

Self-assembling Peptide Discovery: Overcoming Human Bias With Machine Learning

Rohit Batra

Argonne National Lab

Troy Loeffler

SLAC National Accelerator Laboratory

Henry Chan

Argonne National Laboratory <https://orcid.org/0000-0002-8198-7737>

Srilok Srinivasan

Argonne National Lab

Honggang Cui

Johns Hopkins University <https://orcid.org/0000-0002-4684-2655>

Ivan Korendovych

Syracuse University

Vikas Nanda

Rutgers University

Liam Palmer

Northwestern University <https://orcid.org/0000-0003-0804-1168>

Lee Solomon

George Mason University

Harry Fry (✉ hfy@anl.gov)

Argonne National Laboratory

Subramanian Sankaranarayanan

Argonne National Laboratory

Article

Keywords: peptides, tissue engineering, human bias, machine learning

Posted Date: June 21st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-505801/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Chemistry on October 31st, 2022.
See the published version at <https://doi.org/10.1038/s41557-022-01055-3>.

Self-assembling Peptide Discovery: Overcoming Human Bias With Machine Learning

Rohit Batra¹, Troy D. Loeffler^{1,2}, Henry Chan^{1,2}, Srilok Srinivasan¹, Honggang Cui³, Ivan V. Korendovych⁴, Vikas Nanda⁵, Liam C. Palmer⁶, Lee A. Solomon⁷, H. Christopher Fry^{1,*}, and Subramanian KRS Sankaranarayanan^{1,2,*}

¹Center for Nanoscale Materials, Argonne National Laboratory, Lemont, Illinois 60439, USA

²Department of Mechanical and Industrial Engineering, University of Illinois, Chicago, Illinois 60607, USA

³Department of Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA

⁴Department of Chemistry, Syracuse University, Syracuse, New York 13244, USA

⁵Center for Advanced Biotechnology and Medicine, Rutgers University, Piscataway, New Jersey 08854, USA

⁶Department of Chemistry, Northwestern University, Evanston, Illinois 60208, USA

⁷Department of Chemistry and Biochemistry, George Mason University, Manassas, Virginia 20110, USA

* Author to whom correspondence should be addressed.

May 7, 2021

Abstract

Peptide materials have a wide array of functions from tissue engineering, surface coatings to catalysis and sensing. This class of biopolymer is composed of a sequence, comprised of 20 naturally occurring amino acids whose arrangement dictate the peptide functionality. While it is highly desirable to tailor the amino acid sequence, a small increase in their sequence length leads to dramatic increase in the possible candidates (e.g., from tripeptide = 20^3 or 8,000 peptides to a pentapeptide = 20^5 or 3.2 M). Traditionally, peptide design is guided by the use of structural propensity tables, hydrophobicity scales, or other desired properties and typically yields <10 peptides per study, barely scraping the surface of the search space. These approaches, driven by human expertise and intuition, are not easily scalable and are riddled with human bias. Here, we introduce a machine learning workflow that combines Monte Carlo tree search and random forest, with molecular dynamics simulations to develop a fully autonomous computational search engine (named, AI-expert) to discover peptide sequences with high potential for self-assembly (as a representative target functionality). We demonstrate the efficacy of the AI-expert to efficiently search large spaces of tripeptides and pentapeptides. Subsequent experiments on the proposed peptide sequences are performed to compare the predictability of the AI-expert with those of human experts. The AI performs on-par or better than human experts and suggests several non-intuitive sequences with high self-assembly propensity, outlining its potential to overcome human bias and accelerate peptide discovery.

Nature generates innumerable functional materials in living systems in the form of proteins and their supramolecular assemblies. Examples include collagen (extended triple helices forming the fibrous base component of skin, hairs and nails) [1, 2], silk proteins [3], and light harvesting reaction centers [4]. Investigations into such naturally occurring supramolecular assemblies have inspired design of novel biomolecular materials [5, 6, 7]. For instance, the formation of a plaque nucleating region in neurodegenerative diseases (like Alzheimer’s) has been attributed (but not limited to [8]) to the presence of diphenylalanine (FF) amino acid sequence in the A β peptide [9]. As a result, FF containing peptide sequences have been explored in sensing [10] as biocompatible implants [11], piezoelectrics [12, 13] and for drug release [14]. Similarly, other works have employed self-assembling small peptides (< 10 amino acids) for various chemical and biological applications, such as catalysis, light harvesting and conductivity [15, 16, 17, 18, 19]. Importantly, in all cases the emergent functionality of a peptide is an outcome of its self assembled architecture, with the unique property being lost when there is no assembly. The self-assembled architecture and thereby its functionality depends strongly on the amino acid sequence, that has traditionally relied on derivation from natural sequences, human expertise, experience and intuition. Thus, researchers rationally design novel peptide sequences to either replicate, tailor or improve natural properties or investigate emerging functionalities as a result of their assembled structure [6, 20].

Traditional approaches of peptide design utilize hydrophobicity scales determined from the partitioning of amino acids into hydrophilic or hydrophobic environments (e.g. Wimley-White scale [21, 22]) and secondary structure propensity tables obtained from the occurrence of any given amino acid in an α -helix or β -sheet fold (e.g. Chou-Fasman [23]). This often introduces a bias toward high β -sheet propensity amino acids with moderate to high hydrophobicity (e.g. valine, isoleucine, phenylalanine) in the design of supramolecular peptides. Another source of bias comes from the commonly employed patterning strategies like pnnnp or npnpn (p = polar, n = non-polar) that reliably lead to β -sheet rich nanostructured materials. The principal reason why designers knowingly resort to such (biased) approaches is because the design space of peptides can become exorbitantly large—the possible combinations of peptides equals 20^n , where n is the number of amino acids in the peptide chain and the factor 20 arises from the library of commonly available amino acids, as shown in Figure 1. While short sequences such as tripeptides ($n=3$) with 8000 combinations are somewhat tractable, the move to pentapeptides ($n=5$) opens up nearly 3.2 million possibilities. This not only precludes any rigorous experimental study of the complete peptide design space, but also suggests that a large fraction of the possible peptide sequences remain unexplored. While a brute-force computational search based on coarse-grained molecular dynamics (MD) simulations provide a pathway to overcome this search bias and, notably, has been successful in identifying several self-assembling and hydrogelating tripeptides (from a total of 8000 cases) [24], it cannot be extended to larger sequence lengths ($n > 3$) owing to high computational costs.

A major challenge in peptide design lies in efficiently navigating through this elaborate search space of amino acid sequences and propose a subset with the most promising possibilities. Artificial intelligence (AI) and machine learning (ML) based strategies makes this a reality by balancing the *exploration-vs-exploitation* trade-off [25, 26, 27]. Here, we introduce an ‘AI expert’ that combines recent advances in decision trees (Monte Carlo tree search (MCTS) algorithm [28, 29]) with coarse grained MD simulations to identify pentapeptides with high aggregation propensities (AP) in water; see Figure 1. Operating in an autonomous manner, the AI-expert utilizes the MCTS algorithm to make an informed decision on which peptide sequence(s) to evaluate next using the MD simulations, with the score of the modeled peptide(s) provided as the feedback to guide the future search. In contrast to a brute force approach wherein every possibility is investigated, the MCTS streamlines the search by focusing on the most promising areas of the search space, *i.e.*, with high scoring (exploitation) and diverse (exploration) sequences. An additional performance boost to the MCTS algorithm is provided by introducing a novel concept of uniqueness function within the MCTS objective function, and by utilizing a random forest (RF) based surrogate model to bypass some of the expensive MD simulation evaluations. Inspired by the past work on the design of di- and tri-peptides [30, 24], our scoring system consists of the solvent accessible surface areas (SASA) and the Wimley-White scale, respectively to quantify the computational AP and hydrophobicity of a peptide. While the former is based on the structure of a peptide obtained only after time-intensive MD simulation, the latter is compute cost effective and can be evaluated instantly given only the peptide sequence.

