

Interpretable deep learning amyloid nucleation by massive experimental quantification of random sequences

Mike Thompson¹, Mariano Martín², Trinidad Sanmartín Olmo², Chandana Rajesh³, Peter K. Koo³, Benedetta Bolognesi^{*2} and Ben Lehner^{*1,4,5,6}

¹Systems and Synthetic Biology, Centre for Genomic Regulation, The Barcelona Institute for Science and Technology (BIST), Barcelona, Spain

²Institute for Bioengineering of Catalonia (IBEC), The Barcelona Institute of Science and Technology, Barcelona, Spain

³Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

⁴University Pompeu Fabra (UPF), Barcelona, Spain

⁵Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain

⁶Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK

*Email: bbolognesi@ibecbarcelona.eu (B.B) and bl11@sanger.ac.uk (B.L.)

Abstract

Insoluble amyloid aggregates are the hallmarks of more than fifty human diseases, including the most common neurodegenerative disorders. The process by which soluble proteins nucleate to form amyloid fibrils is, however, quite poorly characterized. Relatively few sequences are known that form amyloids with high propensity and this data shortage likely limits our capacity to understand, predict, engineer, and prevent the formation of amyloid fibrils. Here we quantify the nucleation of amyloids at an unprecedented scale and use the data to train a deep learning model of amyloid nucleation. In total, we quantify the nucleation rates of >100,000 20-amino-acid-long peptides. This large and diverse dataset allows us to train CANYA, a convolution-attention hybrid neural network. CANYA is fast and outperforms existing methods with stable performance across diverse prediction tasks. Interpretability analyses reveal CANYA's decision-making process and learned grammar, providing mechanistic insights into amyloid nucleation. Our results illustrate the power of massive experimental analysis of random sequence-spaces and provide an interpretable and robust neural network model to predict amyloid nucleation.

Introduction

Specific insoluble protein aggregates in the form of amyloid fibrils characterize more than fifty clinical conditions affecting more than half a billion people (Fig. 1A)¹. These include common neurodegenerative disorders and the most frequent forms of dementia. Nonetheless, amyloids are present in all kingdoms of life and can have functional roles, including in humans². The importance of amyloids across biological functions and diseases has spurred massive research efforts, yet the determinants and mechanisms of their formation remain quite poorly understood^{3,4}.

Recent advances in cryogenic electron microscopy have allowed the atomic structures of many mature amyloid fibrils to be determined⁵. Amyloids share a cross- β structure wherein hydrogen-bonded β -strands are perpendicularly stacked along the fibril axis, creating β -sheets that face each other and are parallel to the fibril axis^{3,6,7}. The same amino acid (AA) sequence can form multiple different filament structures, and in at least some cases, these amyloid polymorphs appear to be associated with distinct clinical conditions. Amongst humans, amyloid fibrils typically have hydrophobic cores, for which hydrophobicity and β -strand propensity form the basis of many computational methods to predict amyloid propensity from sequence⁸⁻¹⁴. However, other amyloids, for example yeast prions, have very different sequence composition, hinting at a richer diversity of amyloid-forming sequences^{15,16}.

In contrast to the remarkable advances in the structural characterization of mature fibrils, the process of amyloid formation—how soluble proteins overcome a free energy barrier to nucleate fibrils (Fig.1B)—is much less understood. Time-resolved structure determination has been used to study the *in vitro* assembly of amyloids, revealing a striking diversity of intermediate filament folds appearing and disappearing as fibrillation proceeds^{17,18}. However, how this process initiates and why it only occurs for some sequences under physiological conditions remains unclear. Mature amyloid fibrils are very stable and are likely to be the thermodynamically favored state at high protein concentration for many proteins^{19,20}. There is, however, a very high energy barrier to amyloid nucleation for most proteins i.e. the process is under kinetic control²⁰. The kinetic control of amyloid nucleation is, therefore, the key problem to understand: why do only some sequences nucleate amyloid formation on timescales relevant to biology?

We believe that our ability to understand and predict amyloid formation is currently data-limited. Only a small number of sequences with high amyloid propensity are known, restricting the development and benchmarking of predictive computational methods²¹. In contrast to the very small number of characterized amyloids, the number of known and possible protein sequences is vast. For example, there are 20^{20} ($>10^{26}$) different sequences of 20 amino acids. Such a large sequence space can never be substantially explored by experimental or computational techniques, necessitating the development of predictive models.

To address this data gap we have developed a massively parallel selection assay that allows the nucleation kinetics of thousands of different sequences to be quantified in a single experiment^{22,23}. We have previously used this assay to quantify the change in nucleation

kinetics for all possible substitutions, insertions and deletions in the amyloid beta peptide that aggregates as a hallmark of Alzheimer's disease and is mutated in familial forms of the disease. The resulting data agree very well with *in vitro* measured nucleation kinetics and also identify the known familial Alzheimer's disease-causing variants^{22,23}. However, these datasets are limited to small changes to a single sequence, hindering utility for general-purpose model-building.

Here we apply the same assay at a much larger scale, using it to quantify the nucleation of >100,000 peptides with completely random sequences. We use this massive dataset to train CANYA, a convolution-attention hybrid neural network that predicts the propensity of any primary sequence to form amyloids. This fast model outperforms existing predictors of protein aggregation on both internal and out-of-sample datasets, demonstrating the power of massive sequence-space exploration. We then use post-hoc explainable AI (xAI) analyses to provide mechanistic insights into CANYA's decision-making process and learned grammar. CANYA provides a robust and interpretable neural network model for understanding and engineering amyloid-forming proteins. More generally, our results provide a very large and well-calibrated dataset to train and evaluate models beyond CANYA and they demonstrate the utility of massive experimental analysis of random protein sequence-spaces.

Results

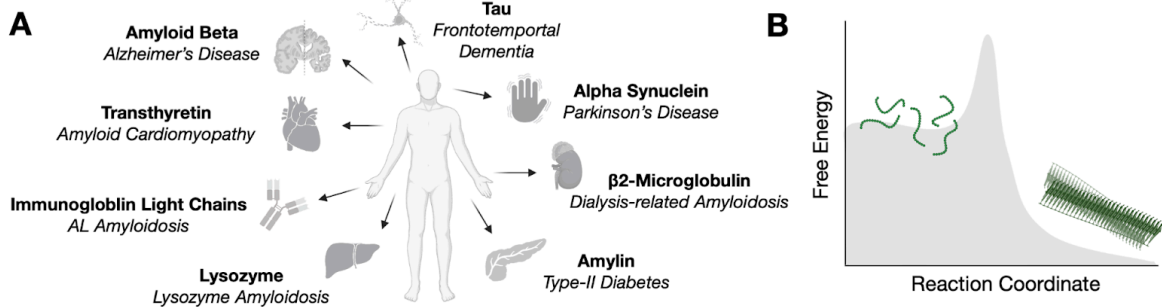
Massively parallel quantification of amyloid nucleation kinetics

To better understand the sequence determinants of amyloid nucleation kinetics, we used an in-cell selection assay to quantify the rate of nucleation of more than a hundred thousand peptides with fully random sequences. We generated four libraries (NNK1-4) of random 20 AA peptides using NNK degenerate codons (where N = A/C/G/T and K = G/T) and expressing them as fusions to the nucleation domain of Sup35, a yeast prion-forming protein that allows fitness-based selection for amyloid nucleation (Fig. 1C)²²⁻²⁴. Briefly, fusion sequences that nucleate amyloids sequester Sup35 resulting in translational readthrough of a premature stop codon in the *ade1* gene so that cells containing those sequences become able to survive in medium lacking adenine. Enrichment or depletion of each sequence after selection can be quantified by deep sequencing, with enrichment scores linearly related to the log of *in vitro* amyloid nucleation rates^{22,23}.

Each library was selected independently and sequencing was used to quantify the relative enrichment ('nucleation score') for each genotype in the library. Sequences in the first three experiments made up our training and testing sets (NNK1-3, ~111,000; Fig. 1D; Extended Data Files 1 and 2), corresponding to about a $1/10^{17}$ fraction of the possible sequence space (20^{20}), while sequences from the fourth experiment (NNK4, ~7,000) were used as a held-out validation data set. After data processing and quality control, the vast majority of sequences had a nucleation score of 0. Consequently, we classified sequences with a nucleation score significantly greater than 0 (one-sided Z-test, FDR adjusted p-value ≤ 0.05) as nucleators (n=21,936), and all other sequences as non-nucleators (n=88,470) (Fig. 1E). Importantly, these nucleation scores are reproducible, as measured by an additional selection experiment on a designed library (replication library) re-quantifying the nucleation of 400

sequences sampled across all four libraries (Pearson correlation range 0.506-0.797, Fig. 1F, Supp. Fig. 1).

Amyloid nucleation occurs in > 50 human diseases



Inferring amyloid nucleation from sequencing-based frequency changes

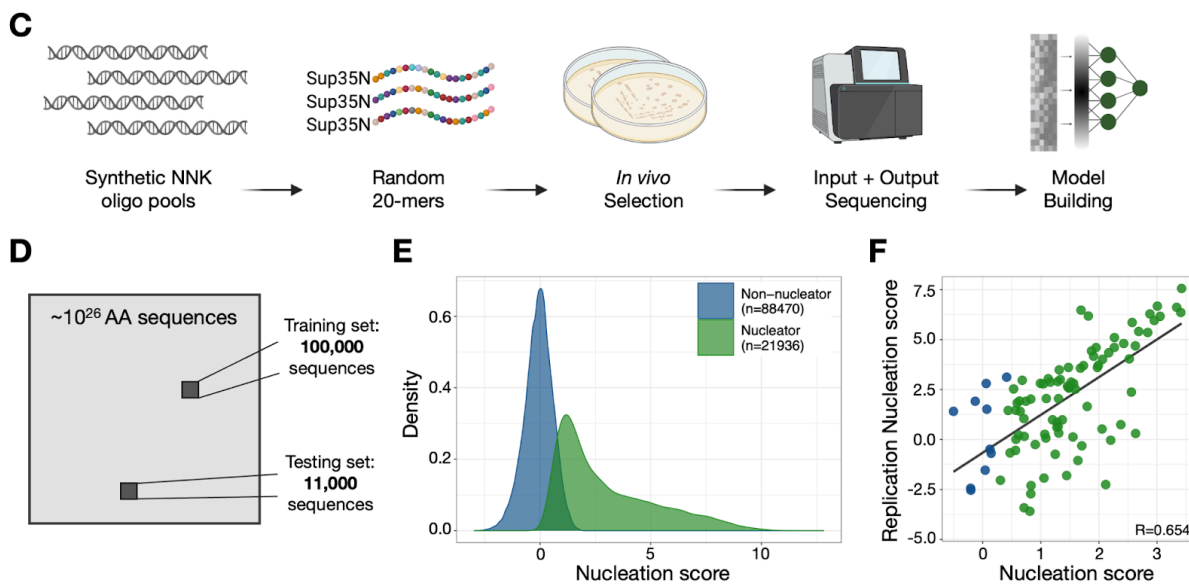


Figure 1. Quantifying the nucleation of >100,000 random peptides. (A) Examples of amyloids in human diseases. (B) The amyloid state is thermodynamically favorable, but requires overcoming a kinetic barrier. (C) Experimental design. (D) While we explore over 110,000 sequences, our dataset is a tiny sample of the possible sequence space. (E) The assayed nucleation scores of sequences labeled “Nucleators” and “Non-nucleators” in our experiment. (F) An example of a follow-up replication experiment using a synthesized library (NNK3; see Supp. Fig. 1 for others; Extended Data File 3).

Nucleating sequences span a large sequence-space and are non-trivial to predict

After classifying sequences as nucleators and non-nucleators, we sought to characterize each class through amino acid composition (Fig. 2A), physicochemical properties (Fig. 2B), and current amyloid prediction tools (Fig. 2C).

