

AmyloGram: prediction of amyloid sequences in R

Jarosław Chilimoniuk^{1*}, Michał Burdukiewicz², Piotr Sobczyk³, Stefan Rödiger⁴,
Małgorzata Kotulska⁵ and Paweł Mackiewicz¹

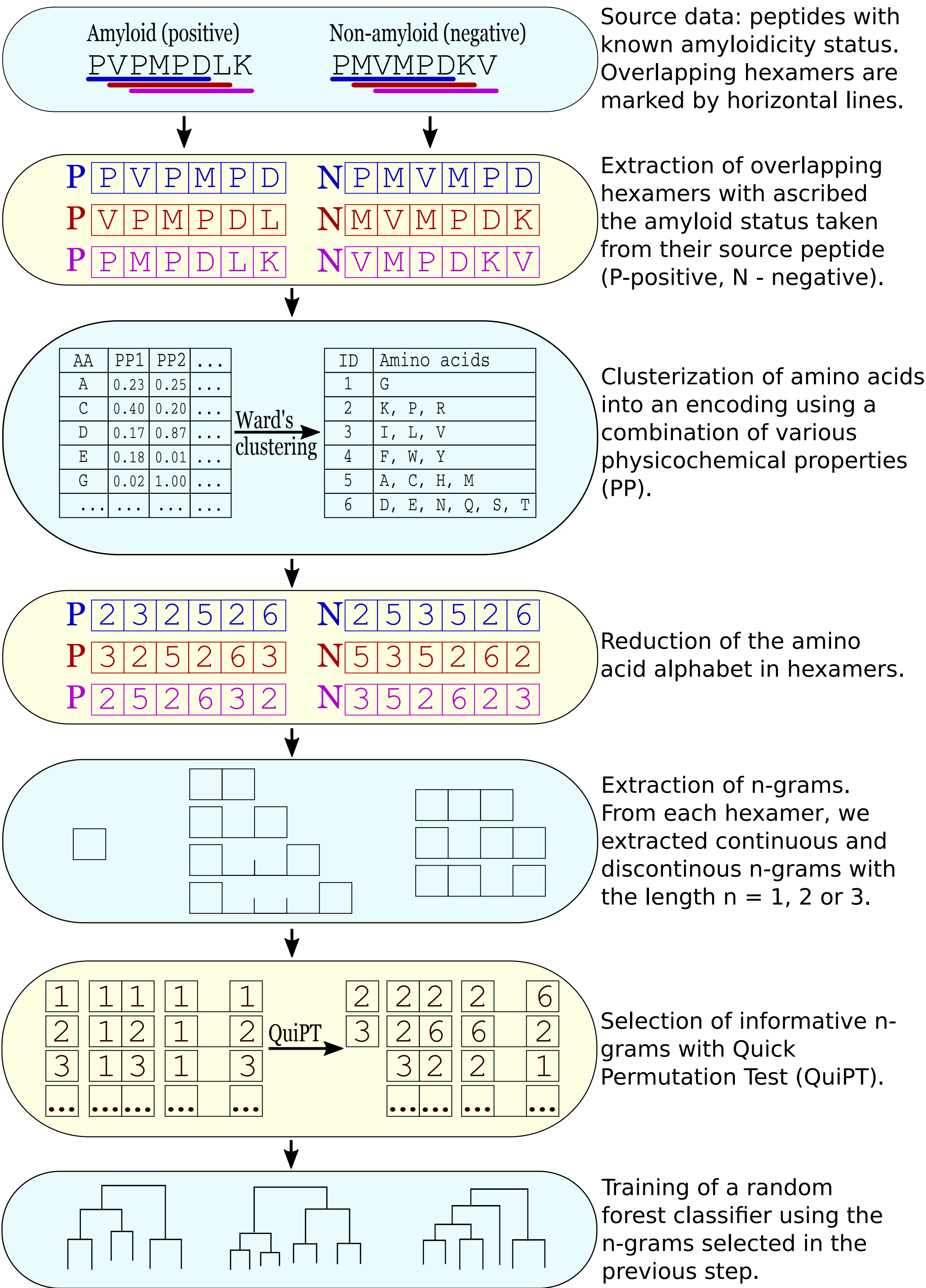
*jaroslaw.chilimoniuk@gmail.com ◇ jarochi@github