Of the 3.2 million possible pentapeptides, the AI-expert sampled and evaluated roughly 6,600 cases using computations (MD simulations). Top 100 pentapeptides from this were further modeled for longer time-

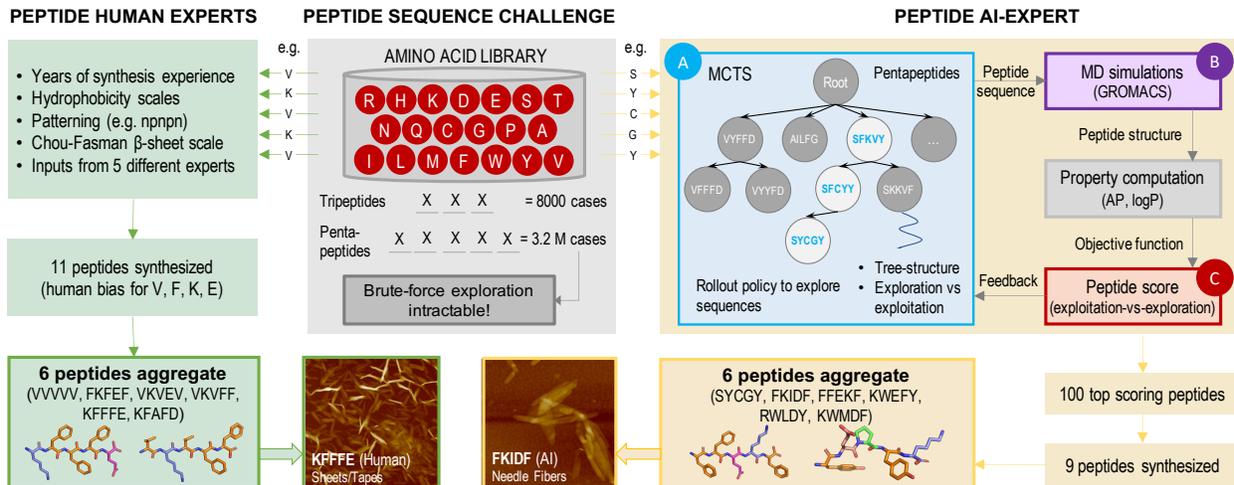


Figure 1: **Workflow adopted to discover self-assembling pentapeptides using inputs from the human experts and the developed AI-expert.** The search space of peptides grows drastically with its sequence length owing to the presence of 20 amino acids. While 8000 possible tripeptides can be explored (computationally) for assembly using a brute-force approach, the space of 3.2 million pentapeptides is intractable. The human experts use rational design approaches, such as hydrophobicity scales, charge balance, patterning (nnpn; n:non-polar, p:polar), and their own individual experiences to design self-assembling pentapeptides. 6 out of the 11 synthesized pentapeptides suggested by the six different human experts were found to aggregate, although the proposed sequences suggested human bias toward V, F, K and E amino acids. In contrast, the developed AI-expert combines (A) Monte Carlo tree search (MCTS), (B) MD simulations and (C) a peptide structure based scoring function to efficiently search for self-assembling peptides. 6 out of 9 synthesized peptides from the AI-expert were found to aggregate. Beyond being able to recover some intuitive sequences (FFEKF), the AI-expert suggested some novel/unusual sequences (SYCGY) involving diverse amino acids (RWLDY) that reflect its advantage to overcome the human bias. Molecular representation and atomic force microscopy images for a few promising pentapeptides from both categories are also shown.

scales (200 ns) using more rigorous MD simulation parameters to improve the AP estimates. 9 top-scoring AI sequences were screened for actual synthesis and experimental investigation, 6 of which were found to aggregate based on light scattering and atomic force microscope measurements. In comparison, 6 out of 11 pentapeptide sequences suggested by the human experts, *i.e.*, peptide synthesis experts with many years of experience, were found to aggregate. We discuss these findings in the context of prevalent bias (and thus similarity) in the sequences proposed by the human experts, the ability of the AI-expert to recover existing materials knowledge by reproducing sequences similar to the human experts, and, most importantly, the power of AI-expert to overcome human bias by discovering previously unknown and completely non-intuitive self-assembling peptide sequences (e.g. SYCGY) in an efficient manner. We provide our perspectives and propose a path forward where the performance of the AI-expert can be enhanced by fusing information from parallel computations and experiments, and by improving the MCTS scoring function to include other structural factors from the existing protein databases.

Results and Discussion

AI-expert for peptide discovery

We develop an AI workflow, henceforth referred to as the AI-expert, to discover self-assembling peptides. The workflow consists of a search algorithm, *i.e.*, MCTS interfaced with GROMACS MD simulation software [31, 32] to model structure-functionality of a peptide, given only its amino acids sequence. The role of MCTS is to intelligently and efficiently generate peptide sequences sampled from the overall search space that can

show high self-assembling scores. The MD simulations *via* GROMACS provide a relatively inexpensive method to estimate AP of the peptide sequence proposed by the MCTS algorithm, and thus provide feedback to improve the quality of the peptide search. It should be noted that the AI-expert autonomously switches between the two stages of peptide generation (MCTS) and evaluation (MD simulations) without any human intervention.

MCTS is a powerful algorithm for planning, optimization and learning tasks owing to its generality, low computational requirements, and a theoretical bound on exploration-vs-exploitation trade-off [33, 34, 29]. It has been particularly successful in problems involving extremely large search space [35, 36, 37], making it the model of choice for this work. While its details are covered in Methods section, we briefly note that it searches in a tree-structured fashion where every node (or tree leaf) contains a unique peptide sequence (e.g. VKVKV) and its associated score; see Figure 1. Moreover, these nodes contain connections in a special manner such that a parent node is connected to several child nodes with slightly different peptide sequences. This gives a meaningful structure to the overall tree, with the high-scoring child node generally belonging to the tree branch that contains other relatively high-scoring parent nodes. To advance the search, MCTS utilizes a tree policy and a rollout policy. The former selects the most promising node, while the latter samples the nearby space (using Monte Carlo trials) of the selected node by introducing small perturbations, referred to as rollouts. The upper confidence bound for parameters (UCP) [38] is a popular choice of tree policy given by:

$$\text{UCP}(\theta_j) = -\min(r_1, r_2, \dots, r_{n_i}) + c \cdot f(\theta_j) \cdot \sqrt{\frac{\ln N_i}{n_i}} \quad (1)$$

where θ_j represents the node j in the MCTS structure, r denotes the score (or reward) of a given rollout, $c(> 0)$ is the exploration constant, n_i is the number of rollout samples taken by node θ_j and all of its child nodes, and N_i is the same value as n_i except for the parent node of θ_j . In this work the scoring (or reward) function was chosen to balance the AP and the hydrophobicity of the peptides using the form, $r_i = \text{AP}'^\alpha * \log \text{P}'^\beta$, where symbol ' denotes the normalized values, and α and β are the coefficients weights. $f(\theta_j)$ is the uniqueness function specifically introduced in this work to drive search towards diverse sequences. The policy in Eq. 1 tries to balance the search between those nodes which have either returned the maximum score (left term) or have not been explored enough (right term). In contrast, the rollout policy introduces random, but controlled, perturbations (from a node) to sample new sequences.

We introduce two modifications to improve the efficiency of MCTS: the uniqueness function $f(\theta_j)$ and a random forest (RF) model guided rollout policy. The uniqueness function enhances the effect of its exploration term in Eq. 1, further motivating MCTS to select those nodes that represent diverse peptides. This pushes the search in new regions (or diverse sequences) that have not been explored before. For this work, we used the Morgan circular fingerprints [39] to numerically represent the peptides followed by the Dice similarity measure to compute the uniqueness of a peptide in relation to others in the MCTS structure (see Methods). The second important modification is the use of a surrogate RF model to quickly predict AP of a peptide given its sequence. This eliminates the need to perform computationally expensive MD simulations during rollouts, especially for cases that are predicted to have very low AP values. However, care should be taken to only partly replace the MD simulations with the RF model as the surrogate model is only approximate and could miss-out on promising cases that are different from the data used for its training. Thus, here, we use the RF model to only guide the rollout policy such that half of the rollouts correspond to cases that are predicted to have high AP, while the remaining half is from random perturbation as in the traditional MCTS setting (see Methods). For both cases, the AP value used in Eq. 1 is obtained only after actual MD simulations. It should be noted that the RF model is trained in an online fashion, with the RF model being regularly updated as more training data from the MD simulations becomes available during the MCTS run. Details on the input features and training parameters of the RF model are provided in Methods section.

