

Fig. 1. The scheme of n-gram extraction. A) Input data - peptides with a known amyloidogenicity status. B) Each peptide sequence was divided into overlapping hexamers. The amyloidogenicity status of a source sequence was used as the amyloidogenicity status of a derived hexamer. C) From each hexamer we extracted continuous 1-, 2- and 3-grams. We selected also gapped 2-grams with the length of gap equal from 1 to 3 residues and gapped 3-grams with a single gap between the first and the second or the second and the third element of the n-gram.

of hot spots for non-amyloid aggregation. Recently, AGGRESCAN3D has been proposed to estimate more accurately aggregation propensity by performing 3D structure based analysis ?.

2 Methods

2.1 Data set

The data used in the study was extracted from AmyLoad data base (Wozniak and Kotulska, 2015). Aside from eight sequences shorter than five residues that were removed from the final data set, we obtained 418 amyloidogenic sequences and 1039 non-amyloidogenic sequences (1457 peptides in total).

Sequences shorter than 6 amino acids and longer than 25 amino acids were removed from the data set. The former were too short to be processed in the devised analysis framework and the latter were too diversified and rare, hampering the proper analysis.

The final data set contained 397 amyloidogenic and 1033 non-amyloidogenic sequences (1430 peptides in total).

2.2 Encodings of amino acids

The amyloidogenicity of a given peptide may not depend on the exact sequence of amino acids, but on its more general properties. To verify this hypothesis, we created 524 284 reduced amino acid alphabets (encodings) with different lengths (from three to six letters) using Ward's clusterization (Jr, 1963) on the selected physicochemical properties from AAIndex database (Kawashima *et al.*, 2008). We handpicked several measures belonging to more general categories important in the amyloidogenicity, such as size, hydrophobicity, solvent surface area, frequency in β -sheets and contactivity. As a rule of thumb, we limited it to properties introduced after 1980 when, thanks to the technological advancements, the measurements were more accurate.

The majority of encodings had at least one duplicate. In such a case, only a single representative was included in the cross-validation. After filtering duplicates, we obtained 18 535 unique encodings.

Since highly correlated or, contrarily, uncorrelated measures create very similar encodings, we further reduced the number of properties to 17 by selecting measures with the absolute value of Pearson's correlation coefficient larger than 0.95 for normalized values.

2.3 Training of learners

During the training phase, we extracted overlapping hexamers from each sequence. Each hexamer was tagged with the same etiquette (amyloid/nonamyloid) as the source peptide. For example, the amyloidogenic sequence of length 6 residues yields only one hexamer tagged as "amyloid". The non-amyloidogenic sequence of 8 residues yields 3 hexamers, all marked as "non-amyloids". (Fig. 1 A and B). Assuming that in longer amyloids only a short part of the sequence is responsible for amyloidogenicity, our method might result in many false positives in the training data set. To evade this problem, we restricted the maximum length of peptides in training data set to fifteen amino acids to easy the extraction of probable hot-spots.

From each hexamer we extracted n-grams of the following length: 1, 2 and 3. In the case of 2- and 3-grams, we separately counted both gapped and continuous n-grams. For 2-grams we counted n-grams with gaps of lengths from 1 to 3 and for 3-grams with a single gap between the first and the second or the second and the third element (see Fig. 1).

To study the length of amyloidogenic signal, we trained three classifiers for each encoding on the sequences of different lengths. We considered separately six-residue sequences, shorter or equal to ten residues and shorter or equal to fifteen residues.

All n-grams extracted from the hexamers in the training data set were filtered using described below Quick Permutation Test with the information gain (mutual information) as the criterion of the importance of a specific n-gram. In the next step, we used n-grams with the p-value smaller than 0.05 to build a random forest (Breiman, 2001) classifier using ranger R package (Wright and Ziegler, 2015).

Furthermore, we repeated the procedure described above on two typical reduced alphabets of amino acids derived from the literature to check if the process of amyloidogenicity does require nonstandard groupings of amino acids. We also added the full amino acid alphabet to assess the advantages of the amino acid encoding.