First, we examined the differences in amino acid frequency between nucleating and non-nucleating sequences. Differences in frequencies were generally modest, however we observed statistically significant differences owing to the large sample size of our data. When looking at composition independent of position, nucleators had higher frequencies of

cysteine (difference in frequency 0.012, $p < 2e-16$), asparagine (0.009, $p < 2e-16$), and isoleucine (0.005, $p < 2e-16$), and lower frequencies of arginine (-0.010, $p < 2e-16$), leucine (-0.008, $p < 2e-16$), and lysine (-0.006, $p < 2e-16$; Fig. 2A, See Supp. Table 1 for full differences). Moreover, both nucleators and non-nucleators covered the beta-sheet propensity and hydrophobicity spaces of the human proteome and known amyloid sequences, and nucleators had slightly higher values of both than non-nucleators on average (difference in means of hydrophobicity=0.130, beta-sheet propensity=0.012, both two-way t-test p-values $< 2e-16$; Fig. 2B). Considering position-specific composition, differences were again modest, ranging from a difference in frequency from -0.06 to 0.03 (Fig. 2D). Subsequently, we grouped amino acids by their physicochemical properties to check for more broad, position-specific differences between the two sequence classes (Fig. 2E). Toward the N-terminus of the random sequence (i.e., closer to Sup35), nucleators were significantly enriched (chi-squared test) for aliphatic residues (min p-value=1.54e-13, position 2 difference=0.033), and significantly depleted for positive (min p-value=1.57e-25, position 9 difference=-0.032) and negative residues (min. p-value=3.14e-11, position 2 difference=-0.016). The differences in charge waned toward the C-terminus (min. p-value above position 15=1.03e-3, position 20 charged difference=0.011), however, and frequency differences in aliphatic residues changed such that nucleators were significantly *depleted* for aliphatic residues relative to non-nucleators (min. p-value=5.77e-39, position 19 difference=-0.058). Several groupings showed other position-sensitive differences, such as an enrichment of aromatic residues toward the C-terminus in nucleators (min. p-value=5.09e-6, position 19 difference=0.015), an enrichment of varying strength for polar residues in nucleators (p-value= 5.57e-8 position 1 difference=0.023, p-value=9.41e-7 position 17 difference=0.020), and the enrichment of cysteines away from the ends of the random construct (min. p-value=1.11e-28, position 10 difference=0.023).

Despite statistical significance, we highlight that differences in sequence space are subtle. In other words, the collection of slight variation in amino acid frequencies offers minimal insight or definitive conclusions around the overall properties or characteristics determining nucleation in our experiment. To attempt to elucidate characteristics that separate the sequence classes and consequently learn important axes of variability, we turned to dimensionality reduction techniques. In addition to manually examining differences within the first several dimensions, we also used the scores in lower-dimensional space as features in a logistic multiple regression task to distinguish nucleators from non-nucleators. Using principal components analysis (PCA), we observed no clear separation between nucleators and non-nucleators whether we used amino acid composition alone (cumulative variance explained from the top 10 PCs = 54.7%, Area Under ROC curve (AUC) using all 10 PC scores=0.601, 95% Confidence Interval (CI)=[0.596, 0.607], Supp. Fig. 2), or maintained positionality of the amino acids when fitting the model (cumulative variance explained from the top 10 PCs = 3.1%, AUC=0.564, 95% CI=[0.559, 0.570], Supp. Fig. 2). This modest separation between classes of sequences was consistent even when using non-linear embedding techniques (first 10 UMAPs AUC=0.584, 95% CI [0.578, 0.589]), or adding amino acid propensities to the dimensionality reduction tools (first 10 PCs AUC=0.614, 95% CI [0.608, 0.619]; Supp. Fig. 2).

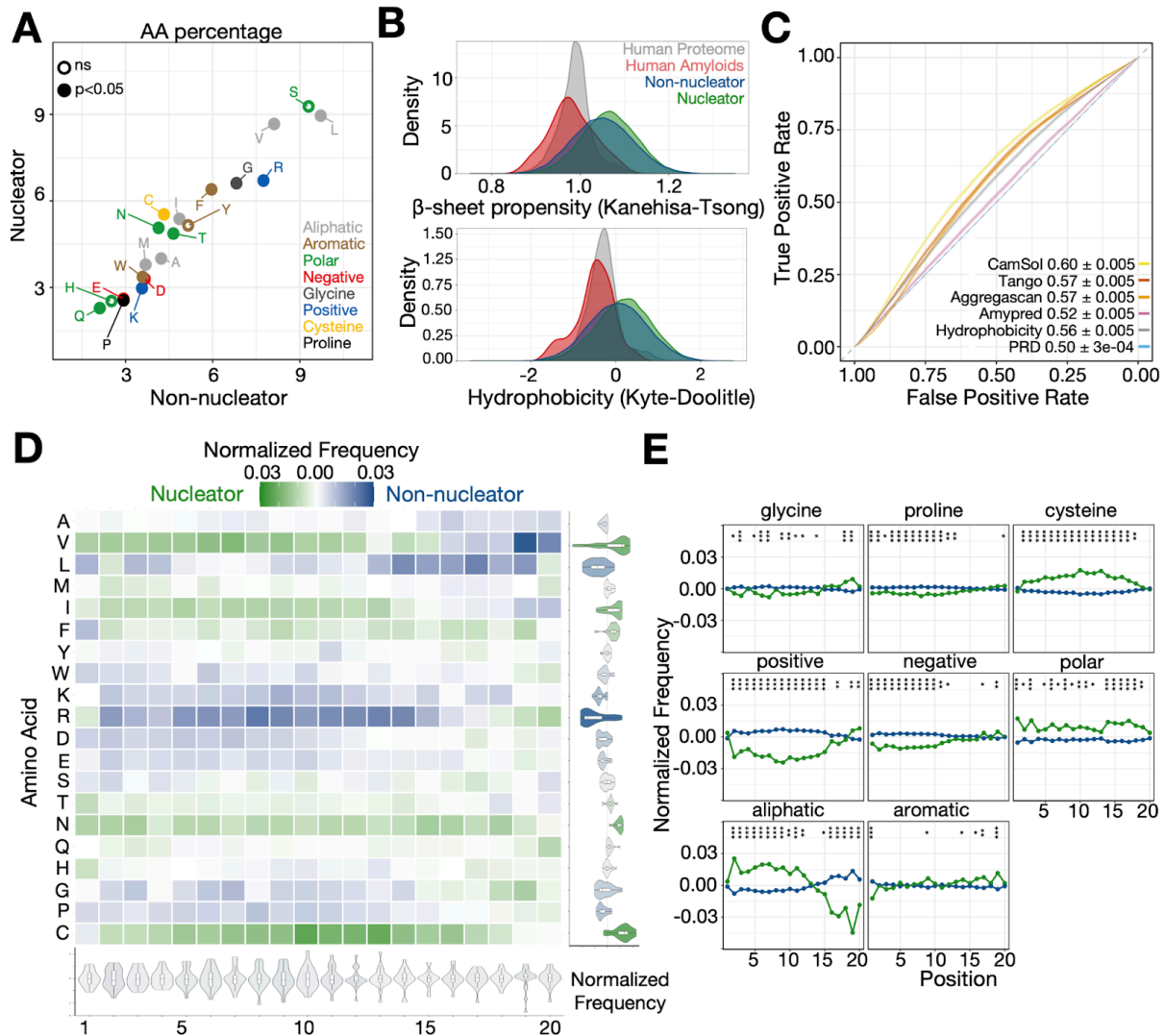
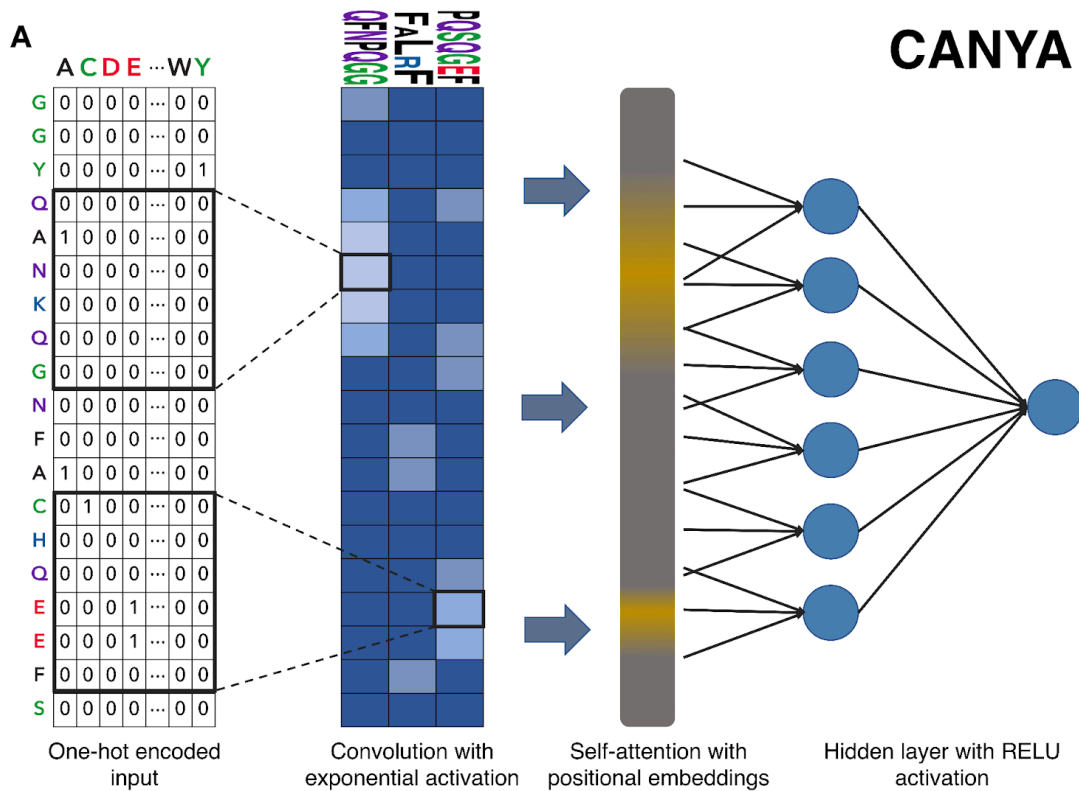


Figure 2. Differences in amino acid composition across nucleating and non-nucleating sequences are subtle. (A) The percent composition of residues grouped by their physicochemical properties in nucleators and non-nucleators. (B) The hydrophobicity and beta-sheet propensity of assayed sequences relative to known human amyloids (Supp. Table 2) and the human proteome. (C) The predictive power of previous amyloid predictors on the random sequences. (D, E) The position-specific differences in amino acid frequencies across nucleating and non-nucleating sequences. Asterisks indicate marginal p-value (chi-square test) lower than 0.05 (**); lower than 0.01 (***) , lower than 0.001 (****).

As dimensionality reduction methods were unable to distinguish the classes of sequences, we next explored whether separation is possible using existing amyloid predictors. Beyond hydrophobicity indices, several of these methods include structural information²⁵ or model biophysical mechanisms¹¹, potentially enabling them to capture more complex features of nucleation. We applied several state-of-the-art amyloid prediction algorithms to our data and found that the methods either failed to generalize to our data or had only modest predictive power (Fig. 2C, CamSol, highest AUC=0.598, 95% CI [0.593, 0.603]). We posit that, since many of these tools have been trained on very small sets of known amyloids or moderate numbers of short hexamer sequences, their applicability to our experimental data may be limited. To understand where the methods underperformed, we examined the scores from

the highest performing methods (CamSol¹² and TANGO¹¹) and found that non-nucleating sequences with a high-predicted nucleation score had higher hydrophobicity (two-sided t-test p-value<2e-16) than all other non-nucleating sequences (Supp. Table 3). We also found that low-predicted nucleators had higher presence of positive (two-sided t-test p-value<2e-16) and negative (p-value<2e-16) residues than all other nucleators (Supp. Table 4).



Performance of CANYA on training datasets

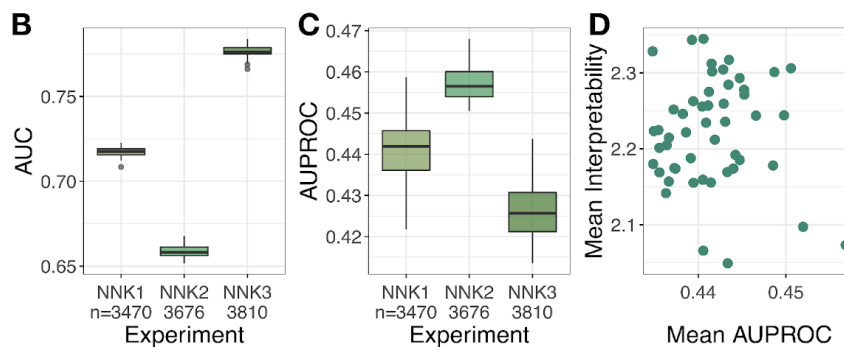


Figure 3. Convolution-Attention Network of amyloid Aggregation (CANYA). (A) CANYA is a 3-layer neural network with 65,491 parameters. The model contains 100 filters, a single attention head with key-length 6, a dense layer with 64 nodes, and finally a sigmoid output layer. (B-D) Evaluation metrics across the top 50 performing (of 100) model fits of CANYA. (B) The area under receiver operating characteristic curve (AUC) for held-out testing sequences. (C) The area under precision recall curve (AUPROC) for held-out testing sequences. (D) The interpretability score (KL divergence; Methods) calculated on all held-out test sequences plotted against the mean AUPROC across experiments. See Supp. Fig. 3 for results on all 100 model fits.