¹University of Wrocław, Department of Bioinformatics and Genomics, ²Warsaw University of Technology, Faculty of Mathematics and Information Science, ³Wrocław University of Science and Technology, Faculty of Pure and Applied Mathematics, ⁴Brandenburg University of Technology Cottbus-Senftenberg, Institute of Biotechnology, ⁵Wrocław University of Science and Technology, Department of Biomedical Engineering

Introduction

Aggregates of amyloid proteins are causes of neurodegenerative disorders. Using easy interpreting features we trained a novel algorithm, AmyloGram, for detection of amyloids. In comparison to other methods predicting amyloids our software achieved the highest performance. Additionally, we have experimentally confirmed that AmyloGram is able to detect false negatives (amyloid proteins wrongly annotated as non-amyloids) in its training data set.

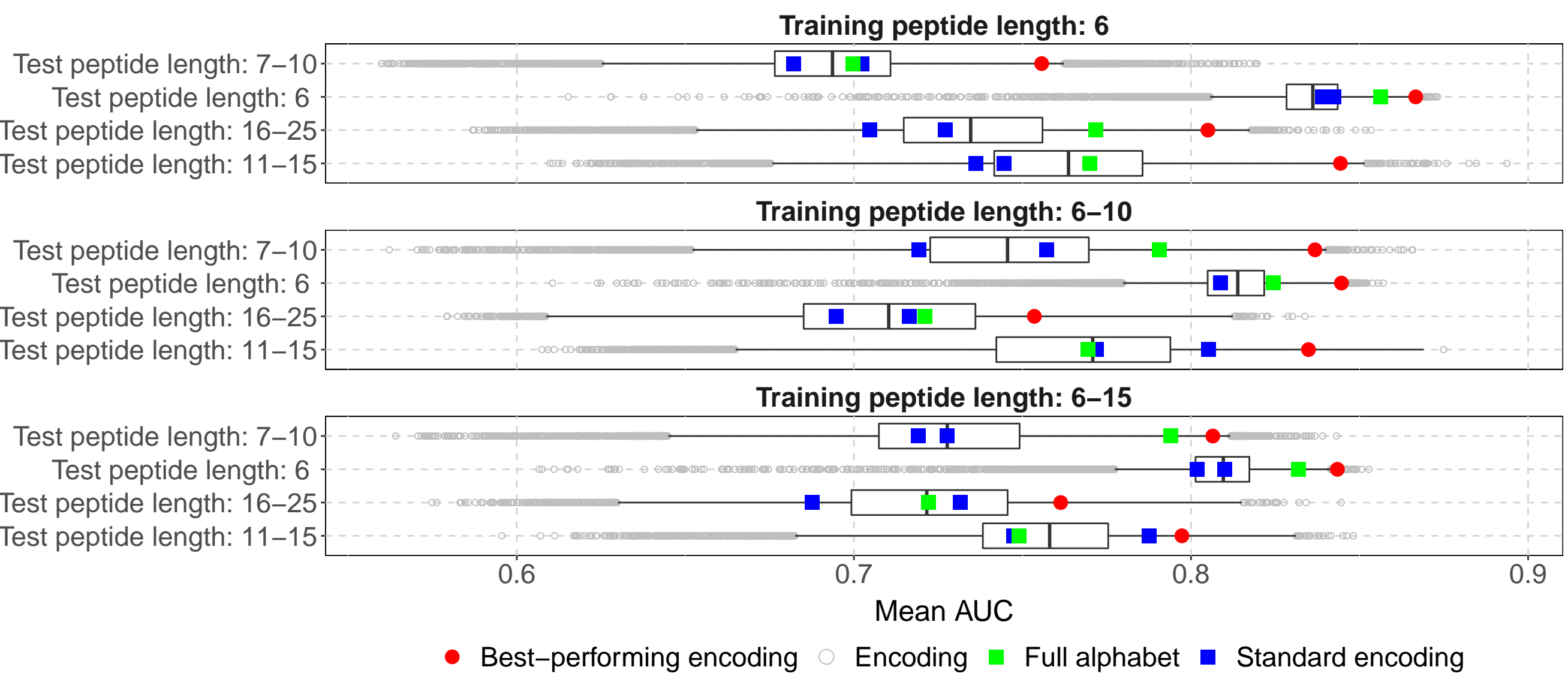
Scheme



Results of cross-validation

The amyloidogenicity of a given peptide may not depend on the exact sequence of amino acids but on its more general properties. Henceforth, we created 524,284 amino acid reduced alphabets (from three to six letters) based on physicochemical properties relevant to amyloidogenicity.

Distribution of mean AUC values of classifiers with various encodings for every possible combination of training and testing data set including different lengths of sequences.