2.4 Quick Permutation Test (QuiPT)

The permutation tests, commonly used for filtering important n-grams, are often one of the most limiting factors for these kinds of analysis requiring significant computational power. The Quick Permutation Test effectively filters n-gram features without performing a huge number of permutations. Let us consider the contingency table for target y and feature x :

target / feature	1	0
1	$n_{1,1}$	$n_{1,0}$
0	$n_{0,1}$	$n_{0,0}$

Under the hypothesis that x and y are independent, the probability of such a contingency table is given by the multinomial distribution. In the permutation test we simply reshuffle feature and target labels, but the total number of positives in both of them is fixed. When we impose this constraint on the multinomial distribution then it depends only on $n_{1,1}$ and is fairly easy to compute. Observe also, that $n_{1,1}$ is from range $[0, \min(n_{1,\cdot}, n_{\cdot,1})]$. After computing Information Gain for each possible value of $n_{1,1}$, we get distribution of Information Gain under hypothesis that target and feature are independent. We reject null hypothesis when Information Gain is big.

Having the analytical formula for the distribution, allows us to perform permutation test much quicker. Furthermore, we get exact quantiles even for extreme tails of the distribution, which is not guaranteed by random permutations. In fact, imagine performing the test with $\alpha = 10^{-8}$, which is not an uncommon value, i.e. when we adjust for multiple testing. Even for huge number of permutations like $m = 10^8$, the standard deviation of quantile estimate in permutation test, $\frac{p(1-p)}{m}$, is roughly equal to α itself.

Category	Property
Contactivity	Average flexibility indices (Bhaskaran and Ponnuswamy, 1988)
Contactivity	14 A contact number (Nishikawa and Ooi, 1986)
Contactivity	Accessible surface area (Radzicka and Wolfenden, 1988)
Contactivity	Buriability (Zhou and Zhou, 2004)
Contactivity	Values of Wc in proteins from class Beta, cutoff 12 A, separation 5 (Wozniak and Kotulska, 2014)
Contactivity	Values of Wc in proteins from class Beta, cutoff 12 A, separation 15 (Wozniak and Kotulska, 2014)
β -frequency	Average relative probability of inner beta-sheet (Kanehisa and Tsong, 1980)
β -frequency	Relative frequency in β -sheet (Prabhakaran, 1990)
β -frequency	Thermodynamic β -sheet propensity (Kim and Berg, 1993)
Hydrophobicity	Hydrophobicity index (Argos <i>et al.</i> , 1982)
Hydrophobicity	Optimal matching hydrophobicity (Sweet and Eisenberg, 1983)
Hydrophobicity	Hydrophobicity-related index (Kidera <i>et al.</i> , 1985)
Hydrophobicity	Scaled side chain hydrophobicity values (Black and Mould, 1991)
Polarity	Polarizability parameter (Charton and Charton, 1982)
Polarity	Mean polarity (Radzicka and Wolfenden, 1988)
Size	Average volumes of residues (Pontius <i>et al.</i> , 1996)
Stability	Side-chain contribution to protein stability (kJ/mol) (Takano and Yutani, 2001)

Table 1. Physicochemical properties used in the creation of reduced amino acid alphabets.

In the context of *k*-mer data we can speed up our algorithm even further. Note that since target *y* is common for testing all *k*-mer features, test statistics depends only on $n_{\cdot,1}$ – the number of positive cases in a feature. Though we test millions of features, there are just few distributions that we need to compute, as usually number of positives in *k*-mer feature is small. We take advantage of this fact and we compute quantiles for just a handful of distributions. Therefore complexity of our algorithm is roughly equal $O(n \cdot p)$.

Lastly, let us point out that QuiPT is very similar for Fisher’s exact test. From the derivation provided in i.e. (?), it becomes obvious that QuiPT is a heuristics for an unsolved problem of a two-tailed Fisher’s exact test. In this heuristics, extremity of a contingency table, is defined by the amount of Information Gain it provides.

2.5 Cross-validation and selection of the best-performing encoding

The ability to correctly predict amyloidogenicity was assessed during the five-fold cross-validation. Since the data set was very heterogeneous, we repeated the cross-validation fifteen times for each classifier to obtain more precise estimates of performance measures for each classifier.

To evaluate if our classifiers are able to use decision rules extracted from sequences of given length to correctly classify longer or shorter sequences, we calculate performance measures separately for four ranges of lengths of sequences: 6, 7-10, 11-15 and 16-25. The number of sequences from the given length range was roughly comparable between folds of cross-validation.