A hybrid neural-network predicts nucleation

Given previous approaches failed to accurately predict nucleation status within our dataset, we built our own model to capture the sequence-nucleation score landscape. Concretely, we developed a hybrid neural network which we term CANYA, or Convolution Attention Network for amyloid Aggregation. Though a neural network may seem inherently less interpretable than simpler models, as we explain below, the architecture of CANYA is not only simple, but also biologically motivated. CANYA builds off the observation that known amyloids are composed of interacting short sequences, such as stacked beta sheets, and treats this information as an inductive bias for the model—first the sequences are passed through a convolutional layer which discovers ‘motifs’, then these motifs are passed through an attention layer to learn positional effects of motifs and to encourage these motifs to interact with each other (Fig.3 A). Moreover, we set the filter lengths of the convolutional layer based on the distribution of secondary structure lengths in 80 known amyloid fibril structures (WALTZ-DB²⁶, Supp. Fig. 4). Though—to our knowledge—this class of models is new to proteins, convolution-attention hybrid models have been used in genomics and found to serve as a sound inductive bias for discovering motifs and their interactions^{27,28}.

We trained CANYA 100 times on over 100,000 synthetic sequences and their respective nucleation status to learn the sequence-nucleation landscape. Unlike massive, computationally intensive neural networks, CANYA comprises only three layers (spanning 65,491 parameters) and requires less than an hour to train on a basic, modern CPU. Despite this simplicity, and having only observed a small fraction of the possible sequence space, CANYA substantially improved the prediction of nucleation status of held-out test sequences (average AUC=0.710, 0.650, 0.769 across NNK experiments 1-3 respectively, Fig. 3B-C) over previous methods (max AUC CamSol, NNK1=0.617, NNK2=0.537, NNK3=0.673).

To understand the differences in performance across methods, we examined the sequence scores between the next best performing method (CamSol) and CANYA. We found that the largest discrepancies for non-nucleating sequences occurred in hydrophobic sequences with tryptophans, and in cysteine- or asparagine-rich sequences with few aliphatic residues in the case of nucleating sequences (Supp. Tables 3 and 4). Our results not only highlight the utility of exploring a vast sequence space, but also suggest that CANYA is able to contextualize physicochemical properties within sequences (e.g., among hydrophobic sequences, CANYA adjusts its score in the presence of bulky, or disruptive residues).

Crucially, we developed CANYA with the goal of interpreting the grammar of nucleation rather than maximizing predictive power. We accordingly scored each trained instance of CANYA using a recently developed interpretability metric to select a model amenable to uncovering this learned grammar²⁹. Briefly, this metric examines the enrichment of motifs utilized when training the model and compares them to the set of all equal-length ($k=3$) kmers in the training sequences (Methods). Strong enrichment (i.e., divergence from the background training sequences) indicates a model may yield clearer resolution in downstream interpretability analyses. Though the area under the precision-recall curve (AUPROC) of test sequences was more consistent than AUC across experiments (average AUPROC NNK1=0.434, NNK2=0.452, NNK3=0.415 ; Fig. 3C), we did not find a correlation between predictive performance and this interpretability metric (correlation of average AUPROC and interpretability score $r=-0.059$, $p\text{-value}=0.6847$, Fig. 3D). We therefore chose

the trained model with the highest interpretability score, conditional on the fact that it scored better than the median-performant model (of 100 training runs; Methods).

CANYA robustly predicts nucleation across external datasets

After establishing that CANYA can predict the experimental nucleation status from primary sequence, we sought to understand whether the nucleation function learned by CANYA is applicable to sequences and contexts outside the experimental dataset from which it was learned. To examine this capability, we evaluated the performance of CANYA on an additional set of random, synthetic sequences, as well as across an independent collection of several public datasets: WaltzDB²⁶, CPAD³⁰, and AmyPred³¹. We also compared CANYA to previous amyloid-prediction approaches on these datasets. Testing each method across a wide range of datasets more concretely enables us to evaluate whether there exist specific conditions or sequences for which one specific predictor is more suitable than another.

First, we evaluated each method on an additional set of ~7,000 unseen random sequences quantified in our nucleation assay (Fig. 4A). Here, we expect CANYA to outperform previous methods as it was trained using data from the same selection assay. However, this serves as a measure of whether CANYA has genuinely learned functional information from its training, as the sequence spaces spanned by the training and test sequences are effectively independent (i.e., $\sim 10^5$ and $\sim 10^3$ samples from a $>10^{22}$ sequence landscape). CANYA remained highly accurate on the 7,000 unseen sequences, significantly outperforming all tested previous methods^{11–14,31,32} (AUC CANYA=0.809, 95% CI [0.798, 0.821] and the next-best performing method AUC=0.707 95% CI [0.694, 0.719]; Fig. 4B, PROC in Supp. Fig. 5). Of the previous methods tested, AggreScan, TANGO, and CamSol significantly outperformed hydrophobicity scales (min AUC=0.679 95% CI [0.665, 0.693], hydrophobicity AUC=0.593 95% CI [0.579, 0.607]).

We next evaluated the methods on 1,400 hexapeptides from WALTZ-DB, one of the largest databases of amyloidogenic and non-amyloidogenic sequences²⁶. However no method significantly outperformed hydrophobicity for classifying aggregating hexamers in WALTZ-DB (AUC=0.813 95% CI [0.791, 0.836]) (Fig. 4C). The hydrophobicity distributions of amyloid and non-amyloid hexamers in WALTZ-DB are indeed very distinct (Supp. Table 5), suggesting biases in this dataset or that hydrophobicity dominates the aggregation potential of such very short peptides. This cautions against the use of such short sequences for model training.

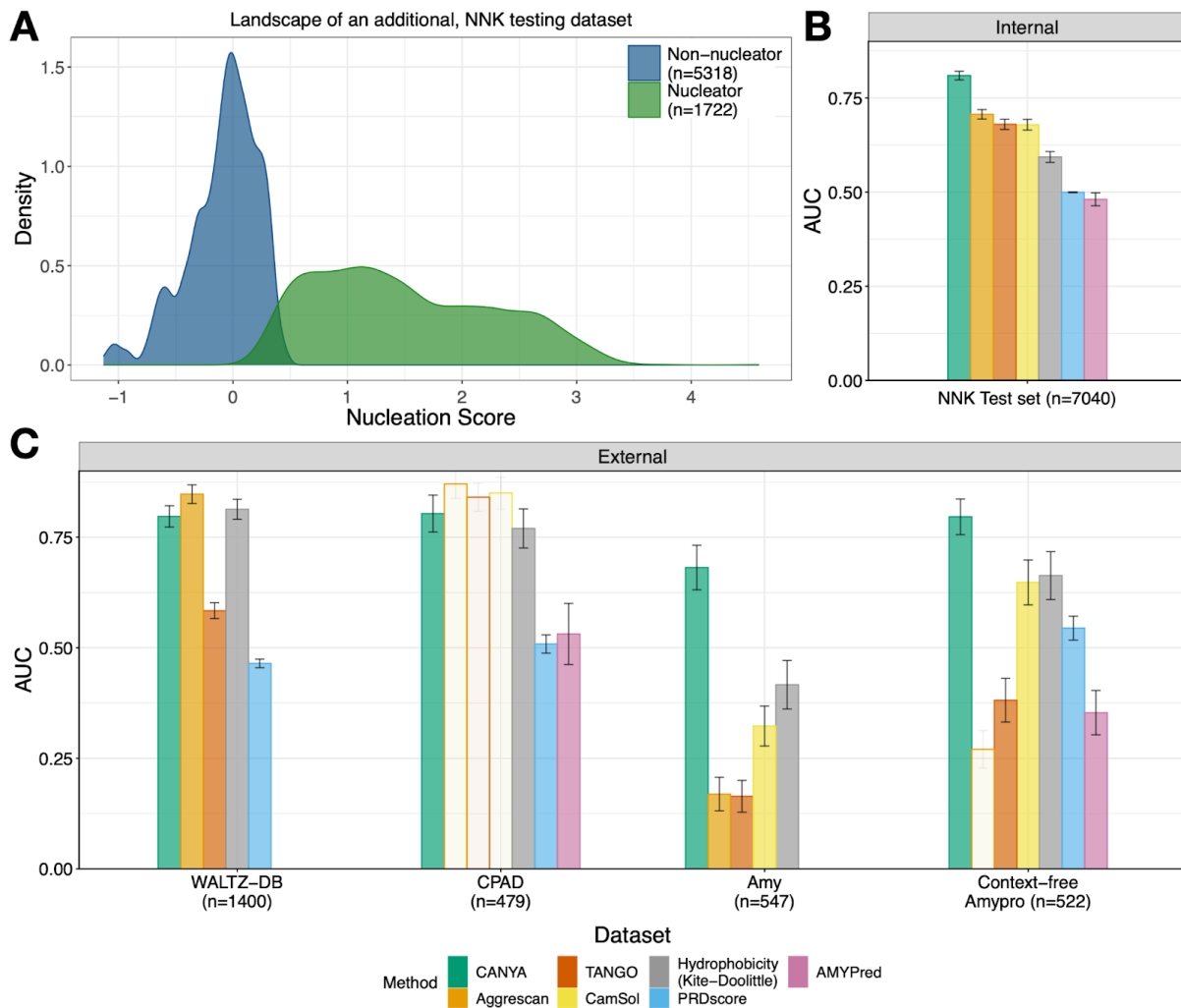


Figure 4 Stable performance of CANYA across diverse prediction tasks. (A) The fitness of a held out nucleation assay experiment with 7040 random length-20 AA sequences (NNK4 dataset). (B) The predictive performance of CANYA relative to previous methods on the held out, test set of 7040 NNK sequences. (C) The AUC of each method on several external datasets. Low-opacity bars represent cases in which the method used data from the testing dataset for training and thus are not valid out-of-sample evaluations. See text for additional descriptions of datasets (Methods, Supp. Table 5) as well as performance reported as area-under precision-recall curve (AUPROC; Supp. Fig. 5).

We next turned to a dataset comprising known amyloid-forming sequences from the Curated Protein Aggregation Database (CPAD). CPAD contains a set of over 2,000 amyloid-forming and non-amyloid-forming sequences curated from literature sources (such as GAP, WALTZ-DB, and AmyLoad)^{26,33–35}. To eliminate overlap with WALTZ-DB and to evaluate the methods on longer sequences, we excluded sequences less than 10 residues long. Here, our evaluation consisted of 479 sequences with median length 16 (Q1 length=10 and Q3 length=22) comprising 304 amyloid-forming sequences and 175 non-amyloid-forming sequences (Supp Table 5). We also note that several of the previous methods (including TANGO, CamSol, and Aggrescan) were directly fitted on sequences within CPAD, violating the ability to evaluate their out-of-sample prediction on this dataset. This complication is exacerbated by several methods (e.g., TANGO, CamSol) also being ensemble methods (or extensions) that leverage several algorithms for prediction—it is not trivial to account for, or remove, these previously seen sequences, as any sequence that was used for training the

main algorithm or their antecedent ensemble methods is not out-of-sample. CANYA performed similarly well to previous methods on CPAD (AUC=0.804, 95% CI [0.762, 0.845], AUPROC=0.855, 95% CI [0.817, 0.890]), including those fine-tuned on CPAD sequences (Fig. 4C).

We also evaluated performance on the Amy dataset³⁶, which is composed of 547 non-redundant, long (>50 residue) sequences of amyloid- and non-amyloid-forming proteins spanning both prokaryotes and eukaryotes. Specifically, Amy contains a set of 382 non-amyloid sequences (median length=708 [Q1=344.5, Q3=1375.5]) gathered from UniProt and 165 amyloid sequences (median length=162 [Q1=77, Q3=443]) from the AmyPro database³⁷. The AmyPro database contains literature-mined amyloid precursor proteins with validated amyloidogenic sequence regions—portions of an amyloid-forming protein that when isolated from the remaining sequence have been confirmed to form amyloids in external experiments. The median length of sequences in Amy (539) is substantially longer than those of the previous datasets (next longest median length=19 in NNK4, or 16 in CPAD), so this prediction task evaluates whether methods can account for context-specific and distal effects when generating their propensity scores. Strikingly, CANYA was the only predictor with both statistically significant AUC and PROC in this task (AUC=0.681 95% CI [0.631, 0.731], Fig. 4C, AUPROC 0.495, 95% CI [0.428, 0.568]). The poor performance of hydrophobicity and previous methods suggests the importance of features beyond sequence composition in determining amyloid propensity in longer protein sequences.