Validation for tripeptides

We first consider the space of tripeptides as a demonstration of the ability of the AI-expert to accelerate the search of self-assembling peptides. Two reasons that dictate this choice are: first, tripeptides exhibit a computationally manageable space of 8000 ($= 20^3$) sequences and second, there already exists a rigorous past

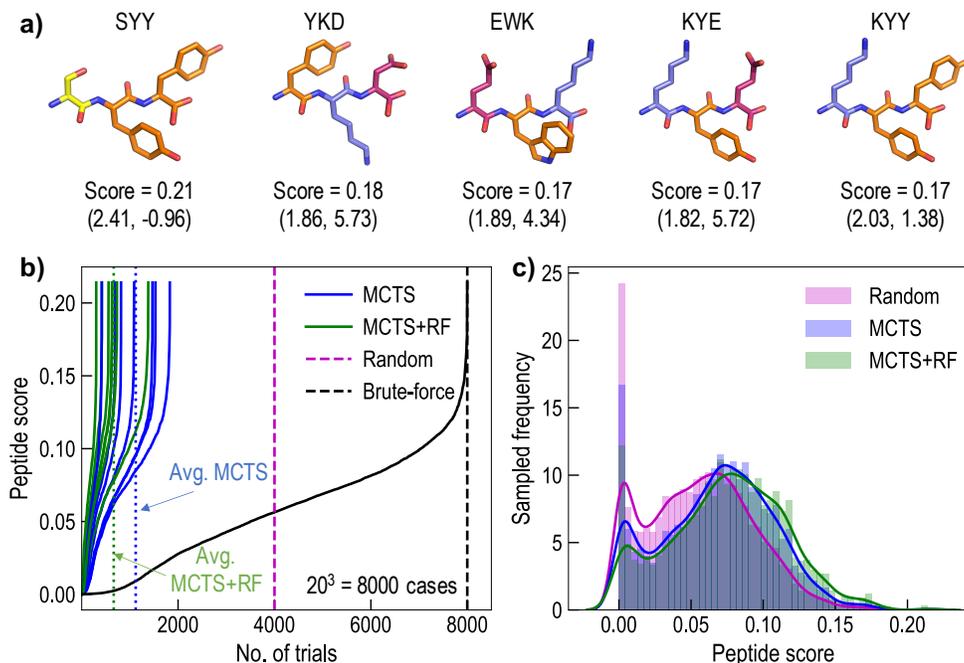


Figure 2: **Performance comparison of the different search strategies on the space of tripeptides.**

(a) Molecular representation of example top scoring tripeptides along with their computed scores. The numbers within the bracket indicate (left) the aggregation propensity (AP) and (right) hydrophobicity (logP) values. (b) Comparison of the number of trials needed to search the highest scoring tripeptide from the complete space of 8000 cases. AI-expert utilizing MCTS or MCTS+RF search strategies on an average take substantially low number trials in comparison to a random or an exhaustive search to find the highest scoring tripeptide (SYY). (c) Comparison of the score of peptide sequences generated using a random, MCTS or MCTS+RF search strategy; solid lines denote the respective normalized density. The AI-expert with MCTS+RF scheme is most efficient in identifying high-scoring peptides as a larger fraction of its generated peptide population have high scores.

study on the use of MD simulations to explore self-assembling tripeptides using a brute-force approach [24]. In fact, the previous works on di- and tri-peptides confirms that the MD simulations based on MARTINI coarse-grained force field [40, 41, 42] are reliable enough to guess self-assembly, when coupled with other metrics of AP and hydrophobicity [30, 24].

To measure the performance gain of the AI-expert over a brute-force approach, we first performed MD simulations for all 8000 cases and sorted them based on their score, $r^{\text{tri}} = \text{AP}^2 \times \log P'$ ($\alpha = 2$, $\beta = 1$, same as that in the past work); see Methods section for details on MD simulations, and computation methodology of AP and logP. Some of the top scoring tripeptides are shown in Figure 2a and Extended Data Table 1. While the specific rank-order of our results may slightly differ to the previous brute force investigation [24], the overall AP vs hydrophobicity trends in both the studies match well. The differences can be traced to the slight variations in the AP computations introduced due to either the stochasticity of the MD simulations or the software choice of the AP computations (see Supporting Information). Nonetheless, the trends in the identified top scoring tripeptides are similar. For example, YKD, EWK, and KYE all follow the general trend of charge balanced peptides (K/D or K/E) plus an aromatic residue, while KYE and SYE are amphiphilic peptides displaying a pair of aromatic residues. We also note that the specific rank-order of the tripeptides does not influence the conclusions made in this study regarding the search efficiency of the AI-expert, as discussed next.

Figure 2b compares the time taken by the different methods to identify the highest scoring tripeptide, i.e., SYE. It can be seen that the AI-expert which uses RF boosted MCTS (labeled MCTS+RF) on an average takes substantially lower number of trials to identify the highest scoring SYE sequence, as compared

to a purely random rollout policy based MCTS. This suggests that the developed rollout policy utilizing RF model indeed helps the AI-expert to efficiently identify high-scoring peptides. Furthermore, the AI-expert, with or without the RF model, performs significantly better than a random or a brute-force search requiring ~ 4000 and 8000 trials, respectively. Similarly, Figure 2c compares the quality of the peptide population generated using a random search, or by the AI-expert utilizing MCTS or MCTS+RF scheme. It is evident that the MCTS+RF scheme samples high-scoring tripeptides most frequently, followed by the MCTS and then the random search. Overall, these results validate that the AI-expert can efficiently identify high-scoring peptides without resorting to time-intensive brute-force search.

Screening of pentapeptides

Having validated the efficiency of the AI-expert for tripeptides, next we use it to discover self-assembling pentapeptides with 3.2 M (20^5) permutations. Such a large search space renders a brute-force search impossible and motivates the need for an AI-guided search. The AI-expert with MCTS+RF scheme was deployed with a slightly different settings of the reward function, i.e., $r^{\text{penta}} = \text{AP}^2 \times \log \text{P}^{0.5}$ ($\alpha = 2$, $\beta = 0.5$), to bias the search towards pentapeptides that are neither too hydrophilic (easily soluble) nor too hydrophobic (difficult to form hydrogels). This adjustment in the reward function is necessary because a majority of the amino acids are hydrophilic, and the naive use of r^{tri} for pentapeptide design will incorrectly assign high scores to hydrophilic pentapeptides (discussed later). Results for the ~ 6600 pentapeptides evaluated using the MCTS+RF (with r^{penta}) are shown in Figure 3a. It can be seen that AI-expert found high occurrence of moderately hydrophobic pentapeptides with $\log \text{P}$ between 0 and -4, although with a broad peak. A list of top 100 peptides from this, based on their reward score and an additional constraint of $-0.6 < \log \text{P} < 2$, were screened for longer MD simulations (200 ns) and the AP estimates were improved (see Supporting Information for the complete list). From a computational viewpoint, significant aggregation was observed in all of the selected 100 cases, with a few example pentapeptides structures shown in the top row of Figure 3b.

Besides the AI-expert, several human-experts were asked to suggest their own sequence of pentapeptides which they expect to assemble. A set of simple guidelines were supplied (see Methods), and in response, a total of 29 pentapeptides were collected. Many literature examples of self-assembling pentapeptides include N- and C- terminii modification (acetylated or carbamidated, respectively) to facilitate assembly [43, 44, 45, 46, 47, 48]. However, in this work the human-experts were directed to leave the pentapeptide terminii unmodified in alignment with the workflow adopted for the AI-expert. Analogous to the AI-expert pentapeptides, AP and $\log \text{P}$ values for these sequences were evaluated using MD simulations and hydrophobicity scales, respectively.