To choose the most adequate reduced amino acid alphabet, we ranked the values of Area under Curve - AUC (with rank 1 for the best AUC, rank 2 for the second best AUC and so on) for each range of sequence length in the testing data set. The encoding with the lowest sum of ranks from all sequence length categories was selected as the best one. For this encoding, we choose the range of peptides length in the training set providing the best AUC in cross-validation.

2.6 Encoding distance

The encoding distance is a measure defining the similarity between two encodings. It has zero value for identical encodings and grows with the differences between encodings. It was introduced to verify if the reduced

alphabets very similar to the best-performing encoding will also have good prediction performance.

We define the encoding distance as the minimum number of amino acids that have to be moved between subgroups of encoding to make *a* identical to *b* (the order of subgroups in the encoding and amino acids in a group is unimportant). This measure is further scaled by a factor reflecting how much moving amino acids between groups altered mean group properties.

To compute the scale factor *s* for the encoding distance between encoding *a* with *n* subgroups and encoding *b* with *m* subgroups we firstly calculate p_i and p_j , mean values of physicochemical properties of all amino acids separately for each subgroup. The factor between *a* and *b* is equal to:

$$s_{ab} = \sum_{i=1}^n \left(\min_{j=1, \dots, m} \left(\sqrt{\sum_{i=1}^l p_i^2} - \sqrt{\sum_{j=1}^l p_j^2} \right) \right)$$

where *l* is equal to the number of physicochemical properties of concern.

2.7 Benchmark of AmyloGram

The best-performing reduced amino acid alphabet chosen during the cross-validation was later used to train AmyloGram, *n*-gram based predictor of amyloidogenicity.

We used *pep424* data set (Walsh *et al.*, 2014) to compare the performance of AmyloGram and other predictors of amyloidogenicity. Since the model of AmyloGram does not cover peptides shorter than five amino acids, we removed them from the total benchmark data set. It should have not affect the comparison as only five sequences were eliminated (around 1% of the original data set). Additionally, we also benchmarked three predictors learned on the *n*-grams extracted from sequences of different length ranges without any amino acid encoding.

All benchmarked classifiers were trained on sequences used during the cross-validation. Since some peptides were common for both *pep424* and AmyLoad, we removed them from the training data set. After purification, the learning data set had 269 positive sequences and 746 negative sequences longer than five residues and shorter than fifteen residues. Aside from the preparation of the training data, we exactly repeated the procedure of *n*-gram extraction as described above (see Fig. 1).