Finally, we evaluated whether each method could identify amyloidogenic regions of each protein in the AmyPro dataset. To do so, we removed from each nucleating sequence its amyloidogenic region and labeled this amyloidogenic sequence as a nucleator, then labeled any remaining, non-overlapping, non-core parts of the nucleating sequence as a non-nucleator (breaking these remaining non-core sequence into non-overlapping sequences of maximum length 100; see Methods). This task evaluates whether methods can distinguish an amyloidogenic region from a non-amyloidogenic region in the absence of any contextual region, for which we term it “Context-Free AmyPro.” CANYA significantly out-performed all previous approaches (AUC=0.796, 95% CI [0.756, 0.837], Fig. 4C), and was the only method to significantly out-perform hydrophobicity (AUC=0.663, 95% CI [0.609, 0.717]). As with the Amy dataset, the amyloid-labeled sequences were much shorter (median length=17 [Q1=10, Q3=35]) than non-amyloids (median length=100 [Q1=81, Q3=100]), however amyloid sequences were generally more hydrophobic (median=-0.09 [Q1=-0.63, Q3=0.72]) than non-amyloid sequences (median=-0.47, [Q1=-0.79, Q3=-0.26]; Supp. Table 5). The improvement of CANYA over hydrophobicity and CamSol suggests CANYA has learned more complicated features of nucleation than hydrophobicity alone, and that these features are informative independent of protein context.

In summary, CANYA’s performance is state of the art and consistent across diverse prediction tasks and protein sizes.

CANYA learns physicochemical nucleation motifs

After establishing the predictive power on CANYA, we performed a series of interpretability analyses to understand how CANYA assigns its nucleation score and to elucidate

difficult-to-see patterns that differentiate the nucleators and non-nucleators in the training data.

First, we establish a set of physicochemical “motifs” learned by the model. To visualize motifs learned by the model, we constructed position-weight-matrices (PWMs) using kmers that activated a given filter at least 75% of the maximum-activating kmer (Methods). We selected a filter length of 3, as this was the mode length of secondary structures in structurally resolved amyloids (Methods; Supp. Fig. 4). Motifs generally showed low information content per position but showed clear physicochemical preferences (Fig. 5). For example, many motifs capture blocks of hydrophobicity (clusters 1 and 2) or charge (clusters 6 and 8). Some motifs showed heterogeneity, or position-preferential effects, such as polar or charged residues being surrounded by hydrophobic (clusters 4 and 5) or aromatic residues (clusters 7, 9 and 10; Fig. 5).

We next turned to a post-hoc interpretability method named Global Importance Analysis (GIA) to learn the effect of each motif³⁸. Briefly, GIA learns effect sizes by embedding a motif of interest in a set of background sequences, then comparing the difference in the models’ predicted nucleation propensity between these background sequences with and without the embedded motif (Fig. 6A). The effects learned by CANYA recapitulated previously known amyloid biology—hydrophobic motifs strongly increased a given sequence’s propensity to nucleate, and charged, proline-containing motifs lowered sequences’ propensity to nucleate (Fig. 5)^{39–41}. Motifs containing residues enriched in yeast prions (Q/N) also increased amyloid propensity (weaker motifs of clusters 1, 2, and 3, stronger motifs of cluster 4), as did motifs enriched in cysteine (cluster 3) or aromatic residues (cluster 2; Fig. 5). Interestingly, tryptophan-containing motifs showed effect sizes in both directions, and these differences corresponded to cases in which the tryptophan was surrounded by charged residues (negative effect; clusters 7, 9, 10; Fig. 5), or hydrophobic, polar, or aromatic non-tryptophan residues (positive effect; clusters 1, 2, 3; Fig. 5). Notably, CANYA also found a set of motifs enriched in hydrophobicity with a positively charged residue (cluster 5, Fig. 5), suggesting it can capture previously uncharted areas of the amyloid sequence space.

We next sought to cluster the motifs to generate a more concise representation of what the model has learned and to reduce the dimensionality of *in silico* analyses to extract further information learned by the model. To do so, we first generated BLOSUM scores (which capture a similarity of amino acids based on evolutionary divergence) for each motif, then performed affinity clustering on the BLOSUM scores to derive a candidate set of clusters⁴². Next, we used the previously calculated GIA scores to investigate whether the effects of motifs corresponded to the same direction of effect of their respective physicochemical cluster. There were seven discrepancies—namely, several motifs containing histidines and tryptophans (Supp. Fig. 6). As our goal was simply to interpret the model and to reduce the number of *in silico* experiments we needed to run, we excluded these seven motifs from any downstream analyses. We verified that this approach results in a sound set of clusters by re-running GIA using the clusters as the feature of interest and confirming that the learned effect size for a cluster was consistent with the motifs of which it is composed (Methods; Supp. Table 6). Summarily, we were left with 10 clusters on which to perform downstream *in silico* experiments, effectively reducing the number of experiments by at minimum one order of magnitude (from 100 filters).

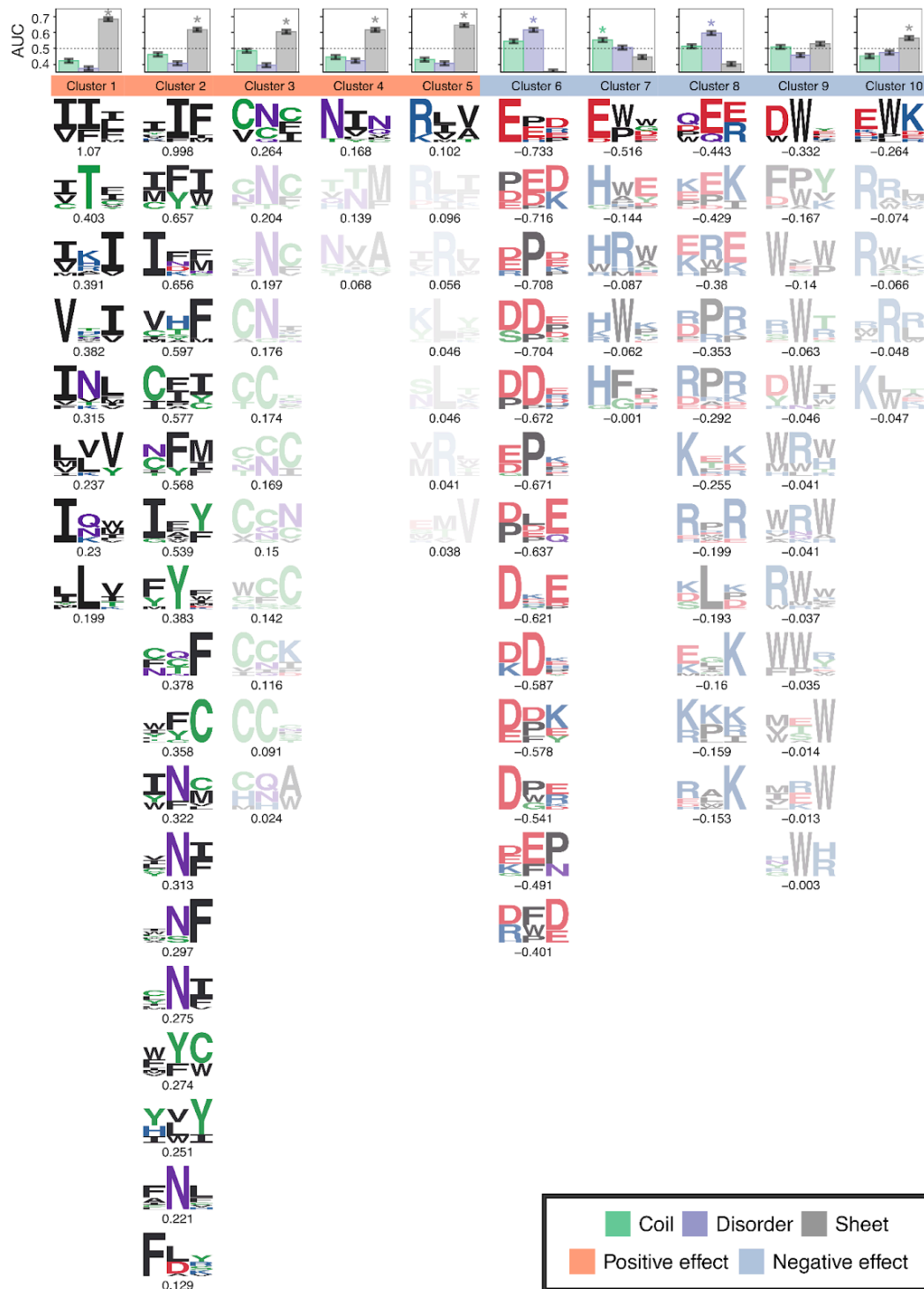


Figure 5 CANYA discovers physicochemical nucleation motifs. The motifs discovered by CANYA, clustered by their physicochemical properties and GIA effect sizes, then sorted based on their effect size magnitude. Translucency represents the ratio of cluster effect size compared to the strongest cluster (Methods). The enrichment (in AUC) of motif-cluster presence in secondary structures of resolved amyloids in Uniprot (Methods; Supplementary Fig. 7). The dashed lines represent an AUC of 0.50 and asterisks represent structures for which the enrichment was significantly higher than both 0.50 and the second most-enriched structure.

Motif activation in known amyloid structures

We examined whether the motif clusters discovered by CANYA showed propensity for secondary structures in known amyloid fibril structures (from the Structural analysis of Amyloid Polymorphs (StAmP) database⁴³). We included in our comparison full-length resolved structures of amyloid fibrils for 114 PDB entries comprising amyloid structures of 23 proteins (Supp. Table 7). Here, we used the activation energy of a cluster across positions to predict whether or not the corresponding position was in a beta-strand, other structured region (coil), or unresolved (disordered, see Methods). The AUC from this task serves as a metric of whether high activation (high matching score) of a motif is associated with a specific structural element. As expected, the clusters with high hydrophobicity and positive effect size were most strongly associated with activating in beta strands (Fig. 5; Supp. Fig. 7). Concretely, positive-importance clusters generally showed tendency toward presence in beta-sheets (max AUC=0.683, 95% CI [0.672, 0.693], cluster 1), whereas the strongest enrichment amongst negative-importance clusters was observed in disordered regions (max AUC=0.617, 95% CI [0.605, 0.628], cluster 6). Interestingly, negative-importance clusters with tryptophan showed varying enrichments in secondary structures (Fig. 5; Supp. Fig 7). Clusters with tryptophans near histidines were moderately enriched in coils (cluster 7; AUC=0.553, 95% CI [0.541, 0.565]), tryptophans adjacent to lysines or arginines showed moderate enrichment for strands (cluster 10; AUC=0.566, 95% CI [0.554, 0.577]), and tryptophans near other tryptophans or aromatic residues (cluster 9) showed no significant enrichment for any structure.

Motif position-dependence and interactions

Treating the motif clusters as input for GIA, we performed an additional set of experiments to evaluate whether CANYA has learned positional information of motif effects and whether motif effects are additive (Fig. 6A).

To learn positional information for each cluster of motifs, we ran an experiment in which we calculated the GIA effect of the cluster at every position of the construct, and compared it to the global, position-averaged effect of the cluster. These comparisons revealed that CANYA indeed learned position-relevant information across each cluster of motifs (Fig. 6D). Generally, the positive-effect clusters showed diminished effects at the ends of the construct and stronger effects at the center (Fig. 6D). The range of percent change was most drastic for cluster 5, potentially due to the presence of a charged residue (%-change in effect from 14.81% 95% CI [9.59, 19.59] to -39.81 %, 95% CI [-43.76, -36.02]; Fig. 6D). Clusters 1, 2, and 4 followed similar trends, however the changes were much more modest (highest percent change position 12, cluster 4=4.42%, 95% CI [2.05, 6.73], lowest percent change position 18, cluster 1=-8.58%, 95% CI [-9.53, -7.63]). Cluster 3, which is marked by its high presence of cysteines, followed a similar trend except that at the C terminus its effect significantly increased (% change=10.76, 95% CI [9.11, 12.37]).

