaia. Peptide therapeutics: targeting the undruggable space. *European journal of*  
128 *medicinal chemistry*, 94:459–470, 2015.
- 129 [3] Ziqing Qian, Patrick G Dougherty, and Dehua Pei. Targeting intracellular protein–protein  
130 interactions with cell-permeable cyclic peptides. *Current opinion in chemical biology*, 38:  
131 80–86, 2017.
- 132 [4] Lei Wang, Nanxi Wang, Wenping Zhang, Xurui Cheng, Zhibin Yan, Gang Shao, Xi Wang, Rui  
133 Wang, and Caiyun Fu. Therapeutic peptides: current applications and future directions. *Signal*  
134 *Transduction and Targeted Therapy*, 7(1):1–27, 2022.
- 135 [5] Norbert Sewald and Hans-Dieter Jakubke. *Peptides: chemistry and biology*. John Wiley &  
136 Sons, 2015.
- 137 [6] Sang-Hoon Joo. Cyclic peptides as therapeutic agents and biochemical tools. *Biomolecules &*  
138 *therapeutics*, 20(1):19–26, 2012.

- 139 [7] Knut Adermann, Harald John, Ludger Ständker, and Wolf-Georg Forssmann. Exploiting natural  
140 peptide diversity: novel research tools and drug leads. *Current opinion in biotechnology*, 15(6):  
141 599–606, 2004.
- 142 [8] Somesh Mohapatra, Nina Hartrampf, Mackenzie Poskus, Andrei Loas, Rafael Gomez-  
143 Bombarelli, and Bradley L Pentelute. Deep learning for prediction and optimization of fast-flow  
144 peptide synthesis. *ACS central science*, 6(12):2277–2286, 2020.
- 145 [9] Alexander J Mijalis, Dale A Thomas, Mark D Simon, Andrea Adamo, Ryan Beaumont, Klavs F  
146 Jensen, and Bradley L Pentelute. A fully automated flow-based approach for accelerated peptide  
147 synthesis. *Nature chemical biology*, 13(5):464–466, 2017.
- 148 [10] Eva M López-Vidal, Carly K Schissel, Somesh Mohapatra, Kamela Bellovoda, Chia-Ling Wu,  
149 Jenna A Wood, Annika B Malmberg, Andrei Loas, Rafael Gómez-Bombarelli, and Bradley L  
150 Pentelute. Deep learning enables discovery of a short nuclear targeting peptide for efficient  
151 delivery of antisense oligomers. *JACS Au*, 1(11):2009–2020, 2021.
- 152 [11] Carly K Schissel, Somesh Mohapatra, Justin M Wolfe, Colin M Fadzen, Kamela Bellovoda,  
153 Chia-Ling Wu, Jenna A Wood, Annika B Malmberg, Andrei Loas, Rafael Gómez-Bombarelli,  
154 et al. Deep learning to design nuclear-targeting abiotic miniproteins. *Nature chemistry*, 13(10):  
155 992–1000, 2021.
- 156 [12] Chien-Hsun Wu, I-Ju Liu, Ruei-Min Lu, and Han-Chung Wu. Advancement and applications  
157 of peptide phage display technology in biomedical science. *Journal of biomedical science*, 23  
158 (1):1–14, 2016.
- 159 [13] Gianni Cesareni, Luisa Castagnoli, and Gianluca Cestra. Phage displayed peptide libraries.  
160 *Combinatorial Chemistry and High Throughput Screening*, 2:1–18, 1999.
- 161 [14] Kristopher Josephson, Alonso Ricardo, and Jack W Szostak. mrna display: from basic principles  
162 to macrocycle drug discovery. *Drug Discovery Today*, 19(4):388–399, 2014.
- 163 [15] Terry T Takahashi, Ryan J Austin, and Richard W Roberts. mrna display: ligand discovery,  
164 interaction analysis and beyond. *Trends in biochemical sciences*, 28(3):159–165, 2003.
- 165 [16] Renaud Prudent, D Allen Annis, Peter J Dandliker, Jean-Yves Ortholand, and Didier Roche.  
166 Exploring new targets and chemical space with affinity selection-mass spectrometry. *Nature  
167 Reviews Chemistry*, 5(1):62–71, 2021.
- 168 [17] Anthony J Quartararo, Zachary P Gates, Bente A Somsen, Nina Hartrampf, Xiyun Ye, Arisa  
169 Shimada, Yasuhiro Kajihara, Christian Ottmann, and Bradley L Pentelute. Ultra-large chemical  
170 libraries for the discovery of high-affinity peptide binders. *Nature communications*, 11(1):1–11,  
171 2020.
- 172 [18] Benjamí Oller-Salvia and Jason W Chin. Efficient phage display with multiple distinct non-  
173 canonical amino acids using orthogonal ribosome-mediated genetic code expansion. *Ange-  
174 wandte Chemie*, 131(32):10960–10964, 2019.
- 175 [19] Titia Rixt Oppewal, Ivar D Jansen, Johan Hekelaar, and Clemens Mayer. A strategy to select  
176 macrocyclic peptides featuring asymmetric molecular scaffolds as cyclization units by phage  
177 display. *Journal of the American Chemical Society*, 144(8):3644–3652, 2022.
- 178 [20] Richard Obexer, Louise J Walport, and Hiroaki Suga. Exploring sequence space: harnessing  
179 chemical and biological diversity towards new peptide leads. *Current opinion in chemical  
180 biology*, 38:52–61, 2017.
- 181 [21] Jr Nestor et al. The medicinal chemistry of peptides. *Current medicinal chemistry*, 16(33):  
182 4399–4418, 2009.
- 183 [22] Miguel García-Ortegón, Gregor NC Simm, Austin J Tripp, José Miguel Hernández-Lobato,  
184 Andreas Bender, and Sergio Bacallado. Dockstring: easy molecular docking yields better  
185 benchmarks for ligand design. *Journal of chemical information and modeling*, 62(15):3486–  
186 3502, 2022.
- 187 [23] Oleg Trott and Arthur J Olson. Autodock vina: improving the speed and accuracy of docking  
188 with a new scoring function, efficient optimization, and multithreading. *Journal of computational  
189 chemistry*, 31(2):455–461, 2010.
- 190 [24] Robert C Edgar and Serafim Batzoglou. Multiple sequence alignment. *Current opinion in  
191 structural biology*, 16(3):368–373, 2006.