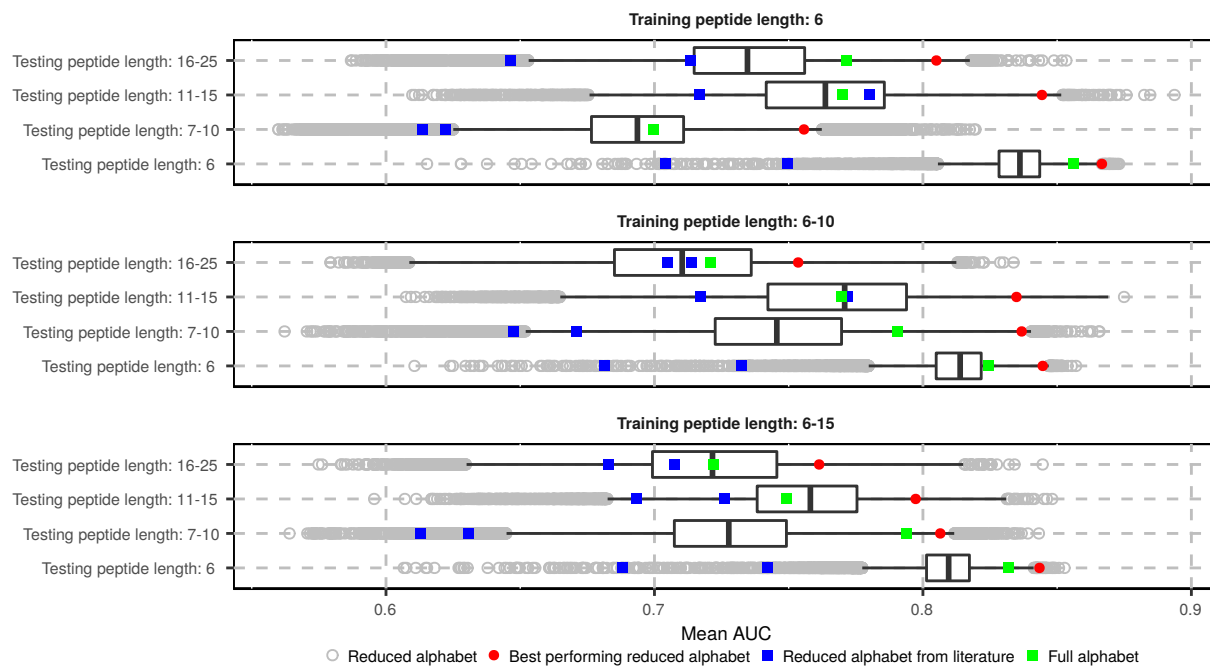


Fig. 2. Distribution of AUC values of different reduced amino acid alphabets for different lengths of sequences in the training and testing data sets.

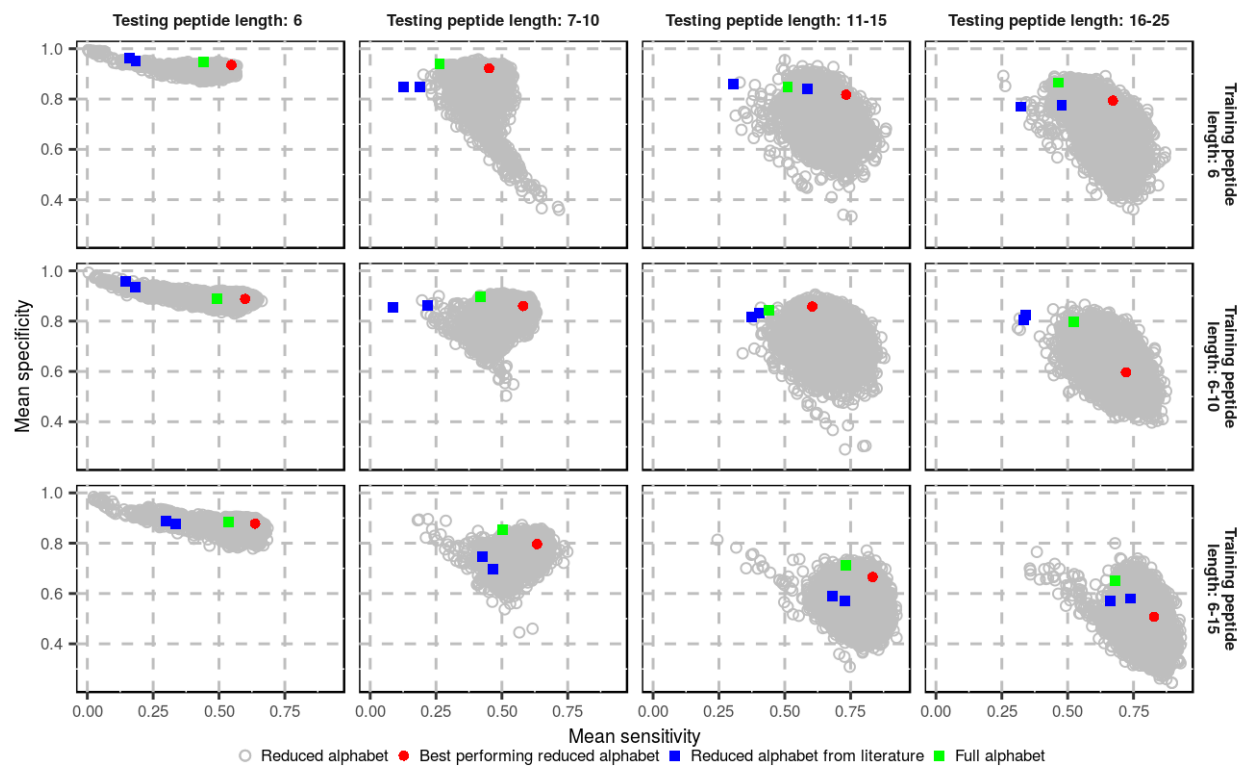


Fig. 3. Sensitivity and specificity of classifiers in cross-validation for different lengths of sequences in the training and testing data sets. The classifier based on the best-performing encoding always have good specificity and sensitivity.