Conversely, the negative-importance clusters all had strengthened effects toward the Sup35 peptide, and all but the negative, proline-rich cluster (cluster 6) had diminished effects toward the C terminus (Fig. 6D). This may be due to the fact that cluster 6 was the most negatively charged cluster, consistent with negative charges in the C-terminus of some amyloid-forming peptides reducing fibril formation⁴⁴.

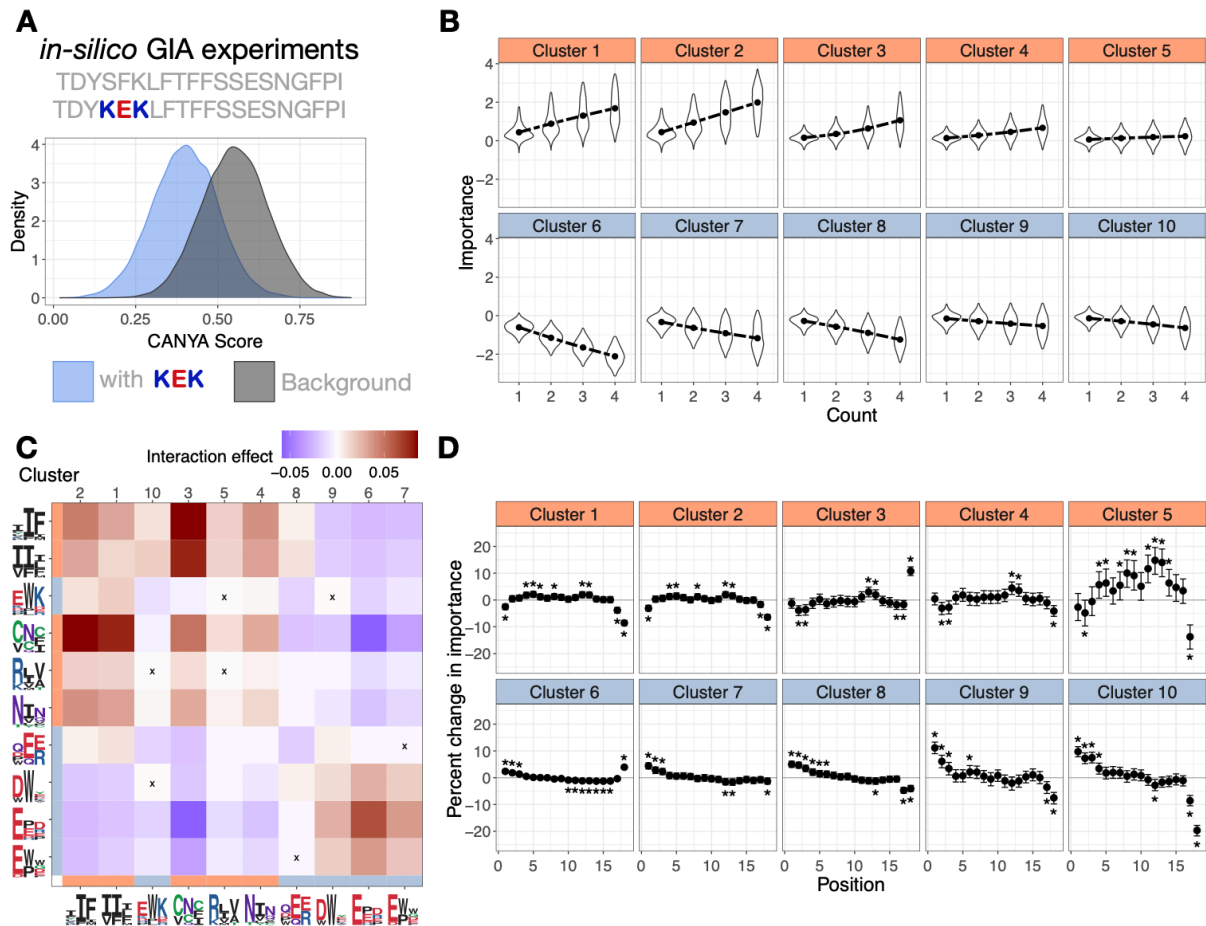


Figure 6. *In-silico* experiments reveal CANYA's learned nucleation grammar. (A) An example of an experiment using GIA, an explainability tool to extract importance (effect sizes) of features in a model. Briefly, model predictions for a background set of sequences are compared to predictions on the same set of sequences with a feature (motif) embedded in them. (B) The distribution of effects from adding 1-4 copies of a cluster-motif to sequences. Points represent importance. (C) Interaction importance from adding motifs from two clusters to sequences. Warmer colors indicate higher CANYA score than from marginally adding the motifs (and their effects) separately to sequences, whereas cooler colors represent a CANYA score lower than expected from adding marginal motif effects. "X" indicates effects that were not significantly different from 0. (D) The position-dependence of motif effects. Plotted is the percent change of a position-specific effect relative to the motif's global, position-averaged effect. Stars represent a significantly non-zero percent change in effect.

Tryptophan-containing motifs (clusters 7, 9, and 10) typically had increased effects closer to Sup35 (greatest % change cluster 9, position 1=11.14% 95% CI [8.96, 13.28]) and dampened effects at the C-terminus (greatest % change cluster 10, position 18=-19.67 [-21.74, -17.80]). Like cluster 5, cluster 10 both contains many positively-charged and hydrophobic residues and had its most dampened effect at position 18. However, cluster 10's effect size is in the opposite direction, likely due to cluster 10's presence of tryptophan.

To learn whether the effects of motifs were additive, we ran an experiment where we embedded motifs in a cluster in non-overlapping positions between 1 and 4 times. Simple additive effects explained nearly all of the variance observed in model predictions (range of R^2 between multiplicity and importance = [0.971, 0.999]; Fig. 6B). However, some clusters

showed evidence of heteroskedasticity in their importance values, which may indicate minor epistatic or background-specific grammars. Accordingly, we used GIA to perform an experiment similar to the one determining additivity of motifs, however we focused on the case in which there are only two motifs, and the embedded motifs are selected from different clusters (Methods). This enables us to learn how interactions between clusters affect nucleation scores. Every cluster showed at least 8 statistically significant interactions (p-value of paired, two-way t-test $<0.05 / (10 \times 10)$ tests); Fig. 6C; Methods), suggesting the importance of modeling sequence context in the prediction of nucleation status. Nonetheless, cluster interaction effects were modest (ranging from -0.058 to 0.085) compared to cluster main effects (-0.608 to 0.448). Almost all clusters exhibited a self-enhancing effect in which their interaction importance was statistically significantly higher than the expected importance from additively combining each marginal effect (maximum importance cluster 6, 0.063). This was not the case for the mixed-charge disorder cluster (cluster 8; importance -0.012 p-value $<2e-16$) interacting with itself, nor the negative-importance charged, hydrophobic cluster (cluster 10; importance -0.011 p-value $<2e-16$). Interestingly, the hydrophobic and aromatic positive-importance clusters (clusters 1 and 2, respectively), showed positive interaction effects with the mixed-charge disorder cluster (cluster 8), whereas the cysteine and asparagine positive clusters (clusters 3 and 4) showed negative interaction effects with both disorder clusters (clusters 6 and 8; Fig. 6C). This is in line with previous reports wherein disordered regions (like cluster 8) are posited to facilitate amyloid fibril formation in contexts of hydrophobic, core-like regions (like clusters 1 and 2)⁴⁵.

Discussion

Amyloid protein aggregation is a hallmark of many human diseases and a major problem in biotechnology. However, relatively few protein sequences are known to nucleate amyloids under physiological conditions, and this shortage of data likely limits our ability to understand, predict, engineer, and prevent the formation of amyloid fibrils.

Here we have quantified the nucleation of amyloids at an unprecedented scale and used the data to train a fast and interpretable deep learning model of amyloid nucleation. In total we quantified the nucleation rates of $>100,000$ 20-amino-acid peptides, which enabled us to explore a vast range of diverse sequences very different to the small number of known amyloids. We used this unique data set to train an open-source and robust predictor of nucleation, CANYA, an inherently interpretable model whose architecture is inspired by biology. In demonstrating the utility of CANYA, we simultaneously provide a principled evaluation of existing amyloid predictors^{11–14,31} using our own and existing datasets^{26,30,31,37}, serving as a guideline for the community. We also adapted state-of-the-art explainable AI (xAI) techniques from genomic neural networks to the protein space^{28,29,38,46}. This not only reveals insight into the decision-making process of our model, but also illustrates how xAI techniques developed for genomic neural networks can provide intelligible information from neural networks that model protein function.

The consistent performance of CANYA across evaluation tasks suggests CANYA does in fact learn an accurate approximation of the sequence-nucleation landscape, despite only

training on random, synthetic peptides. The performance of previous methods compared to that of hydrophobicity scales suggests that the use of limited dataset sizes and short peptides may have limited the amount of additional nucleation-relevant information these approaches could learn. This underscores the importance of using longer sequences and high-throughput assays to profile previously unexplored regions of the sequence-nucleation landscape.

Using interpretability analyses we identified physicochemical motifs that impact CANYA nucleation scores. Motifs with hydrophobic residues, cysteines and asparagines generally had positive effects on nucleation, whereas negative, mixed charge, and proline-containing motifs had negative effects. Positive-effect motifs comprised sequences reminiscent of those found in beta sheets of known amyloids, whereas negative-effect motifs were enriched in the unstructured disordered regions of known fibrils. We also found position-specific effects of motifs—hydrophobic motifs generally had their strongest effects in the center of constructs, whereas charged, tryptophan-containing motifs were strongest toward the N-terminus and weakest toward the C-terminus. This polarity of sequence effects deserves further attention in future work. Finally, motif effects were mostly additive, with only subtle motif-motif interactions.

While we believe that both CANYA and our dataset represent important advances, we note several limitations of our approach. Primarily, we only tested the nucleation of sequences of 20 amino acids and in one particular context. We show through several evaluations that the information learned by modeling these 20 amino acids can offer accurate predictions of nucleation status across a wide range of protein lengths, however, there likely remains additional predictive power to be harvested by experimentally testing at scale and modeling longer sequences and consequently longer-range interactions. Moreover, we limited our neural network architecture to a relatively simple class of models as our focus was on interpretability. Recent literature suggests that leveraging protein embeddings—in lieu of one-hot encoding sequences—may boost our predictive power^{47–54}, though such an approach will likely pose difficulties when performing post-hoc xAI experiments as done here⁵⁵. Further, our model comprises a modest 65,000 parameters and leverages sparsity despite having over 100,000 sequences on which to learn. Many models of protein structure employ much more complex architectures, with both substantially larger numbers of layers and parameters^{47,52,54,56–58}. Future investigations may build off of the work presented here by generating longer sequences, or exploring more complex architectures.

The pairing of massive scale experimental data generation using random sequences with interpretable models has led to insights into genomic regulatory functions⁵⁹. However, to the best of our knowledge, it has been little utilized in the space of proteins to probe mechanisms beyond short motifs. We believe the approach deserves wider adoption, whenever sequences are functional at sufficient frequencies to allow their identification in practical library sizes.

Finally, we note that CANYA is simply an approximation of our assay—namely, that the learned grammar and predictions of CANYA need only be faithful to the experiment itself and may be partially distinct from the process of nucleation in different experimental or *in vivo* contexts. Nonetheless, given CANYA's performance across tasks, we are optimistic that its predictions and insights will assist in advancing our understanding of how proteins nucleate

to form amyloids. Moreover, we expect the dataset presented here to be used to train and evaluate many additional models, and the predictions and outputs of these models to loop back into additional large-scale experimental explorations of sequence space.

Data Availability

All datasets generated from this study are provided under Gene Expression Omnibus (GEO) accession number GSE268261.

External datasets can be found under their respective repositories, which we list here:

AmyPred https://pmlabstack.pythonanywhere.com/dataset_AMYPredFRL

AmyPro <http://www.amypro.net/>

CPAD <https://web.iitm.ac.in/bioinfo2/cpad2/index.html>

WALTZ-DB <http://waltzdb.switchlab.org/sequences>.

Code availability

CANYA is open-source, free to use, and available at the following link <https://github.com/lehner-lab/canya>.

Author contributions

BB, BL, and MT conceived the project. MT and MM performed computational analyses with input from CR, PKK, BB and BL. TSO performed all experiments. MT conceived and wrote the model with input from CR, PKK, BB and BL. MT, MM, BB and BL wrote the manuscript, with input from all authors.

Competing interests

The authors have declared no competing interests.