Results for the top 100 pentapeptides from the AI-expert (purple markers) and the 29 sequences from the human-experts (green markers) are shown in Figure 3b. Also, captured are the results of the candidates that were synthesized (black markers) and those that were found to aggregate (solid markers) based on light scattering and microscopy measurements. A detailed comparison of the pentapeptides suggested by the AI-expert and the human experts is discussed later. Here, however, we make the following observations: first, the top sequences screened by the AI-expert lie in a relatively smaller $\log \text{P}$ range as compared to those proposed by the human experts. This is because the AI-expert screen candidates only on the basis of the scoring function, while the humans rely on a multitude of factors, such as patterning, hydrophobicity scales and individual past experiences. Second, the AI-expert suggested sequences, in general, show higher AP values as compared to those of the human experts. This implies that at least from a computational modeling viewpoint, the AI-expert indeed found pentapeptide sequences with higher degree of aggregation. This is also evident from the example pentapeptide structures obtained after longer MD simulations (200 ns) comparing the AI-expert (top row) and the human experts (bottom row) sequences in Figure 3c. Third, many sequences that were computationally found to have high AP values did not display any assembly upon experimental synthesis. These cases highlight the limitations of the MARTINI force field to capture accurate aggregation behavior in peptides, or the inadequacy/simplicity of the reward function used in this work which consists of just the AP and $\log \text{P}$ values. Lastly, pentapeptides that were (experimentally) observed to aggregate belonged to a narrow range of $\log \text{P}$ (3 to -5) and AP (1.5 to 2.5), signaling the importance of these theoretically derived values in identifying novel peptide sequences for self-assembly. The observed narrow range of $\log \text{P}$ agrees well with the convention of balancing the hydrophobic and hydrophilic

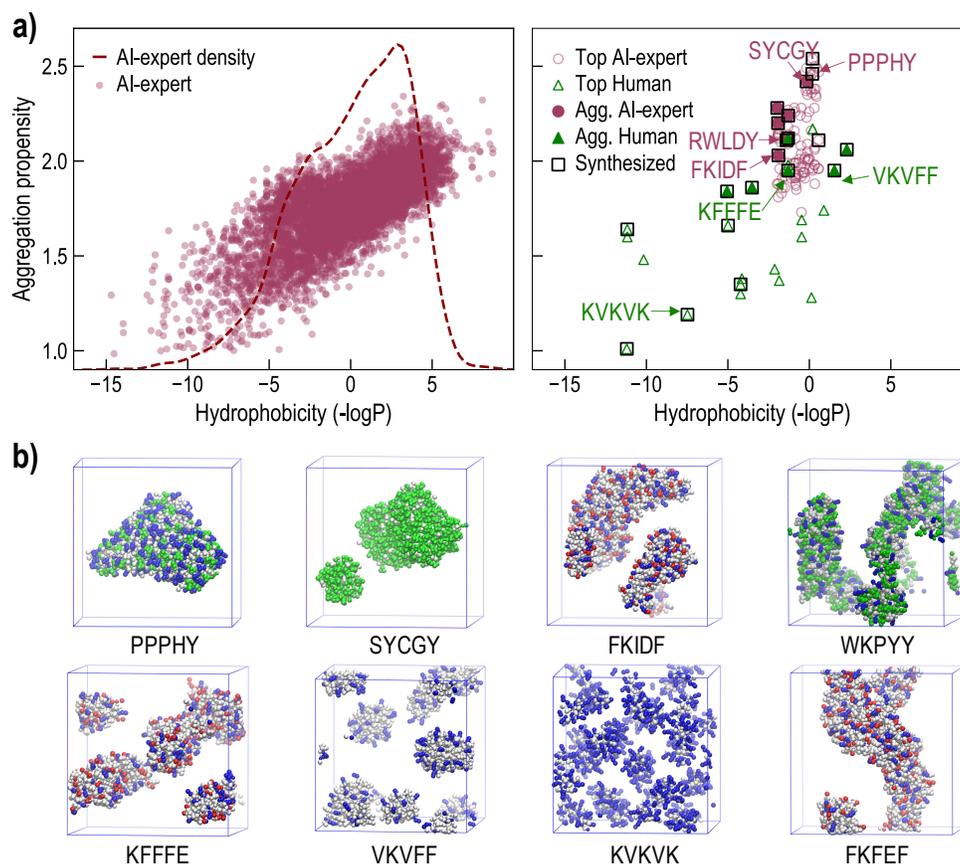


Figure 3: **Screening of pentapeptides from the AI-expert and the human experts.** (a) (left) Results of the MCTS+RF based computational search of the AI-expert using the scoring function r^{Penta} . A broad peak with $\log P$ between 0 and -4 indicate generation of moderately hydrophobic peptides that display a good balance between aggregation propensity (AP) and hydrophobicity. The AP results are based on shorter MD simulations (50 ns). (Right) Top peptides that were screened by the AI-expert (top 100 using r^{Penta}), suggested by the human experts, and those that were selected for synthesis. The AP results are based on longer MD simulations (200 ns). (b) MD simulation results (200 ns) for example top scoring pentapeptides, showing different levels of aggregation.

content in peptide sequences; if logP values are too negative (hydrophobic), peptides begin to precipitate out of solution or are rendered entirely insoluble in water even at low concentrations, and if logP values are too positive (hydrophilic), then they remain as water soluble monomers. Similarly, computed AP value is also a good indicator of peptide assembly, as no peptide sequence with a low AP value was observed to assemble.

Discovery of self-assembling pentapeptides

This section covers details on the 20 synthesized pentapeptides and the observed self-assembled structures. 11 of the 29 sequences from the human-experts and 9 of the 100 AI-expert suggested sequences were prepared using solid-phase peptide synthesizer (see Methods) with the termini of the peptides kept unmodified, i.e. amine, and carboxyl groups of the final product were unprotected. Since, we are interested in the ability of peptides to aggregate and/or assemble, it is important here to distinguish between two seemingly analogous terms: aggregation and assembly. Aggregation implies the lack of noticeable structure and assembly implies the presence of nano-, meso-, and microscale features like micelles, vesicles, fibers, and sheets. Thus, for a detailed analysis of the aggregated/assembled structure, and to find experimental quantities analogous to the computed logP and AP values, liquid chromatography, mass spectrometry, infrared spectrometry, atomic force microscopy (AFM) and opacity measurements were made for each of the synthesized pentapeptide.

First, all the synthesized peptides were analyzed for purity and mass by high performance liquid chromatography and mass spectrometry. The retention time (RT) was recorded and was found to correlate with hydrophobicity (logP); a linear relationship between RT and logP is visible in Figure 4a, although some deviations are also noted. Importantly, most of the peptides that show aggregation (solid circles) displayed high RT.

To determine assembly, all 20 peptides were dissolved in water at 2 wt% and the pH was adjusted to 7. After 24 hrs, the solutions either remained clear, grew cloudy, or gelled upon adjusting the pH (see Figure 4d). The samples opacity (absorption at 800 nm) was monitored with a plate reader (OD800nm) and peptides with OD800nm >0.1 (for water, OD800nm =0.04) were considered to aggregate (see Extended Data Table 2). Overall, 6 candidates, i.e., VVVVV, FKFEF, VKVEV, VKVFF, KFFFE and KFAFD, from the humans, and 6 candidates, i.e., SYCGY, FKIDF, FFEKF, KWEFY, RWLDY and KWMDF, from the AI-expert were found to aggregate. One peptide from the human experts, RVSVD, yielded high opacity values after 1 week and was not considered to be a “positive hit”. Only a modest match between the measured OD800nm and the computed AP was observed as shown in Figure 4a. Peptides that yielded OD800nm >0.1 had an AP value >1.8, although some peptides with low opacity (OD800nm < 0.05) also had an AP >1.8 indicating that the computed AP is not always a good predictor for aggregation. In this regard, we also caution that the opacity measurements do not always indicate assembly. For example, micelles at the nanoscale remain translucent and would not yield high values at OD800nm.