3 Results and discussion

3.1 Selection of the best-performing encoding

The selected encoding performed better than other reduced alphabets considering all sequence length ranges in training and testing data set.

It had the AUC always in the fourth quartile (Fig. 2). For the best-performing encoding the most problematic was correct prediction of the amyloidogenicity in the longest peptide data set 16-25, not when learned

on very short sequences, but for longer ones (6-10 and 6-15). Such a behavior was typical for most of the analyzed encodings.

The highest values of AUC the chosen encoding reaches in the relatively the simplest cases of predicting amyloidogenicity of the shortest sequences (6 residues). Of course, the decision rules were the easiest to extract when also the learning data set had only sequences of the same lengths, but the results were quite comparable even if the training set contained longer sequences. Also in this situation, the majority of reduced amino acid alphabets behaved like the selected encoding and also had higher AUC than in other scenarios.

Considering only the AUC value, the most problematic testing data were sequences longer than 15 amino acids. Surprisingly, the classifiers trained on the shortest sequences available (six residues) were performing the best. That might indicate that our n-gram approach extracts the important features the better, the shorter are sequences in training data set. For example, in the amyloidogenic peptide of length 15, only a very specific region of residues might be responsible for the creation of harmful aggregates. In this case, when in our framework overlapping n-grams are extracted, only part of them may carry the true signal of amyloidogenicity, but all are marked as amyloids. Despite this problem, the overall prediction of classifiers learned on the long sequences was still adequate, with the median values of AUC higher than 0.7 for every testing set.

In addition to the high AUC, the best encoding had also very good sensitivity and specificity regardless of the length of sequences in training and testing set (Fig. 3). The classifiers trained on the peptides of length 6 tend to have the best specificity, while predictors learned on the long sequences have the best sensitivity. Albeit the classifiers trained on six-residue-long sequences have generally better AUC, training on the sequences ranging from six to ten residues seem to yield the most balanced classifiers with optimal sensitivity and specificity.

We also evaluated classifiers based on the full amino acid alphabet. In most cases they were also placed in the fourth quartile considering their AUC value (Fig. 2). Nevertheless, they never predicted amyloidogenicity better than the selected encoding. It suggests that despite the predictive power of n-grams based on full alphabet, the encoding allows better generalization of the prediction rules and in the consequence a better performance.

Similarly to the best-performing encoding, the sensitivity of classifiers based on the full amino acid alphabet decreased with the length of sequences in the training data set (Fig. 3). Furthermore, these classifiers always seemed to have one of the worst sensitivities among all analyzed predictors, especially when tested on longer amyloids. It means that using the full amino acid alphabet it was easier to point the non-amyloidogenic sequences instead of recognizing amyloids.

The encodings found in the literature performs substantively worse than other analyzed amino acid alphabets in all categories. It indicates that classical divisions of amino acids do not create groups suitable for the recognition of amyloids. This observation is well-supported by the specificity-sensitivity plot (Fig. 3), where classifiers trained with this encodings of amino acids groups with the worst performers.

3.2 The best-performing encoding and important n-grams

The reduced amino acid alphabet chosen during the analysis has six subgroups (Tab. 2). Two subgroups, 3 and 4 contain strongly hydrophobic amino acids, while 1 and 5 are mostly hydrophilic. In addition to this, subgroup 1 and 4 have amino acid with the highest average flexibility. The least flexible amino acids are in subgroup 5. The most polarizable amino acids belong to subgroup 4, while the lowest polarizability is typical for glycine, the only amino acid in subgroup 1. The glycine has also the highest thermodynamic beta sheet propensity, which has the lowest values for subgroups 3 and 4.