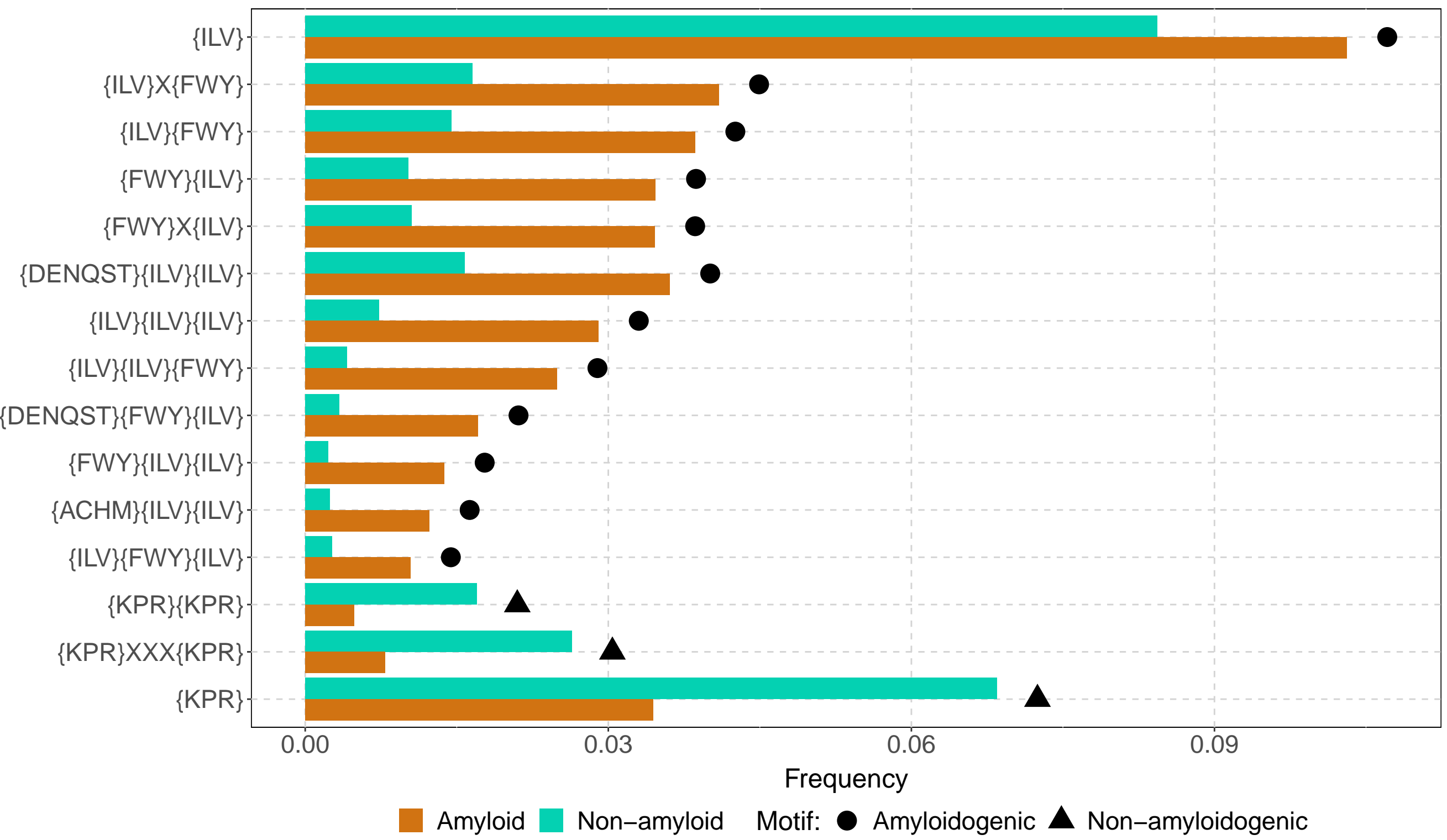


The predictor based on the best-performing encoding reached the highest AUC (0.8667) in classification of the shortest sequences (with the length of 6 residues).

Classifiers based on the full (i.e., unreduced) amino acid alphabet never predicted amyloidogenicity better than the best classifier based on the reduced alphabet.

The standard encodings found in the literature performed worse than other analyzed encodings in most categories.

Informative n-grams



The frequency of important n-grams used by the best-performing classifier in amyloid and non-amyloid sequences. The elements of n-grams are amino acids encoded using the best-performing reduced amino acid alphabet. X represents any amino acid. Dots and triangles denote n-grams occurring in motifs found in respectively amyloidogenic and non-amyloidogenic sequences (Paz and Serrano, 2004).

