
Peptide Featurization Methods in Drug Discovery

Grace Adalia Rico

Master of Science in Bioinformatics
Bioinformatics Research Center
Aarhus University

June 2026

In collaboration with Raven Biosciences

Abstract

Predicting the biological properties of peptides from sequences is a central problem in drug discovery *in silico*, and the quality of those predictions depends critically on the numerical representation provided by featurizers. This thesis benchmarks new and existing featurization methods spanning small molecule descriptors, peptide specific encodings, protein language model embeddings, and structure-based descriptors. These are evaluated on major histocompatibility complex binding affinity and hemolysis tasks. The novel methods include a circular fingerprint adaptation and a semi-supervised self-distillation model. Of the featurizers tested, small molecule and simple peptide featurizers performed on par with deep learning methods, contrary to expectations that deep learning methods would dominate. Featurizers were less successful on binding affinity than hemolysis which aligns with the comparative complexities of these two tasks. The novel deep learning model tested here was as successful as, if not more successful than other featurizers. Findings suggest that featurizers' performance is task dependent, and there remain tasks such as binding affinity where even the most sophisticated models are yet to capture the relevant features of peptides.

Acknowledgements

I would like to thank my colleagues and mentors at Raven Biosciences for their guidance and expertise. I would especially like to thank Lucy for her insight and support.

Contents

Abstract	i
Acknowledgements	ii
1 Introduction	1
2 Methods	6
2.1 Data	6
2.2 JADBio - Featurizer Evaluation	9
2.3 Featurization	10
2.3.1 Small Molecule Featurizers	11
2.3.2 1D and 2D Peptide Featurizers	12
2.3.3 Deep Learning Featurizers	14
3 Results and Discussion	22
3.1 Small-Molecule Featurizers	24
3.2 1D and 2D Peptide Featurizers	25
3.3 Deep Learning Featurizers	26
3.3.1 DINO + SupCon	27
3.4 Case Study: Attention Analysis	31
3.5 General Discussion	36
3.5.1 Task Dependence of Featurizer Rankings	36
3.5.2 Limitations	36
3.5.3 Future Directions	39
4 Conclusion	40
A Appendix	49
A.1 Use of Generative Artificial Intelligence	52