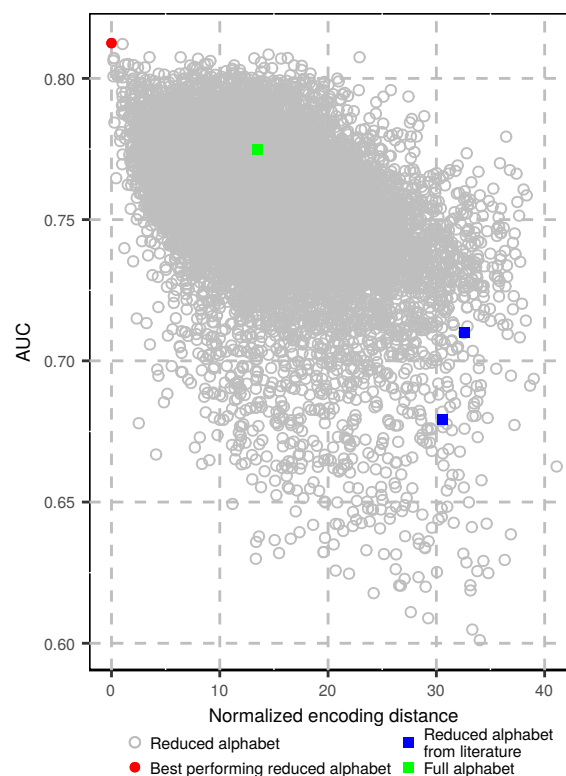


Fig. 4. The encoding distance and AUC of reduced alphabets studied in the cross-validation. Classifiers with the smallest encoding distance to the best classifier have the highest AUC.

Table 2. The best-performing encoding.

Subgroup ID	Amino acids
1	G
2	K, P, R
3	I, L, V
4	F, W, Y
5	A, C, H, M
6	D, E, N, Q, S, T

In total, eleven combinations of physicochemical properties created the best performing reduced amino acid alphabet. Only four features appeared in all combinations: hydrophobicity index (Argos *et al.*, 1982), average flexibility indices (Bhaskaran and Ponnuswamy, 1988), polarizability parameter (Charton and Charton, 1982) and thermodynamic β -sheet propensity (Kim and Berg, 1993).

We calculated encoding distances between the best-performing reduced amino acid alphabet and other encodings (Fig. 4). To compute the scale factor, we used described above features common for all encodings identical to the best-performing one. The value of AUC is significantly lower for more distant encodings (-0.4370 Pearson's correlation coefficient, p -value smaller than 2.2×10^{-16}). Such relationship between AUC and encoding distance confirms that only reduced alphabets sufficiently similar to the best-performing encoding are able to remove unnecessary diversity while preserving information important for recognition of amyloids.

We selected 65 n-grams that have p -values smaller than 0.05 in the all folds in all repetitions of cross-validation regardless of the length of sequences in the training set (see Fig. 5). The frequency of the n-grams