Thus, we further analyzed the secondary structures of the aggregated peptides using FTIR spectroscopy. In the amide I region of the spectrum, peptide/protein in D₂O show a signature FTIR vibrations of 1675 cm⁻¹ and 1627 cm⁻¹ that are representative of a β -sheet conformation [49]. While we observed (Figure 4c) the former peak in almost all the samples investigated (2 wt% in D₂O, pH 7), the latter peak, which is attributed to the β -sheet composition, was observed in the following peptides: FKIDF, KFFFE, FKFEF, VVVVV and VKVFF. In addition to the samples in solution, dried films cast from diluted stock solutions (10 μ L of a 0.2 wt% solution onto a CaF₂ window) yielded the β -sheet signature of amide I vibration at 1627 cm⁻¹ in more peptides (see Supporting Information). 9/11 peptides from the users indicated β -sheet formation as opposed to 3/9 AI peptides. This reflects the bias of the users towards peptides with β -sheet formation, as will be discussed later.

To investigate the morphology of the designed pentapeptides, atomic force microscopy (AFM) was employed on the dried sample films (see Methods). We note that the dried films could contain structural artifacts, but in most cases the microscopy data correlate well with our solution studies—only three cases are believed to form aggregates owing to the drying, KVKVK, RVSVD and VKVKV. Among the cases recommended by the AI-expert, SYCGY yielded microscale length needles with widths on the order of 100s of nm (see Figure 4b). Nanoscale structures were found in the peptides rich in aromatic amino acids, tryptophan and phenylalanine: KWEFY (fibers), FFEKF (spheres), and FKIDF (platelets). This is a similar design feature in the past study on tripeptide series explored by the Ulijn and Tuttle groups [24]. Two sequences that introduce aliphatic amino acids (Leu and Met) in the middle of the sequences, RWLDY and

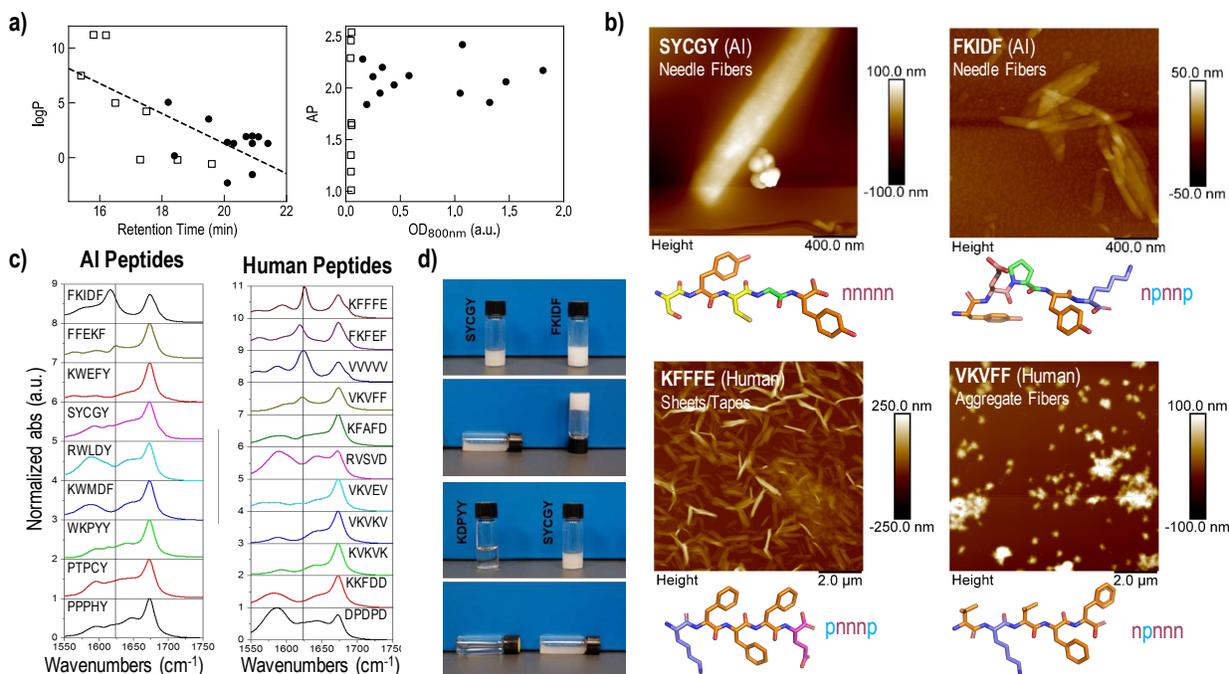


Figure 4: **Experimental measurements of self-assembly in pentapeptides.** (a) Retention time (RT) and opacity (OD800nm) measurements for the 20 synthesized pentapeptides as suggested by the AI and the human experts. Peptides that were found to aggregated are marked with solid circles. While RT was found to linearly correlate with logP, OD800nm was analogous to the computed aggregation propensity (AP) values. (b) Atomic force microscopy images for example pentapeptides synthesized in this work, along with their molecular representation and pattern of mixed polar (p) and non-polar (n) amino acids. Aggregates in form of fibres, sheets/tapes, and other irregular shapes are visible. (c) Infrared spectroscopy measurements for the 20 synthesized pentapeptides as suggested by the AI and the human experts. The peak near 1600 cm^{-1} highlights the formation of secondary structures (β -sheets) in many of the human-expert selected systems, which was largely missing in the AI recommended systems. (d) Photographs of example gel, solution and suspensions formed by different peptides.

KWMDF, yield spherical structures with varying degrees of aggregation. Interestingly, WKPYI indicated no assembly via our experimental protocol, but large aggregates were observed in the solution as well as the AFM. Thus, among the selected 9 cases from this list, only PPPHY and PTPCY did not indicate any discernible structure and resembled dried organic matter on a substrate.

Human experts designed phenylalanine rich peptides, similar to those identified by the AI-expert, demonstrated a high propensity for forming nanostructures. Nanoscale fibers were discovered for KFAFD, FKFEF, VKVFF (spherical bundles) and KFFFE (nanoplates). Many of the valine rich peptides from the human experts form fibrous structures upon drying (VKVKV, KVKVK and RVSVD) while no evidence for solution structures were observed. Interestingly, large platelets with 25 nm height were observed for relatively hydrophobic ($\log P = -2.3$) VVVVV. Even though it is highly hydrophilic ($\log P = 5.05$), VKVEV gelled upon increasing the pH to 7. Thus, 6 of the 11 pentapeptides suggested by human experts formed nanostructures, most of which formed β -sheet conformation.

Performance comparison of AI-expert and human experts

In terms of the overall ability of the AI-expert to predict assembly in pentapeptides, it performs at par or slightly better than the human experts. As shown in Figure 5a, the success rate of the AI-expert (using r^{penta}) is 66.67 % as compared to 54.5 % in the case of human experts. We, however, argue that aggregation success rate alone is not a sufficient metric to evaluate performance. It is important to realize that, in

contrast to the human experts, the AI-expert had no direct feedback from any actual experiments and merely relied on computationally derived quantities, such as AP and logP, to make predictions. So, if all the synthesized peptides are rank-ordered on the basis of their computational score (i.e., r^{penta}), 7 of the top 8 peptides are from the AI-expert and only 1 are from the human experts (see Figure 5b). This means that in the ideal case where the computational scoring function is a perfect indicator of assembly, the AI-expert would have performed much superior to the human experts. Thus, efforts are needed to either modify the scoring function (e.g., addition of other structural factors or manipulating weighting scheme), or improve the performance of the force fields to accurately relate the computational score with peptide assembly.

a) Aggregation performance

Expert	Agg.	Total	% Success
Human	6	11	54.5
AI (r^{penta})	6	9	66.7

b) Top scoring pentapeptides: Computational score

Rank	Peptide	Expert	r^{penta}
1	PTPCY	AI	0.49
2	PPPHY	AI	0.44
3	SYCGY	AI	0.43
4	KWMDF	AI	0.38
5	WKPY	AI	0.34
6	KWEFY	AI	0.33
7	FFEKF	AI	0.31
8	FKFEF	Human	0.28