Acknowledgements

This work was funded by the La Caixa Research Foundation project 'DeepAmyloids' (LCF/PR/HR21/52410004). Work in the lab of B.L. was also funded by an European Research Council (ERC) Advanced (883742) grant, the Spanish Ministry of Science and Innovation (LCF/PR/HR21/52410004, EMBL Partnership, Severo Ochoa Centre of Excellence), the Bettencourt Schueller Foundation, the AXA Research Fund, Agència de Gestió d'Ajuts Universitaris i de Recerca (AGAUR, 2017 SGR 1322), the CERCA Program/Generalitat de Catalunya and Wellcome (Grant reference: 220540/Z/20/A, 'Wellcome Sanger Institute Quinquennial Review 2021-2026'). A.J.F. was funded by a Ramón y Cajal fellowship (RYC2021-033375-I) financed by the Spanish Ministry of Science and Innovation (MCIN/AEI/10.13039/501100011033) and the European Union (NextGenerationEU/PRTR). Work in the lab of BB was also funded by the Spanish Ministry of Science, Innovation and Universities (PID2021-127761OB-I00, RYC2020-028861-I funded by MCIN/AEI/ 10.13039/501100011033, "ERDF A way of making Europe" and "ESF Investing in your future") and by the European Union (ERC Consolidator, Glam-MAP,

101125484). MT was funded by EMBO Fellowship ALTF 266-2023. P.K. and C.R. were funded by NIH grants R01HG012131 and R01GM149921.

Methods

Plasmid library construction

Libraries of random sequences (NNK1-4) were synthesized by Integrated DNA Technologies (IDT) as ultramers of 20 NNK codons (60 nucleotides). A library containing 400 sequences selected from the previous four random libraries was synthesized as an oligopool by IDT for validation and replication (Extended Data Files 3 and 4). In both cases, sequences were flanked by constant regions of 25 nt upstream and 21 nt downstream for cloning. The NNK ultramers and the replication oligo pool were extended in a 1-cycle PCR (Q5 high-fidelity DNA polymerase, NEB) with primers TSO_2 and TSO_65 (Extended Data File 5). The resulting products were treated with 2 µl/tube of ExoSAP (ExoSAP-IT, Applied Biosystems) for 30 minutes at 37 °C and 20 minutes at 80 °C and purified through a MinElute column (Qiagen). In parallel, the PCUP1-Sup35N plasmid was linearized by PCR (Q5 high-fidelity DNA polymerase, NEB; primers TSO_3 and TSO_4, Extended Data File 5). The products were purified from a 1% agarose gel (QIAquick Gel Extraction Kit, Qiagen) and ligated by Gibson with 3 h of incubation at 50°C followed by dialysis for 3 h on a membrane filter (MF-Millipore 0.025 µm membrane, Merck) and vacuum concentration. The resulting (NNK1-4) libraries were transformed into 10-beta Electrocompetent *E. coli* (NEB), by electroporation with 2.0 kV, 200 Ω, 25 µF (BioRad GenePulser machine). Cells were recovered in SOC medium for 30 min and grown overnight in 50 ml of LB ampicillin medium. A small amount of cells was also plated on LB ampicillin plates to assess transformation efficiency. Total transformants were estimated (Extended Data File 6), 50 ml of overnight culture were harvested to purify each library with a midi prep (Plasmid MIDI Kit, Qiagen). Libraries NNK1-4 were bottlenecked to ~1 million transformants, while for the replication library we estimated 625,000 transformants.

Large-scale yeast transformation of random libraries

Saccharomyces cerevisiae GT409 [psi-pin-] (MAT α ade1–14 his3 leu2-3,112 lys2 trp1 ura3–52) provided by the Chernoff lab was used in all experiments in this study²⁴. Yeast cells were transformed with the above plasmid library midipreps. After an overnight pre-growth culture in 25 ml of YPDA medium at 30°C, cells were diluted to OD₆₀₀ = 0.3 in 175 ml YPDA and incubated at 30°C 200 rpm for ~4 hr. When cells reached the exponential phase, they were harvested, washed with milliQ, and resuspended in sorbitol mixture (100 mM LiOAc, 10 mM Tris pH 8, 1 mM EDTA, 1M sorbitol). After a 30 min incubation at room temperature (RT), 4 µg of plasmid library and 175 µl of ssDNA (UltraPure, Thermo Scientific) were added to the cells. PEG mixture (100 mM LiOAc, 10 mM Tris pH 8, 1 mM EDTA pH 8, 40% PEG3350) was also added and cells were incubated for 30 min at RT and heat-shocked for 15 min at 42°C in a water bath. Cells were harvested, washed, resuspended in 250 ml recovery medium (YPD, sorbitol 0.5M, 70 mg/L adenine) and incubated for 1.5 hr at 30°C 200 rpm. After recovery, cells were resuspended in 350 ml -URA plasmid selection medium and allowed to grow for 50 hr. Transformation efficiency was calculated for each of the four

transformations by plating an aliquot of cells in -URA plates (Extended Data File 6). Two days after transformation, the culture was diluted to OD₆₀₀ = 0.08 in 500 ml -URA medium and grown until exponential phase. At this stage, cells were harvested and stored at -80°C in 25% glycerol. In yeast, libraries NNK1-4 were bottlenecked to 0.5-1 million transformants (Extended Data File 6).

Small-scale yeast transformation of replication library

Yeast cells were transformed with the library containing 400 sequences in three biological replicates. An individual colony was grown overnight in 3 ml YPDA medium at 30 °C and 4 g. Cells were diluted in 60 ml to OD₆₀₀ = 0.25 and grown for 4–5 h. When cells reached the exponential phase (OD~0.7–0.8), cells were harvested at 400 × g for 5 min, washed with milliQ, and resuspended in 1 ml YTB (100 mM LiOAc, 10 mM Tris pH 8.0, 1 mM EDTA). They were harvested again and resuspended in 72 µl YTB. 100 ng of plasmid library were added to the cells, together with 8 µl of salmon sperm DNA (UltraPure, Thermo Scientific) previously boiled, 60 µl of dimethyl sulfoxide (Merck) and 500 µl of YTB-PEG (100 mM LiOAc, 10 mM Tris pH 8.0, 1 mM EDTA, 40% PEG 3350). The cells incubated at room temperature for 30 minutes at 4g. Heat-shock was performed at 42 °C for 14 min in a thermo block. Finally, cells were harvested and resuspended in 50 ml plasmid selection medium (-URA, 20% glucose), allowing them to grow for 50 h at 30 °C and 4 g. A small amount of cells was also plated in plasmid selection medium to assess transformation efficiency. We estimated 70,000 transformants per replicate (Extended Data File 6). Two days after transformation, the culture was diluted to OD₆₀₀ = 0.08 in 500 ml -URA medium and grown until exponential phase. At this stage, cells were harvested and stored at -80°C in 25% glycerol.

Selection experiments

Cells were thawed from -80 °C in 50 ml plasmid selection medium at OD = 0.05 and grown until exponential for 15 h. At this stage, cells were harvested and resuspended in 300 ml protein induction medium (-URA, 2% glucose, 100 µM Cu₂SO₄) at OD = 0.1. After 24 h the 250 ml input pellets were collected, and cells were plated on -ADE-URA selection medium in 145-cm² plates (Nunc, Thermo Scientific). Plates were incubated at 30 °C for 7 days. Finally, colonies were scraped off the plates with PBS 1x and harvested by centrifugation to collect the output pellets. Both input and output pellets were stored at -20 °C before DNA extraction. For each random library experiment, one input sample and three technical replicates of the output pellet were processed for sequencing. Selection experiments for the replication library were instead performed in three biological replicates, following the same steps as above. Three input and three output samples were processed for sequencing.

DNA extraction and sequencing library preparation

Input and output pellets were thawed and resuspended in 1.5 ml extraction buffer (2% Triton-X, 1% SDS, 100 mM NaCl, 10 mM Tris pH 8, 1 mM EDTA pH 8), and underwent two cycles of freezing and thawing in an ethanol-dry ice bath (10 min) and at 62°C (10 min). Samples were then vortexed together with 1.5 ml of phenol:chloroform:isoamyl 25:24:1 and 1.5 g of glass beads (Sigma). The aqueous phase was recovered by centrifugation and mixed again with 1.5 ml phenol:chloroform:isoamyl 25:24:1. DNA precipitation was performed by adding 1:10 V of 3M NaOAc and 2.2 V of 100% cold ethanol to the aqueous phase and incubating the samples at -20°C for 1 hr. After a centrifugation step, pellets were

dried overnight at RT. Pellets were resuspended in 900 μ l resuspension buffer (10 mM Tris pH 8, 1 mM EDTA pH 8) and treated with 7.5 ml RNase A (Thermo Scientific) for 30 min at 37°C. The DNA was finally purified using 30 μ l of silica beads (QIAEX II Gel Extraction Kit, Qiagen), washed and eluted in 22 μ l of elution buffer. Plasmid concentrations were measured by quantitative PCR with SYBR green (Merck) and primers annealing to the origin of replication site of the PCUP1-Sup35N plasmid at 58 °C for 40 cycles (TSO_05 and TSO_06, Extended Data File 5). The library for high-throughput sequencing was prepared in a two-step PCR (Q5 high-fidelity DNA polymerase, NEB). In PCR1, 160 million of molecules were amplified for 15 cycles at 68°C with frame-shifted primers with homology to Illumina sequencing primers (primers TSO_7 to TSO_20, Extended Data File 5). The products were purified with ExoSAP treatment (Affymetrix) and by column purification (MinElute PCR Purification Kit, Qiagen). They were then amplified for 10 cycles in PCR2 with Illumina-indexed primers (primers TSO_21 to TSO_54, Extended Data File 5). The library was sequenced by 150 bp paired-end sequencing in an Illumina NextSeq500 sequencer at the CRG Genomics core facility.

Sequence data preprocessing

We processed each of the 4 NNK experiments separately using DiMSum⁶⁰. Briefly, DiMSum comprises an end-to-end pipeline for processing deep mutational scanning datasets from raw reads to measured sequences and their associated assay scores (plus errors). DiMSum was run with the following parameters: `cutadaptMinLength="60"`; `cutadaptErrorRate="0.2"`; `vsearchMinQual="30"`; `vsearchMaxee="0.5"`; `startStage="0"`; `fitnessMinInputCountAny="0"`; `maxSubstitutions="20"`; `mixedSubstitutions="TRUE"`; `experimentDesignPairDuplicates="TRUE"`. We then removed sequences with fewer than 100 reads in the input sequencing experiment. Next, we centered the fitness estimates (nucleation scores) of each dataset individually by adding or subtracting the corresponding mode fitness of the non-nucleating sequences. After centering each sequence, we next labeled sequences as “nucleators” (or “non-nucleators”) by transforming their fitness estimate to a Z-score composed of the fitness estimate scaled by the DiMSum error, and performing a one-sided hypothesis test to check whether standardized score was significantly larger than 0. We treated sequences whose p-values after FDR adjustment were ≤ 0.05 as “nucleators”, and remaining sequences as “non-nucleators.” A proportion of sequences produced no reads after the selection experiments, thus leading to NA scores from DiMSum. We labeled these sequences as “non-nucleators.” If a sequence contained a stop codon, we used only the component of the sequence preceding the stop for model training. For cases in which this resulted in duplicate sequences (e.g. FN*VILRDEGHGSYGFNNN and FN*FVVMHTCIMVVFCLGDI are both mapped to “FN”), we summarized the truncated sequence by taking its mean nucleation score or mode nucleation status across observations. If a given truncated sequence had an equal number of nucleator and non-nucleator status observations, we discarded this truncated sequence. As a result we classified > 35,000 sequences for libraries NNK1-3 (35,456; 37,578; 38,893 respectively) and 7,040 for NNK4.

The architecture of CANYA

CANYA is a biologically motivated hybrid-neural network designed to discover motifs and their interactions. More concretely, the architecture of CANYA is inspired by recent work that suggests stacked convolution and attention layers serve as a reasonable inductive bias for motif and motif-interaction discovery. The hyperparameters of CANYA were influenced by

summary statistics of interacting secondary structure elements in amyloids within the PDB (Supp. Fig. 4). Summarily, we chose the simplest architecture of our model such that it is expressive, interpretable, and importantly, principled in biological knowledge.