- 192 [25] Kostas Papadopoulos, Kathryn A Giblin, Jon Paul Janet, Atanas Patronov, and Ola Engkvist.  
193 De novo design with deep generative models based on 3d similarity scoring. *Bioorganic &*  
194 *Medicinal Chemistry*, 44:116308, 2021.
- 195 [26] Didier Rognan. Binding site similarity search to identify novel target–ligand complexes.  
196 *Computational Chemogenomics*, page 171, 2013.
- 197 [27] Greg Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry,  
198 and predictive modeling. *Greg Landrum*, 2013.
- 199 [28] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical*  
200 *information and modeling*, 50(5):742–754, 2010.
- 201 [29] Dávid Bajusz, Anita Rácz, and Károly Héberger. Why is tanimoto index an appropriate choice  
202 for fingerprint-based similarity calculations? *Journal of cheminformatics*, 7(1):1–13, 2015.
- 203 [30] Yuan Wang. *Molecular Complexity Effects and Fingerprint-Based Similarity Search Strategies*.  
204 PhD thesis, Universitäts-und Landesbibliothek Bonn, 2009.

## 205 Checklist

- 206 1. For all authors...
- 207 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
208 contributions and scope? [Yes]
- 209 (b) Did you describe the limitations of your work? [Yes]
- 210 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 211 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
212 them? [Yes]
- 213 2. If you are including theoretical results...
- 214 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 215 (b) Did you include complete proofs of all theoretical results? [N/A]
- 216 3. If you ran experiments...
- 217 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
218 mental results (either in the supplemental material or as a URL)? [No] Code can be  
219 shared upon the final publication of the work. Data and instructions have been provided  
220 to reproduce the work.
- 221 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
222 were chosen)? [Yes]
- 223 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
224 ments multiple times)? [N/A]
- 225 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
226 of GPUs, internal cluster, or cloud provider)? [Yes]
- 227 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 228 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 229 (b) Did you mention the license of the assets? [N/A]
- 230 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 231
- 232 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
233 using/curating? [N/A]
- 234 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
235 information or offensive content? [N/A]
- 236 5. If you used crowdsourcing or conducted research with human subjects...
- 237 (a) Did you include the full text of instructions given to participants and screenshots, if  
238 applicable? [N/A]
- 239 (b) Did you describe any potential participant risks, with links to Institutional Review  
240 Board (IRB) approvals, if applicable? [N/A]
- 241 (c) Did you include the estimated hourly wage paid to participants and the total amount  
242 spent on participant compensation? [N/A]