Top scoring pentapeptides: Experimental score

Rank	Peptide	Expert	ExpScore
1	SYCGY	AI	0.30
2	KFAFD	Human	0.26
3	FFEFK	AI	0.24
4	VVVV	Human	0.21
5	KFFFE	Human	0.09
6	VKVEV	Human	0.05
7	RWLDY	AI	0.03
8	KWEFY	AI	0.03

c)

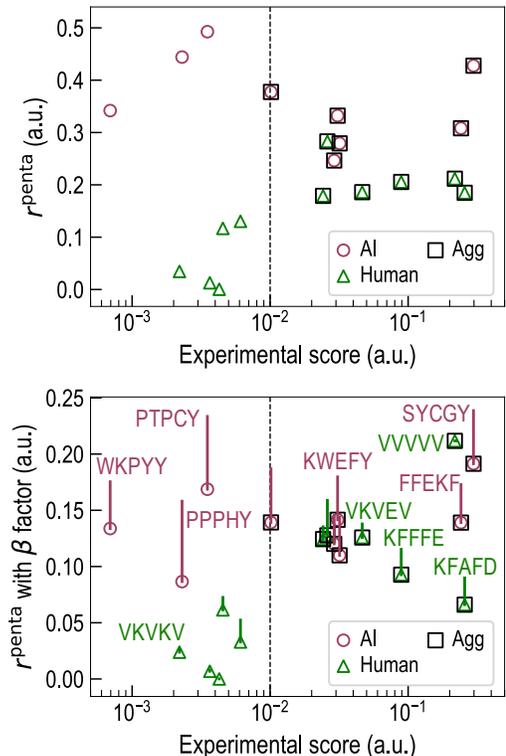


Figure 5: **Performance comparison of the AI-expert against human experts.** The performance of the AI-expert and the human experts are evaluated in terms of the (a) aggregation success rate, and (b) computational and experimental scores of the proposed peptides. (c) Co-relation between the computational and the experimental score, with (bottom) and without (top) the β -sheet factor. Although a high computational score not necessarily indicate aggregation, the experimental score beyond a threshold of 0.01 (dotted lines) captures the peptide aggregation extremely well.

As another form of performance evaluation, we devised an experimental scoring (ExpScore) metric that incorporates information from the characterization measurements. ExpScore was defined as the product of the normalized RT (analogous to logP) and the normalized OD800nm (analogous to AP) and captured the aggregation in peptides extremely well; see Methods. As seen in Figure 5b, even on the basis of ExpScore, the AI-expert performs at par with the human experts with each of the two suggesting 4 of the top 8 scoring peptides. Nevertheless, the the highest scoring case (SYCGY) was predicted by the AI-expert, further corroborating the ability of the AI-expert for peptide discovery.

The diversity of the proposed peptide sequence is another important performance metric. Much of the sequences proposed by the human experts revolved around the use of only four amino acids, i.e., phenylalanine (F), valine (V), lysine (K) and glutamic acid (E). This is not only reflective of the human bias, but in some

sense also the reason for good success rate from the human experts (for instance, many similar pairs, VKVFF, VKVEV, and KFFFE, FKFEF, are all counted as positive hits). The AI-expert, on the other hand, suggested very diverse sequences covering more than 10 distinct amino acids. Furthermore, it suggested sequences which are very unusual, such as SYCGY, or include many distinct amino acids, e.g., FKIDF, RWLDY and KWEFY. None of these sequences is likely to be recommended by a human expert and truly suggests the power of the AI-expert to overcome human bias, identify novel sequences and plausibly unearth new protein chemistry.

Besides discovering novel chemistry, the AI-expert automatically (re-)discovered a few known rational design approaches. For instance, a significant percentage of the peptides ($\sim 50\%$) determined by the AI-expert, were either charge neutral or charge balanced such that a positively charged amino acid like lysine or arginine was frequently paired with a negatively charged amino acid like glutamic or aspartic acid. Such inclusion of salt-bridges/electrostatic pairing in peptide design is a standard practice. Another unique area where both the human-experts and the AI-expert agreed on was the incorporation of phenylalanine rich peptides balanced with the charged pair of lysine and glutamic/aspartic acid, e.g. KFFFE, FKFEF (humans), and FKIDF, FFEKF (AI). Phenylalanine is well known to form β -sheet rich structures that further stabilize supramolecular assemblies via $\pi - \pi$ interactions. The generation of such sequences via the AI-expert is rather encouraging.

However, there are many deviations between the human experts and the AI-expert. Using r^{penta} , some unusual high scoring sequences were observed. These sequences lacked electrostatic pairs and incorporated uncharged/polar amino acids, e.g. SYCGY and PPPHY. Incorporation of cysteine can be challenging due to its ability to cross-link and form cysteine bridges. For self-assembly, this can be beneficial, but the simulations do not predict disulfide bond formation rendering the experimental and simulation results difficult to compare. Another deviation between the AI-expert and the rational design approach is with respect to the amino acids valine and proline. Very few peptides suggested by the AI-expert contained valine, yet human experts show a proclivity for it (e.g. KVKVK, RVSVD, VKVEK). Generally, valine is employed due to its high β -sheet propensity which often leads to long-range ordering in self-assembly [50]. The trend was opposite in the case of proline; while several of the AI-expert sequences were dominated by this amino acid, none of the human experts believed it to help with any sort of assembly. The added observation that none of the proline containing pentapeptide aggregated, points to a limitation of the AI-expert.

AI-expert improvement opportunities

A critical component behind the success of the AI-expert is its scoring function, which needs to be designed very carefully. To show how dramatically it can influence the performance of the AI-expert, we again use the AI-expert to design pentapeptides, but this time with the reward as r^{tri} . This mostly yielded highly hydrophilic candidates with a logP value between 2 and 6, as shown in Supporting Information, Figure S2. Although scored high (based on r^{tri}), these candidates are expected to be soluble in water and not show any assembly. Analogous screening procedure was followed, wherein 10 candidates were synthesized from the list of top 100 cases based on longer MD simulations (200 ns). Only 2 of 10 cases, i.e., KFFFDY and FFEKF, yielded aggregates with a success rate of only 20% (see Supporting Information). Thus, the selection of an instructive scoring function is quintessential for the AI-expert, and the weighting parameters α and β needs to be carefully adjusted according to the sequence length, n , of the peptide.

We noted previously that the AI-expert did not receive any feedback from actual experiments and only a modest correlation between the computational (r^{penta}) and the experimental score was found, as shown in Figure 5c. This observation, along with the general tendency of the human experts to incorporate amino acids with high β -sheet propensity (e.g. valine and phenylalanine) and unfitting proclivity of the AI-expert to proline, provide an opportunity to improve the performance of the AI-expert.

Chou and Fasman have reported β -sheet propensity as a statistical distribution of amino acid conformers in the protein data bank [23] and the approach has been periodically updated [50]. Being a quantitative measure, this can be used to modify the scoring function, i.e., $r^{\text{score}} = AP'^{\alpha} \times \log P'^{\beta} \times \log \sigma'^{\gamma}$, where σ is the reported β -sheet propensity factor of weight γ . Figure 5c compares the correlation between the computational and the experimental score with ($\gamma=1$) and without ($\gamma=0$) the β -sheet propensity factor. Further, the vertical lines in the bottom panel of Figure 5c are proportional to change in the peptide score upon addition of the β -sheet propensity factor and denote its effect on the peptide score. It can be seen

that with the β -sheet, the computational score for the proline containing sequences, WKPY, PTPCY and PPPHY, which did not assemble, decreased substantially. Similarly, the human experts suggested peptides that had low computational and experimental scores, and did not aggregate (VKVKV), continued to display low scores even after the inclusion of the β -sheet propensity factor. Further, many of the peptides that were found to aggregate either did not show any change (VVVVV, VKVEV) or their scores decreased marginally (KFFFE, SYCGY, FFEKF). This overall results in an improved ranking of the peptides, and we believe, this improved scoring system (r^{Score}) could be used for peptide discovery in the future. However, caution should be exercised in the selection of the γ factor, as a high γ value will dominate the search towards selected amino acids (e.g., V, I and F) or to peptides with only β -sheet conformations, and thus introduce unwanted bias in the search and minimize amino acid diversity.