CANYA takes as input an amino acid sequence of length limit up to 145 residues, and outputs a score related to the sequence's propensity to form amyloids. Prior to passing the sequence to the input layer, we first one-hot encode it, allowing only the 20 canonical amino acids. As we use filters of length 3 (See below for justification; Supplementary Fig. 4), we pad the sequence with two 0s both up- and downstream the sequence. Finally, if this padded sequence is not of length 149, we add a mask with values of -1 downstream the sequence until it reaches length 149. The input length restrictions of CANYA arise from the fact that a given sequence in the assay is fused to a Sup35N construct of length 125, is (up to) length 20, and is padded with two 0s on each side. Explicitly, the training data of CANYA looks as follows:

`00[one-hot encoded Sup35N][one-hot encoded random sequence]00`

when there is no masking or stop codons, and as follows if so:

`00[one-hot encoded Sup35N][one-hot encoded random sequence]00[-1]`

where the number of -1 values is the required quantity such that the sequence is length 149.

The input layer of CANYA correspondingly accepts a matrix of size 149x20 representing a one-hot encoded, padded, and potentially masked peptide sequence. The output layer is a single unit with sigmoid activation. The hidden layers of CANYA are:

1. Convolution (100 filters, size 3, stride 1, exponential activation)
2. Self-attention (1 attention head, key-length 6)
3. Fully-connected layer (64 units, ReLU activation)

We selected an exponential activation function for the convolutional layer as this type of activation is generally more robust for motif discovery⁶¹. We chose filters of length 3 as this was the mode length of beta-sheets in amyloid sequences with resolved structures in Uniprot (Supp. Fig. 4). We utilize dropout with probability 0.1 after the convolution and attention layers, and 0.4 after the fully-connected layer. We use an elastic net regularization (with value 0.01) when learning the weights between the attention and fully-connected layers. Finally, to encourage the model to learn positional information, we do not perform pooling after the convolution layer, and we include positional encodings prior to taking the softmax in the attention layer. We trained CANYA for 100 epochs using the adam optimizer with default values and the binary Kullback-Leibler divergence as a loss function. We limited the learning rate of the model during training by monitoring the validation area under precision-recall curve, decaying at a factor of 0.2 with patience 4, and performed early stopping by monitoring the validation area under precision-recall curve with patience 10. For sequences with length greater than 145, we collect the CANYA score at every overlapping length-145 window of the sequence, then use its minimum CANYA score as its final score (under the logic that nucleation-forming propensity is limited by a sequence's most nucleation-disrupting region).

Compilation of external datasets

We first collected 6-mers from the WALTZ-DB dataset²⁶. Here, we assigned all sequences whose "Classification" field was "amyloid" as a 1, and all other sequences as 0. We next

collected the collection of aggregating peptides from the CPAD repository³⁰. We used sequences from the “Peptide” field, filtering for sequences of at least length 10 and for sequences that did not contain a space in their sequences. We assigned sequences with “Classification” field “Amyloid” a 1, and all other sequences 0. We then collected the Amy dataset from the AMYPred-FRL server³¹. Here, we assigned all sequences in the negative sets a label of 0, and all sequences in the positive sets a label of 1. Many of the sequences had lengths greater than 145, we therefore applied a sliding window approach to these sequences in which we score every overlapping length-145 region of a sequence, and assigned a final score to the entire sequence as the minimum of the length-145 regional scores. The final external dataset we used was from the AmyPro database³⁷. Though AmyPro contains overlapping sequences with the Amy dataset, we treated this task differently than the previous tasks. Namely, all sequences in the AmyPro dataset were amyloids, and so we sought to evaluate methods’ abilities to distinguish the amyloidogenic region from the non-amyloidogenic regions of the sequences. First, we collected all sequences from the “regions” field in the dataset. Next we removed each of these “region” sequences from the main peptide sequence and concatenated the remaining two portions of the main sequence together, comprising a set of positive sequences (labeled 1) from the “region” field and negative sequences (labeled 0) from the remaining peptide sequences. Finally, we limited the length of all sequences to 100 by breaking sequences longer than length 100 into non-overlapping subsequences of at most length 100. While this task evaluates un-natural sequences, it evaluates the ability of each method to distinguish amyloid cores from non-amyloid cores while also making the problem more amenable to previous approaches, which generally underperformed on long sequences. We list descriptive summary statistics (e.g. length, sample sizes, hydrophobicity) in Supplementary Table 5.

Aggregation predictors

Aggregation predictors or physicochemical scales (Tango¹¹, Amypred³¹, CamSol¹², PLAAC¹³, Aggrescan¹⁴) were used to calculate a score for each sequence. When appropriate, individual residue-level scores were summed to obtain a single score per sequence. CamSol, Amypred and Aggrescan were run with the default parameters. PLAAC was run using a core of length 6 and weightings from input sequences. Tango was run with pH 7.2, no protection of termini, ionic strength = 0.1 and T = 298K (25°C). Some of the predictors present sequence length limitations: Amypred runs only for sequences longer than 10 amino acids, CamSol for sequences longer than 6 amino acids, and Aggrescan cannot be run for sequences longer than 2004 amino acids.

Selecting a model for interpretability analyses

We trained CANYA with random weight initialization 100 times and recorded for each fitted model the area under the curve (AUC) of the test data, area under the precision-recall curve (AUPROC) of the test data, and interpretability score adapted from a recently developed approach for interpretability analyses of genomic neural networks²⁹. Briefly, Majdandzic et al. propose an approach to quantify the consistency of the attribution maps of a trained model by comparing the entire set of kmers in the training sequences to the set of kmers in (adjusted⁴⁶) attributed positions in the training sequences. These two distributions of kmers—in the case of CANYA, 3-mers—are compared using the Kullbeck-Leibler (KL) divergence, where a higher KL divergence suggests greater amenability to downstream interpretability analyses. To calculate an interpretability score for each trained instance of

CANYA, we used this same approach, but rather than using kmers of nucleotides, we used kmers from the input amino acids. As we saw that the test AUPROC was more consistent across experiments, we used a models' mean AUPROC across experiments and interpretability score as model selection criteria. More rigorously, we selected the model with the highest interpretability score, conditional on the fact that its mean AUPROC across datasets was greater than the median of these mean scores across model training instances.

Visualization of filters (motifs)

Notably, the use of random sequences in amino acid space poses difficulties for observing a typical, lexicographic motif, and consequently, observing convergence toward a lexicographic motif in first-layer convolutional filters. We elaborate as follows: using a filter length of 3, there is a 1 in 8,000 (20^3) chance of observing a given kmer. Ideally, for the model to learn a stable feature, this kmer must not only exist in a sizable proportion of sequences, but its effect must also not be masked out by surrounding contextual information. Even if we were to ignore contextual information, this motif would need to occur independently multiple times, an event whose probability quickly converges to 0. Consequently, we are much stricter than previous approaches when generating a position weight matrix (PWM) for a given filter. For interpretability's sake, we limit the kmers comprising a PWM for a filter to the minimum of either the 10 most-activating kmers of a filter, or the collection of kmers whose activation is at least 75% of the maximum-activating kmer. Summarily, a filter is both visualized and represented numerically by its PWM composed of at most the top 10 strongest activating kmers.

Motif clustering

Following the above logic, CANYA must learn physicochemical properties of amino acids and understand how these properties interact amongst each other when constructing its features at the convolution layer. Moreover, these physicochemical 3-mers, or motifs, may often capture redundant physicochemical information, but independent sequences—for example, two different motifs capturing hydrophobicity may separately comprise sequences of “IVF” or “ALM.” To further improve interpretability and reduce the dimensionality of downstream experiments leveraging the learned motifs of CANYA, we performed clustering on the PWM matrices. More concretely, we calculated BLOSUM scores for each filter by taking the dot product between its PWM and BLOSUM score matrix⁴². We next performed affinity propagation on these calculated motif BLOSUM scores to cluster the motifs. Affinity propagation discovered 10 clusters of motifs. However, after performing Global Importance Analysis (GIA) experiments³⁸, we found 7 discrepancies when evaluating whether a given motif had the same effect size (importance score) direction compared to the effect size of the motif with the greatest absolute effect within the cluster. As our goal was to interpret model decisions and physicochemical clusters, we removed these 7 filters from their corresponding clusters so that each cluster contained only filters with the same effect size direction. We show the original and changed cluster assignments in Supplementary Figure 6.

Global Importance Analysis (GIA) Experiments

To learn the effect of motif presence on CANYA's decision-making, we turned to Global Importance Analysis (GIA) *in-silico* experiments³⁸. Briefly, GIA is a post-hoc interpretability method applied to genomic neural networks that enables users to learn importance scores (i.e. effect sizes) of a given sequence feature on a model's output score. The importance

score is derived from taking the average difference in model score between a set of background sequences, and this same set of background sequences but with a functional element, such as a motif, placed in the background sequence (sequence length is maintained, i.e. a window of the sequence is replaced by the functional element). For all experiments, we limited our analyses to 25,000 randomly selected, full-length (length-20 and absent of stop codons) training sequences that were confidently predicted by CANYA. We defined “confidently predicted” as nucleators with CANYA score above 0.3 and non-nucleators with CANYA score below 0.2 (see Supp. Fig. 8 for prediction score distributions). Finally, we emphasize that owing to the random nature of our experiment, the training sequences serve as a valid set of background sequences for GIA as they span an extremely wide range of contexts.

In the first set of GIA experiments, we sought to characterize the importance score of each filter individually. To do so, we first randomly selected 25,000 sequences from the training set, comprising sequences from across all three experiments. Next, for a given filter, we collected the activation energy of each kmer used to represent the PWM, and used the ratio of the activation energy of each kmer to the activation energy of the kmer with the maximum activation energy in this PWM to generate kmer sampling probabilities. For each sequence, we randomly sampled one kmer using this normalized ratio as the kmer’s sampling probability, and embedded this kmer into the sequence. Afterward, we calculated for all 25,000 background sequences and all 25,000 modified sequences the CANYA nucleation score prior to applying the softmax function. We calculated each filter’s importance score as the mean paired difference in scores between the 25,000 background and modified sequences.

After clustering the learned motifs, we next wished to validate whether the clusters could be utilized to simplify further interpretability analyses by reducing the scale of *in-silico* experiments performed. To do so, we conducted a GIA experiment within each cluster to determine a cluster-level importance score. The experiment follows the same logic as the original, filter-level GIA experiment, only that we first randomly selected a filter within a cluster prior to sampling a kmer from its PWM. The filters were randomly selected according to the ratio of their absolute GIA importance score to the maximum absolute GIA importance score across filters of the corresponding cluster. Indeed, cluster-level scores recapitulated the scores of the motifs from which they were composed (Supp. Table 5). We therefore performed all following GIA analysis at the cluster level, using this filter-first, kmer-second sampling scheme.

We next performed an experiment to evaluate the additivity of motif-clusters on nucleation propensity. Here, we collected 25,000 background sequences from the training dataset, then embedded into these background sequences 1 to 4 kmers in non-overlapping positions where each of the 4 kmers was sampled using the filter-first, kmer-second sampling scheme. Each sequential kmer addition (from kmers 2-4) was embedded in the sequence such that the sequence with antecedent kmer multiplicity maintained the kmer(s) at its (their) original embedded position(s). We calculated the cluster importance score for a given multiplicity by taking the mean difference in prediction score between the sequences with the injected kmer(s) and their corresponding background sequences—in other words, each importance score is generated by taking the mean difference between 25,000 background sequences and 25,000 modified background sequences with either 1, 2, 3, or 4 embedded kmers.

To evaluate whether CANYA learned position-specific importance of motifs, we performed an additional GIA experiment in which we systematically embedded a motif-cluster at each position of a random sequence. In these experiments, we performed a single GIA experiment with 25,000 background sequences and 25,000 modified sequences for each position from positions 1-18, so that the entire 3-mer could be contained within the sequence.

In a final GIA experiment, we characterized interaction effects between motif-clusters. For a given motif-cluster pair, we sampled a kmer (as mentioned above) from each cluster as well as a corresponding position randomly from positions 1-18 in which to embed each kmer. We evaluated the CANYA score for the background sequence, the background sequence with the kmer from the first cluster at the first sampled position, the background sequence with the kmer from the second cluster at the second sampled position, and the background sequence with both kmers at both positions. We called the interaction importance as the result of subtracting the sum of CANYA predictions of the sequences with each marginal kmer embedding from the sum of the CANYA predictions of the background sequence and sequence with both motifs. The final importance was calculated as the mean interaction importance across 25,000 sequences.