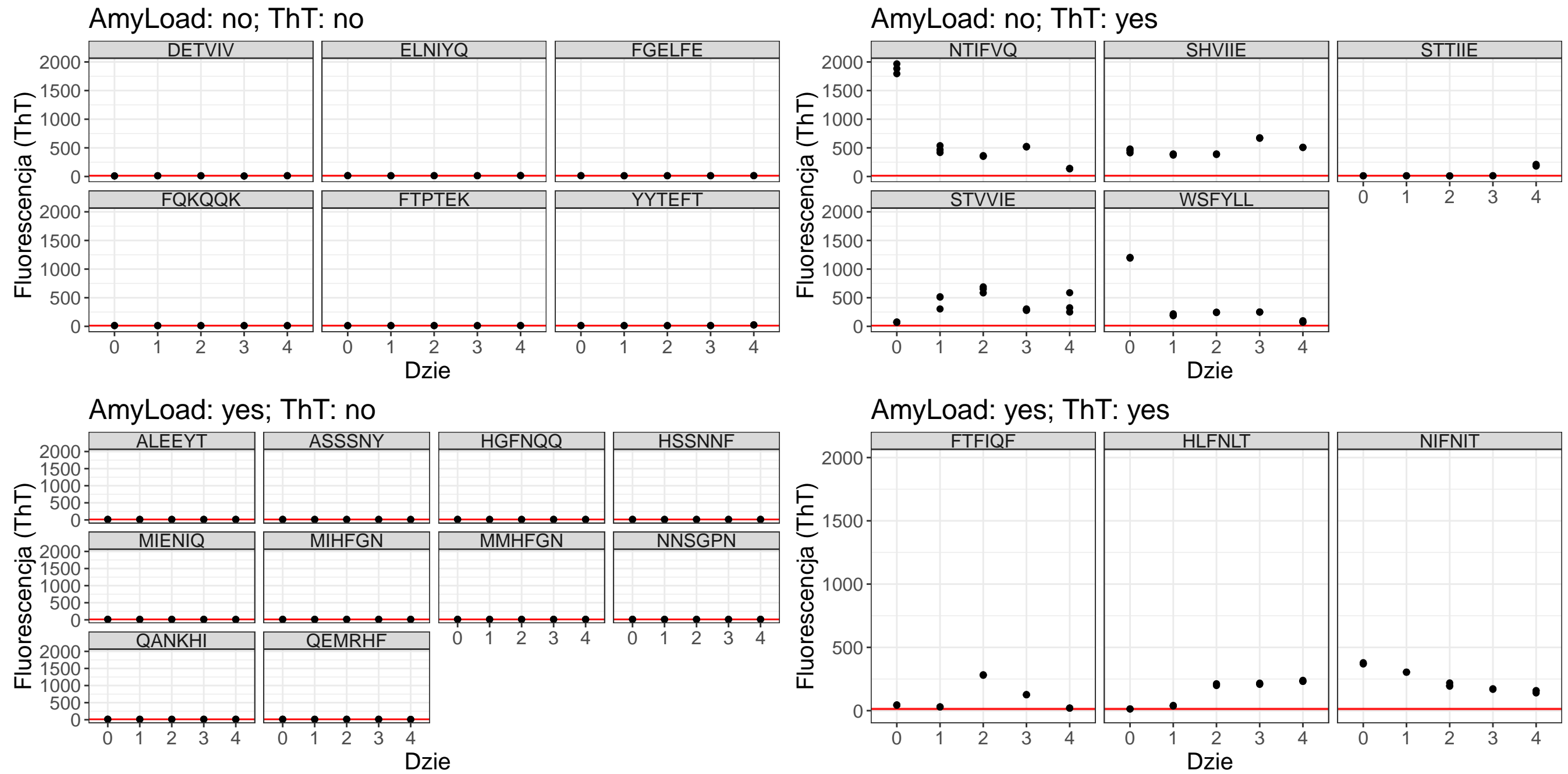
Benchmark results

Classifier	AUC	MCC	Sens.	Spec.
AmyloGram	0.8972	0.6307	0.8658	0.7889
PASTA (Walsh et al., 2014)	0.8550	0.4291	0.3826	0.9519
FoldAmyloid (Garbuzynskiy et al., 2010)	0.7351	0.4526	0.7517	0.7185
APPNN (Família et al., 2015)	0.8343	0.5823	0.8859	0.7222

The predictor based on the best-performing alphabet, called AmyloGram, was benchmarked against the most popular tools for the detection of amyloid peptides using an external data set *pep424*.

Experimental validation

Using AmyloGram we analyzed all 34 peptides from AmyLoad database. 23 peptides were predicted by AmyloGram otherwise than described in AmyloGram. Hehamers were validated experimentally using ThT assay.



15 out of 23 peptides had the same amyloid properties as predicted by AmyloGram.

Summary and funding

Thanks to the reduction of the amino acid alphabet and description of peptides by short sub-sequences (n-grams), we were able to create the efficient predictor of amyloidogenic sequences called AmyloGram.

Our software is available as a web-server:

www.smorfland.uni.wroc.pl/shiny/AmyloGram/ and R package:

<https://cran.r-project.org/package=AmyloGram>.

This research was partially funded by the KNOW Consortium and National Science Center (2015/17/N/NZ2/01845).

Bibliography

- Família, C., Dennison, S. R., Quintas, A., and Phoenix, D. A. (2015). Prediction of Peptide and Protein Propensity for Amyloid Formation. *PLOS ONE*, 10(8):e0134679.
- Garbuzynskiy, S. O., Lobanov, M. Y., and Galzitskaya, O. V. (2010). FoldAmyloid: a method of prediction of amyloidogenic regions from protein sequence. *Bioinformatics (Oxford, England)*, 26(3):326–332.
- Paz, M. L. d. I. and Serrano, L. (2004). Sequence determinants of amyloid fibril formation. *Proceedings of the National Academy of Sciences*, 101(1):87–92.
- Walsh, I., Seno, F., Tosatto, S. C. E., and Trovato, A. (2014). PASTA 2.0: an improved server for protein aggregation prediction. *Nucleic Acids Research*, page gku399.