Our future vision is the development of a fully autonomous peptide design platform where the AI-expert interacts with a robotic platform capable of synthesizing and characterizing new sequences, whose feedback is directly digested by the AI-expert to suggest new sequences and the search is progressed in an iterative manner. To accelerate this process, inputs from simulations could also be utilized to avoid low scoring peptides.

Conclusions

AI methodologies are incredibly useful for guiding scientists towards identifying novel short self-assembling peptides and are being considered as the future of synthesis and molecular design. AI facilitated peptide discovery is necessary because of the intractable search space (20^n , where n is the peptide length). Here, we developed an AI-expert to evaluate the aggregation propensity of 6,600 out of 3.2 million possible pentapeptides using MD simulations and hydrophobicity ($\log P$) scales. In addition, we queried expert peptide designers to provide promising sequences. Top 9 sequences from the AI-expert and 11 candidates from the human experts were synthesized and characterized. An experimental scoring system (sample opacity vs HPLC retention time) that reflected the AI scoring system (aggregation propensity vs hydrophobicity) was critical in identifying failures and successes in the approaches of both the AI-expert and the human experts. Overall, the AI-expert performed at par or slightly better with a 66.7% success rate to human experts' 54.5%. Not only did the AI-expert recovered known design strategies, such as identification of charge balanced phenylalanine rich peptides (AI - FKIDF and FFEKF, human - KFFFE and FKFEF), it also found novel sequences that deviate significantly from the traditional approach (e.g. SYCGY).

Human bias was demonstrated to favor pentapeptides with high β -sheet propensity scores and was used as an opportunity to improve the AI scoring metric. Including the β -sheet factor to the AI score shifted the rankings in the correct direction, but still did not fully resemble the experimental ranking. Future efforts will focus on the application of high throughput peptide synthesis coupled to the developed experimental scoring system to provide an experimental feedback loop to the AI-expert beyond the currently implement theoretical metrics (AP and $\log P$). Similar AI strategy could be extended to screen small libraries of peptides for more specific applications. Although this study demonstrates the success of the AI-expert in discovering self-assembling peptides, it can be extended to discover functional peptide assemblies for applications involving light harvesting, catalysis, mechanical stability and conductivity.

Acknowledgements

This work was performed, in part, at the Center for Nanoscale Materials, a U.S. Department of Energy Office of Science User Facility, and supported by the U.S. Department of Energy, Office of Science, under Contract No. DE-AC02-06CH11357. Supported by the University of Chicago and the Department of Energy under Department of Energy Contract No. DE-AC02-06CH11357 awarded to UChicago Argonne, LLC, operator of Argonne National Laboratory. We gratefully acknowledge the computing resources provided on Bebop, high performance computing clusters operated by the Laboratory Computing Resource Center (LCRC) at Argonne National Laboratory. SKRSS acknowledges the support from the UIC faculty start-up fund. Authors acknowledge Dr. Tell Tuttle for sharing their computational data on tripeptides.

Code Availability

The codes, scripts and the AI-expert framework are available from the authors upon reasonable request.

Data Availability

The data that support the findings of this study are available from the authors upon reasonable request.

References

- [1] Shichen Zhu, Qijuan Yuan, Tao Yin, Juan You, Zhipeng Gu, Shanbai Xiong, and Yang Hu. Self-assembly of collagen-based biomaterials: preparation, characterizations and biomedical applications. *J. Mater. Chem. B*, 6(18):2650–2676, 2018.
- [2] Anna Sorushanova, Luis M Delgado, Zhuning Wu, Naledi Shologu, Aniket Kshirsagar, Rufus Raghunath, Anne M Mullen, Yves Bayon, Abhay Pandit, Michael Raghunath, et al. The collagen suprafamily: from biosynthesis to advanced biomaterial development. *Adv. Mater.*, 31(1):1801651, 2019.
- [3] Randolph V Lewis. Spider silk: ancient ideas for new biomaterials. *Chem. Rev.*, 106(9):3762–3774, 2006.
- [4] Gregory D Scholes, Graham R Fleming, Alexandra Olaya-Castro, and Rienk Van Grondelle. Lessons from nature about solar light harvesting. *Nat. Chem.*, 3(10):763–774, 2011.
- [5] Quan Luo, Chunxi Hou, Yushi Bai, Ruibing Wang, and Junqiu Liu. Protein assembly: versatile approaches to construct highly ordered nanostructures. *Chem. Rev.*, 116(22):13571–13632, 2016.
- [6] Gang Wei, Zhiqiang Su, Nicholas P Reynolds, Paolo Arosio, Ian W Hamley, Ehud Gazit, and Raffaele Mezzenga. Self-assembling peptide and protein amyloids: from structure to tailored function in nanotechnology. *Chem. Soc. Rev.*, 46(15):4661–4708, 2017.
- [7] Rein V Ulijn and Andrew M Smith. Designing peptide based nanomaterials. *Chem. Soc. Rev.*, 37(4):664–675, 2008.
- [8] Anupama Lakshmanan, Daniel W Cheong, Angelo Accardo, Enzo Di Fabrizio, Christian Riekkel, and Charlotte AE Hauser. Aliphatic peptides show similar self-assembly to amyloid core sequences, challenging the importance of aromatic interactions in amyloidosis. *Proc. Natl. Acad. Sci.*, 110(2):519–524, 2013.
- [9] Sayanti Brahmachari, Zohar A Arnon, Anat Frydman-Marom, Ehud Gazit, and Lihi Adler-Abramovich. Diphenylalanine as a reductionist model for the mechanistic characterization of β -amyloid modulators. *ACS Nano*, 11(6):5960–5969, 2017.
- [10] Miri Yemini, Meital Reches, Judith Rishpon, and Ehud Gazit. Novel electrochemical biosensing platform using self-assembled peptide nanotubes. *Nano Lett.*, 5(1):183–186, 2005.
- [11] Tayebeh Zohrabi, Neda Habibi, Ali Zarrabi, Maryam Fanaei, and Lai Yeng Lee. Diphenylalanine peptide nanotubes self-assembled on functionalized metal surfaces for potential application in drug-eluting stent. *J. Bio. Mater. Res. A*, 104(9):2280–2290, 2016.
- [12] Xuehai Yan, Pengli Zhu, and Junbai Li. Self-assembly and application of diphenylalanine-based nanostructures. *Chem. Soc. Rev.*, 39(6):1877–1890, 2010.
- [13] Andrei Kholkin, Nadav Amdursky, Igor Bdikin, Ehud Gazit, and Gil Rosenman. Strong piezoelectricity in bioinspired peptide nanotubes. *ACS Nano*, 4(2):610–614, 2010.
- [14] Xuehai Yan, Qiang He, Kewei Wang, Li Duan, Yue Cui, and Junbai Li. Transition of cationic dipeptide nanotubes into vesicles and oligonucleotide delivery. *Angew. Chem., Int. Ed.*, 119(14):2483–2486, 2007.