Secondary structure enrichment scoring of motifs

To examine whether certain motifs were characteristically similar to sequences found in specific secondary elements of amyloids, we examined activation energies of filters across secondary structure elements in a set of amyloids with resolved structures in the PDB. Concretely, we collected 114 entries from the STAMP dataset⁴³, then downloaded their structural information from the PDB (see Supp. Table 7 for entries and corresponding proteins). Next, we passed all sequences through CANYA, and extracted their filter activation energies (i.e., output from the convolution layer). At each position, we summarized a cluster's activation energies as the maximum activation energy across filters within a cluster, generating a vector of maximum activation energies for each cluster. Next, we encoded each secondary structure (coil, beta strand, or disorder) as a binary vector where 1 indicated positions in the corresponding secondary structure, and 0 indicated otherwise. We collected this set of secondary structure vectors and activation energy vectors for all sequences, then concatenated them across sequences. Finally, we generated secondary structure enrichment scores by calculating the AUC between a given secondary structure element and cluster activation energy across all sequences.

Extended Data

Extended Data File 1 All sequences recorded spanning each experiment with reported fitnesses, error, and nucleation status

Extended Data File 2 Sequences used to train and test CANYA

Extended Data File 3 Sequences used in replication experiments with their original measured fitness and fitness from the replication experiment

Extended Data File 4 Validation sequences and their corresponding nucleotide sequences

Extended Data File 5 Oligo pool and primer sequences for the NNK experiments

Extended Data File 6 Transformants measured across each experiment

References

1. Chiti, F. & Dobson, C. M. Protein Misfolding, Amyloid Formation, and Human Disease: A Summary of Progress Over the Last Decade. *Annu. Rev. Biochem.* **86**, 27–68 (2017).
2. Fowler, D. M., Koulov, A. V., Balch, W. E. & Kelly, J. W. Functional amyloid--from bacteria to humans. *Trends Biochem. Sci.* **32**, 217–224 (2007).
3. Ke, P. C. *et al.* Half a century of amyloids: past, present and future. *Chem. Soc. Rev.* **49**, 5473–5509 (2020).
4. Dobson, C. M., Knowles, T. P. J. & Vendruscolo, M. The Amyloid Phenomenon and Its Significance in Biology and Medicine. *Cold Spring Harb. Perspect. Biol.* **12**, (2020).
5. Scheres, S. H. W., Ryskeldi-Falcon, B. & Goedert, M. Molecular pathology of neurodegenerative diseases by cryo-EM of amyloids. *Nature* **621**, 701–710 (2023).
6. Sabaté, R. & Ventura, S. Cross- β -sheet supersecondary structure in amyloid folds: techniques for detection and characterization. *Methods Mol. Biol.* **932**, 237–257 (2013).
7. Eisenberg, D. S. & Sawaya, M. R. Structural Studies of Amyloid Proteins at the Molecular Level. *Annu. Rev. Biochem.* **86**, 69–95 (2017).
8. Shi, Y. *et al.* Structure-based classification of tauopathies. *Nature* **598**, 359–363 (2021).
9. Yang, Y. *et al.* Cryo-EM structures of amyloid- β 42 filaments from human brains. *Science* **375**, 167–172 (2022).
10. Schweighauser, M. *et al.* Structures of α -synuclein filaments from multiple system atrophy. *Nature* **585**, 464–469 (2020).
11. Fernandez-Escamilla, A.-M., Rousseau, F., Schymkowitz, J. & Serrano, L. Prediction of sequence-dependent and mutational effects on the aggregation of peptides and proteins. *Nat. Biotechnol.* **22**, 1302–1306 (2004).
12. Sormanni, P., Aprile, F. A. & Vendruscolo, M. The CamSol method of rational design of protein mutants with enhanced solubility. *J. Mol. Biol.* **427**, 478–490 (2015).
13. Lancaster, A. K., Nutter-Upham, A., Lindquist, S. & King, O. D. PLAAC: a web and command-line application to identify proteins with prion-like amino acid composition. *Bioinformatics* **30**, 2501–2502 (2014).
14. Conchillo-Solé, O. *et al.* AGGRESCAN: a server for the prediction and evaluation of ‘hot spots’ of aggregation in polypeptides. *BMC Bioinformatics* **8**, 65 (2007).
15. Wickner, R. B. *et al.* Yeast Prions Compared to Functional Prions and Amyloids. *J. Mol. Biol.* **430**, 3707–3719 (2018).
16. Wickner, R. B. Yeast and Fungal Prions. *Cold Spring Harb. Perspect. Biol.* **8**, (2016).
17. Wilkinson, M. *et al.* Structural evolution of fibril polymorphs during amyloid assembly. *Cell* **186**, 5798–5811.e26 (2023).
18. Lövestam, S. *et al.* Disease-specific tau filaments assemble via polymorphic intermediates. *Nature* **625**, 119–125 (2024).
19. Knowles, T. P. J., Vendruscolo, M. & Dobson, C. M. The amyloid state and its association with protein misfolding diseases. *Nat. Rev. Mol. Cell Biol.* **15**, 384–396 (2014).
20. Baldwin, A. J. *et al.* Metastability of native proteins and the phenomenon of amyloid formation. *J. Am. Chem. Soc.* **133**, 14160–14163 (2011).
21. Navarro, S. & Ventura, S. Computational methods to predict protein aggregation. *Curr. Opin. Struct. Biol.* **73**, 102343 (2022).

22. Seuma, M., Faure, A. J., Badia, M., Lehner, B. & Bolognesi, B. The genetic landscape for amyloid beta fibril nucleation accurately discriminates familial Alzheimer's disease mutations. *Elife* **10**, (2021).
23. Seuma, M., Lehner, B. & Bolognesi, B. An atlas of amyloid aggregation: the impact of substitutions, insertions, deletions and truncations on amyloid beta fibril nucleation. *Nat. Commun.* **13**, 7084 (2022).
24. Chandramowlishwaran, P. *et al.* Mammalian amyloidogenic proteins promote prion nucleation in yeast. *J. Biol. Chem.* **293**, 3436–3450 (2018).
25. Louros, N., Orlando, G., De Vleeschouwer, M., Rousseau, F. & Schymkowitz, J. Structure-based machine-guided mapping of amyloid sequence space reveals uncharted sequence clusters with higher solubilities. *Nat. Commun.* **11**, 3314 (2020).
26. Louros, N. *et al.* WALTZ-DB 2.0: an updated database containing structural information of experimentally determined amyloid-forming peptides. *Nucleic Acids Res.* **48**, D389–D393 (2020).
27. Ullah, F. & Ben-Hur, A. A self-attention model for inferring cooperativity between regulatory features. *Nucleic Acids Res.* **49**, e77 (2021).
28. Ghotra, R. S., Lee, N. K. & Koo, P. K. Uncovering motif interactions from convolutional-attention networks for genomics. *NeurIPS 2021 AI for Science Workshop* (2021).
29. Majdandzic, A. *et al.* Selecting deep neural networks that yield consistent attribution-based interpretations for genomics. *Proc Mach Learn Res* **200**, 131–149 (2022).
30. Thangakani, A. M. *et al.* CPAD, Curated Protein Aggregation Database: A Repository of Manually Curated Experimental Data on Protein and Peptide Aggregation. *PLoS One* **11**, e0152949 (2016).
31. Charoenkwan, P. *et al.* AMYPred-FRL is a novel approach for accurate prediction of amyloid proteins by using feature representation learning. *Sci. Rep.* **12**, 7697 (2022).
32. Kyte, J. & Doolittle, R. F. A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **157**, 105–132 (1982).
33. Thangakani, A. M., Kumar, S., Velmurugan, D. & Gromiha, M. M. Distinct position-specific sequence features of hexa-peptides that form amyloid-fibrils: application to discriminate between amyloid fibril and amorphous β -aggregate forming peptide sequences. *BMC Bioinformatics* **14 Suppl 8**, S6 (2013).
34. Thangakani, A. M., Kumar, S., Nagarajan, R., Velmurugan, D. & Gromiha, M. M. GAP: towards almost 100 percent prediction for β -strand-mediated aggregating peptides with distinct morphologies. *Bioinformatics* **30**, 1983–1990 (2014).
35. Wozniak, P. P. & Kotulska, M. AmyLoad: website dedicated to amyloidogenic protein fragments. *Bioinformatics* **31**, 3395–3397 (2015).
36. Niu, M., Li, Y., Wang, C. & Han, K. RFAmyloid: A Web Server for Predicting Amyloid Proteins. *Int. J. Mol. Sci.* **19**, (2018).
37. AmyPro database.
38. Koo, P. K., Majdandzic, A., Ploenzke, M., Anand, P. & Paul, S. B. Global importance analysis: An interpretability method to quantify importance of genomic features in deep neural networks. *PLoS Comput. Biol.* **17**, e1008925 (2021).
39. Sawaya, M. R., Hughes, M. P., Rodriguez, J. A., Riek, R. & Eisenberg, D. S. The expanding amyloid family: Structure, stability, function, and pathogenesis. *Cell* **184**, 4857–4873 (2021).
40. Murray, K. A. *et al.* Identifying amyloid-related diseases by mapping mutations in

- low-complexity protein domains to pathologies. *Nat. Struct. Mol. Biol.* **29**, 529–536 (2022).
41. Kanchi, P. K. & Dasmahapatra, A. K. Polyproline chains destabilize the Alzheimer's amyloid- β protofibrils: A molecular dynamics simulation study. *J. Mol. Graph. Model.* **93**, 107456 (2019).
 42. Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. U. S. A.* **89**, 10915–10919 (1992).
 43. Louros, N., van der Kant, R., Schymkowitz, J. & Rousseau, F. StAmP-DB: a platform for structures of polymorphic amyloid fibril cores. *Bioinformatics* **38**, 2636–2638 (2022).
 44. Izawa, Y. *et al.* Role of C-terminal negative charges and tyrosine residues in fibril formation of α -synuclein. *Brain Behav.* **2**, 595–605 (2012).
 45. Tompa, P. Structural disorder in amyloid fibrils: its implication in dynamic interactions of proteins. *FEBS J.* **276**, 5406–5415 (2009).
 46. Majdandzic, A., Rajesh, C. & Koo, P. K. Correcting gradient-based interpretations of deep neural networks for genomics. *Genome Biol.* **24**, 109 (2023).
 47. Li, F.-Z., Amini, A. P., Yue, Y., Yang, K. K. & Lu, A. X. Feature Reuse and Scaling: Understanding Transfer Learning with Protein Language Models. *bioRxiv* 2024.02.05.578959 (2024) doi:10.1101/2024.02.05.578959.
 48. Flamholz, Z. N., Biller, S. J. & Kelly, L. Large language models improve annotation of prokaryotic viral proteins. *Nat Microbiol* **9**, 537–549 (2024).
 49. Thumhuri, V., Almagro Armenteros, J. J., Johansen, A. R., Nielsen, H. & Winther, O. DeepLoc 2.0: multi-label subcellular localization prediction using protein language models. *Nucleic Acids Res.* **50**, W228–W234 (2022).
 50. Teufel, F. *et al.* SignalP 6.0 predicts all five types of signal peptides using protein language models. *Nat. Biotechnol.* **40**, 1023–1025 (2022).
 51. Detlefsen, N. S., Hauberg, S. & Boomsma, W. Learning meaningful representations of protein sequences. *Nat. Commun.* **13**, 1914 (2022).
 52. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
 53. Elnaggar, A. *et al.* ProtTrans: Toward Understanding the Language of Life Through Self-Supervised Learning. *IEEE Trans. Pattern Anal. Mach. Intell.* **44**, 7112–7127 (2022).
 54. Yang, K. K., Fusi, N. & Lu, A. X. Convolutions are competitive with transformers for protein sequence pretraining. *Cell Syst* **15**, 286–294.e2 (2024).
 55. Tang, Z. & Koo, P. K. Evaluating the representational power of pre-trained DNA language models for regulatory genomics. *bioRxiv* (2024) doi:10.1101/2024.02.29.582810.
 56. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
 57. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* **630**, 493–500 (2024).
 58. Lin, Z. *et al.* Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv* 2022.07.20.500902 (2022) doi:10.1101/2022.07.20.500902.
 59. Liao, S. E., Sudarshan, M. & Regev, O. Deciphering RNA splicing logic with interpretable machine learning. *Proc. Natl. Acad. Sci. U. S. A.* **120**, e2221165120 (2023).
 60. Faure, A. J., Schmiedel, J. M., Baeza-Centurion, P. & Lehner, B. DiMSum: an error

model and pipeline for analyzing deep mutational scanning data and diagnosing common experimental pathologies. *Genome Biol.* **21**, 207 (2020).

61. Koo, P. K. & Ploenzke, M. Improving representations of genomic sequence motifs in convolutional networks with exponential activations. *Nat Mach Intell* **3**, 258–266 (2021).