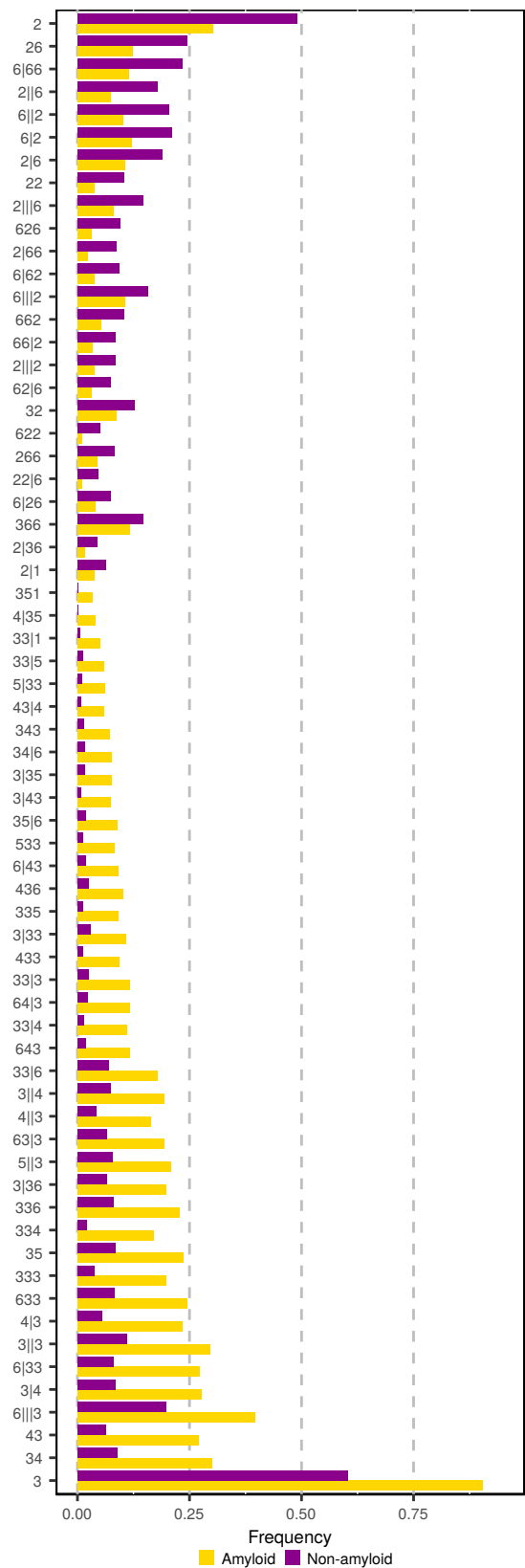


Fig. 5. The frequency of important n-grams used by the best-performing classifier in amyloid and non-amyloid sequences. The elements of n-grams are amino acids encoded using the best-performing reduced amino acid alphabet (see Tab. 2). A vertical bar represents a gap in a n-gram.

was computed for all sequences derived from AmyLoad database. The n-grams typical for amyloidogenic sequences (with the highest frequency in amyloids) incorporate mostly highly hydrophobic, typical for β -structures amino acids from subgroups 3 and 4. The important n-grams occurring frequent in amyloids often have repeats of 3, suggesting that the presence of amino acids belonging to this subgroup might be one of the most effective predictors of amyloidogenicity.

N-grams typical for non-amyloidogenic peptides have mostly elements belonging to subgroups 2 and 6. The amino acids of these subgroups are strongly hydrophilic and highly flexible, in the opposite to the residues typical for amyloids.

3.3 Benchmark of AmyloGram

Table 3. Results of benchmark on pep424 data set for AmyloGram, PASTA2, FoldAmyloid and random forest predictor learned on n-grams extracted without any amino acid encoding from the sequences of the length specified in the brackets.

Classifier	AUC	MCC	Sensitivity	Specificity
AmyloGram	0.8972	0.6307	0.8658	0.7889
PASTA2	0.8550	0.5227	0.7987	0.7444
FoldAmyloid	0.7351	0.4526	0.7517	0.7185
full alphabet (6)	0.8411	0.5427	0.4966	0.9593
full alphabet (6-10)	0.8581	0.5698	0.7517	0.8259
full alphabet (6-15)	0.8610	0.5490	0.8188	0.7519

The benchmark covered Amylogram as well as two best-performing peer-reviewed predictors of amyloidogenicity (PASTA2 (Walsh *et al.*, 2014) and FoldAmyloid (Garbuzynskiy *et al.*, 2010)). We analyzed Area Under the Curve (AUC), Matthew’s Correlation Coefficient (MCC), Sensitivity and Specificity (see Tab. 3).

Interestingly, the n-gram extraction method was efficient enough to produce classifiers able to outperformed published methods. Two of three n-gram based classifiers trained on the full alphabet have AUC higher than PASTA2 and all three were more successful than FoldAmyloid. They also maintained they high Specificity as seen previously during cross-validation.

Although the proposed n-gram extraction creates efficient classifiers, the encoding of amino acids further increases the prediction of amyloidogenicity. AmyloGram has the highest AUC, MCC and Sensitivity among all tested classifiers. Is has lower specificity than two classifiers trained on the full alphabet, but still outperforms other published method in this category. It is important to highlight that AmyloGram is the most of analyzed classifiers, having the best Specificity/Sensitivity trade-off, as indicated by the value of MCC.

4 Conclusion

Thanks to the reduction of amino acid alphabet, we were able to create the efficient predictor of amyloidogenic sequences called AmyLoad. One of the strength of our approach is its highly interpretable outcome, which hopefully sheds new light on the process of amyloid aggregation.

The idea of using the reduced amino acid alphabets is not new, but we employed innovative framework to generate and validate several thousands of possible amino acid encodings. Due to this approach, we are able to specify important physicochemical properties that define the best-performing alphabet. We confirmed the relevance of properties commonly associated with amyloidogenicity as hydrophobicity and discover new ones, as flexibility.

Our analysis was completed with the extraction of important n-grams, short motifs highly determining the amyloidogenicity or nonamyloidogenicity of the peptide. 65 important n-grams revealed that mostly alifatic and nonpolar amino acids as isoleucine, leucine and valine are responsible for the hydrophobic character of amyloids. The presence of aromatic amino acids is also a very significant signal of potential for amyloidogenicity.

The n-grams consisting of elements from the best-performing reduced amino acid alphabet separate efficiently amyloids and non-amyloids. Our classifier outperforms existing software for amyloid detection in all performance measures.

Acknowledgments and funding

Calculations were carried out in Wroclaw Center for Networking and Supercomputing (<http://www.wcss.pl>), grant No. 347 and by KNOW Consortium.

References

- Argos, P., Rao, J. K., and Hargrave, P. A. (1982). Structural prediction of membrane-bound proteins. *European journal of biochemistry / FEBS*, **128**(2-3), 565–575.
- Bhaskaran, R. and Ponnuswamy, P. (1988). Positional flexibilities of amino acid residues in globular proteins. *International Journal of Peptide and Protein Research*, **32**(4), 241–255.
- Black, S. D. and Mould, D. R. (1991). Development of hydrophobicity parameters to analyze proteins which bear post- or cotranslational modifications. *Analytical Biochemistry*, **193**(1), 72–82.
- Breiman, L. (2001). Random Forests. *Machine Learning*, **45**(1), 5–32.
- Charton, M. and Charton, B. I. (1982). The structural dependence of amino acid hydrophobicity parameters. *Journal of Theoretical Biology*, **99**(4), 629–644.
- 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.
- Jr, J. H. W. (1963). Hierarchical Grouping to Optimize an Objective Function. *Journal of the American Statistical Association*, **58**(301), 236–244.
- Kanehisa, M. I. and Tsong, T. Y. (1980). Local hydrophobicity stabilizes secondary structures in proteins. *Biopolymers*, **19**(9), 1617–1628.
- Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., and Kanehisa, M. (2008). AAIindex: amino acid index database, progress report 2008. *Nucleic Acids Research*, **36**(suppl 1), D202–D205.
- Kidera, A., Konishi, Y., Oka, M., Ooi, T., and Scheraga, H. A. (1985). Statistical analysis of the physical properties of the 20 naturally occurring amino acids. *Journal of Protein Chemistry*, **4**(1), 23–55.
- Kim, C. A. and Berg, J. M. (1993). Thermodynamic beta-sheet propensities measured using a zinc-finger host peptide. *Nature*, **362**(6417), 267–270.
- Nishikawa, K. and Ooi, T. (1986). Radial locations of amino acid residues in a globular protein: correlation with the sequence. *Journal of Biochemistry*, **100**(4), 1043–1047.
- Pontius, J., Richelle, J., and Wodak, S. J. (1996). Deviations from Standard Atomic Volumes as a Quality Measure for Protein Crystal Structures. *Journal of Molecular Biology*, **264**(1), 121–136.
- Prabhakaran, M. (1990). The distribution of physical, chemical and conformational properties in signal and nascent peptides. *The Biochemical Journal*, **269**(3), 691–696.
- Radzicka, A. and Wolfenden, R. (1988). Comparing the polarities of the amino acids: side-chain distribution coefficients between the vapor phase, cyclohexane, 1-octanol, and neutral aqueous solution. *Biochemistry*, **27**(5), 1664–1670.
- Sweet, R. M. and Eisenberg, D. (1983). Correlation of sequence hydrophobicities measures similarity in three-dimensional protein structure. *Journal of Molecular Biology*, **171**(4), 479–488.
- Takano, K. and Yutani, K. (2001). A new scale for side-chain contribution to protein stability based on the empirical stability analysis of mutant proteins. *Protein Engineering*, **14**(8), 525–528.
- 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.
- Wozniak, P. P. and Kotulska, M. (2014). Characteristics of protein residue-residue contacts and their application in contact prediction. *Journal of Molecular Modeling*, **20**(11).
- Wozniak, P. P. and Kotulska, M. (2015). AmyLoad: website dedicated to amyloidogenic protein fragments. *Bioinformatics (Oxford, England)*.
- Wright, M. N. and Ziegler, A. (2015). ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R. *arXiv:1508.04409 [stat]*. arXiv: 1508.04409.
- Zhou, H. and Zhou, Y. (2004). Quantifying the effect of burial of amino acid residues on protein stability. *Proteins*, **54**(2), 315–322.