- [15] Xiubo Zhao, Fang Pan, Hai Xu, Mohammed Yaseen, Honghong Shan, Charlotte AE Hauser, Shuguang Zhang, and Jian R Lu. Molecular self-assembly and applications of designer peptide amphiphiles. *Chem. Soc. Rev.*, 39(9):3480–3498, 2010.
- [16] Mischa Zelzer and Rein V Ulijn. Next-generation peptide nanomaterials: molecular networks, interfaces and supramolecular functionality. *Chem. Soc. Rev.*, 39(9):3351–3357, 2010.
- [17] Honggang Cui, Matthew J Webber, and Samuel I Stupp. Self-assembly of peptide amphiphiles: from molecules to nanostructures to biomaterials. *Peptide Science: Original Research on Biomolecules*, 94(1):1–18, 2010.
- [18] Caroline M Rufo, Yurii S Moroz, Olesia V Moroz, Jan Stöhr, Tyler A Smith, Xiaozhen Hu, William F DeGrado, and Ivan V Korendovych. Short peptides self-assemble to produce catalytic amyloids. *Nat. Chem.*, 6(4):303–309, 2014.
- [19] Lee A Solomon, Matthew E Sykes, Yimin A Wu, Richard D Schaller, Gary P Wiederrecht, and H Christopher Fry. Tailorable exciton transport in doped peptide–amphiphile assemblies. *ACS Nano*, 11(9):9112–9118, 2017.
- [20] Liam C Palmer and Samuel I Stupp. Molecular self-assembly into one-dimensional nanostructures. *Acc. Chem. Res.*, 41(12):1674–1684, 2008.
- [21] Stephen H White and William C Wimley. Hydrophobic interactions of peptides with membrane interfaces. *Biochim. Biophys. Acta Biomembr.*, 1376(3):339–352, 1998.
- [22] William C Wimley, Trevor P Creamer, and Stephen H White. Solvation energies of amino acid side chains and backbone in a family of host-guest pentapeptides. *Biochemistry*, 35(16):5109–5124, 1996.
- [23] Peter Y Chou and Gerald D Fasman. Prediction of protein conformation. *Biochemistry*, 13(2):222–245, 1974.
- [24] Pim WJM Frederix, Gary G Scott, Yousef M Abul-Haija, Daniela Kalafatovic, Charalampos G Pappas, Nadeem Javid, Neil T Hunt, Rein V Ulijn, and Tell Tuttle. Exploring the sequence space for (tri-) peptide self-assembly to design and discover new hydrogels. *Nat. Chem.*, 7(1):30, 2015.
- [25] Rohit Batra, Le Song, and Rampi Ramprasad. Emerging materials intelligence ecosystems propelled by machine learning. *Nat. Rev. Mater.*, pages 1–24, 2020.
- [26] Prasanna V Balachandran, Benjamin Kowalski, Alp Sehirlioglu, and Turab Lookman. Experimental search for high-temperature ferroelectric perovskites guided by two-step machine learning. *Nat. Commun.*, 9(1):1–9, 2018.
- [27] Turab Lookman, Prasanna V Balachandran, Dezheng Xue, John Hogden, and James Theiler. Statistical inference and adaptive design for materials discovery. *Curr. Opin. Solid State Mater. Sci.*, 21(3):121–128, 2017.
- [28] Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction, 2011.
- [29] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Trans. Comput. Intell. AI Games*, 4(1):1–43, 2012.
- [30] Pim WJM Frederix, Rein V Ulijn, Neil T Hunt, and Tell Tuttle. Virtual screening for dipeptide aggregation: toward predictive tools for peptide self-assembly. *J. Phys. Chem. Lett.*, 2(19):2380–2384, 2011.
- [31] H Bekker, HJC Berendsen, EJ Dijkstra, S Achterop, R Van Drunen, D Van der Spoel, A Sijbers, H Keegstra, B Reitsma, and MKR Renardus. Gromacs: A parallel computer for molecular dynamics simulations. In *Physics Computing*, volume 92, pages 252–256. World Scientific Singapore, 1993.

- [32] Mark James Abraham, Teemu Murtola, Roland Schulz, Szilárd Páll, Jeremy C Smith, Berk Hess, and Erik Lindahl. Gromacs: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX*, 1:19–25, 2015.
- [33] Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *5th International conference on computers and games*, pages 72–83. Springer, 2006.
- [34] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *15th European conference on machine learning*, pages 282–293. Springer, 2006.
- [35] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [36] Thaer M Dieb, Shenghong Ju, Junichiro Shiomi, and Koji Tsuda. Monte carlo tree search for materials design and discovery. *MRS Commun.*, 9(2):532–536, 2019.
- [37] Srilok Srinivasan, Rohit Batra, Henry Chan, Ganesh Kamath, Mathew J. Cherukara, and Subramanian Sankaranarayanan. Artificial intelligence guided de novo molecular design targeting COVID-19. *ChemRxiv*, 2020.
- [38] Yun-Ching Liu and Yoshimasa Tsuruoka. Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theor. Comput. Sci.*, 644:92–105, 2016.
- [39] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *J. Chem. Inf. Model.*, 50(5):742–754, 2010.
- [40] Luca Monticelli, Senthil K Kandasamy, Xavier Periole, Ronald G Larson, D Peter Tieleman, and Siewert-Jan Marrink. The martini coarse-grained force field: extension to proteins. *J. Chem. Theory Comput.*, 4(5):819–834, 2008.
- [41] Gurpreet Singh and D Peter Tieleman. Using the wimley–white hydrophobicity scale as a direct quantitative test of force fields: the martini coarse-grained model. *J. Chem. Theory Comput.*, 7(7):2316–2324, 2011.
- [42] Djurre H de Jong, Xavier Periole, and Siewert J Marrink. Dimerization of amino acid side chains: lessons from the comparison of different force fields. *J. Chem. Theory Comput.*, 8(3):1003–1014, 2012.
- [43] James D Tang, Cameron Mura, and Kyle J Lampe. Stimuli-responsive, pentapeptide, nanofiber hydrogel for tissue engineering. *J. Am. Chem. Soc.*, 141(12):4886–4899, 2019.
- [44] David E Clarke, Christopher DJ Parmenter, and Oren A Scherman. Tunable pentapeptide self-assembled β -sheet hydrogels. *Angew. Chem. Int. Ed.*, 57(26):7709–7713, 2018.
- [45] Meital Reches, Yair Porat, and Ehud Gazit. Amyloid fibril formation by pentapeptide and tetrapeptide fragments of human calcitonin. *J. Bio. Chem.*, 277(38):35475–35480, 2002.
- [46] Tom Guterman, Maayan Levin, Sofiya Kolusheva, Davide Levy, Nadav Noor, Yael Roichman, and Ehud Gazit. Real-time in-situ monitoring of a tunable pentapeptide gel–crystal transition. *Angew. Chem.*, 131(44):16016–16022, 2019.
- [47] Paraskevi L Tsiolaki, Stavros J Hamodrakas, and Vassiliki A Iconomidou. The pentapeptide lqvvr plays a pivotal role in human cystatin c fibrillization. *FEBS Lett.*, 589(1):159–164, 2015.
- [48] Marta J Krysmann, Valeria Castelletto, Antonios Kelarakis, Ian W Hamley, Rohan A Hule, and Darrin J Pochan. Self-assembly and hydrogelation of an amyloid peptide fragment. *Biochemistry*, 47(16):4597–4605, 2008.
- [49] Jilie Kong and Shaoning Yu. Fourier transform infrared spectroscopic analysis of protein secondary structures. *Acta Biochim. Biophys. Sin.*, 39(8):549–559, 2007.

- [50] Kazuo Fujiwara, Hiromi Toda, and Masamichi Ikeguchi. Dependence of α -helical and β -sheet amino acid propensities on the overall protein fold type. *BMC Struct. Biol.*, 12(1):1–15, 2012.
- [51] RDKit open source toolkit for cheminformatics.
- [52] Alberto Gobbi and Dieter Poppinger. Genetic optimization of combinatorial libraries. *Biotechnol. Bioeng.*, 61(1):47–54, 1998.
- [53] William Humphrey, Andrew Dalke, and Klaus Schulten. Vmd: visual molecular dynamics. *J. Mol. Graph.*, 14(1):33–38, 1996.
- [54] martinize.py.
- [55] Herman JC Berendsen, JPM van Postma, Wilfred F van Gunsteren, ARHJ DiNola, and Jan R Haak. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.*, 81(8):3684–3690, 1984.
- [56] Hess BPLINCS. A parallel linear constraint solver for molecular simulation. *J. Chem. Theory Comput.*, 4:116–122, 2008.
- [57] Siewert J Marrink, H Jelger Risselada, Serge Yefimov, D Peter Tieleman, and Alex H De Vries. The martini force field: coarse grained model for biomolecular simulations. *J. Phys. Chem. B*, 111(27):7812–7824, 2007.
- [58] Siewert J Marrink, Alex H De Vries, and Alan E Mark. Coarse grained model for semiquantitative lipid simulations. *J. Phys. Chem. B*, 108(2):750–760, 2004.
- [59] Semen O Yesylevskyy, Lars V Schäfer, Durba Sengupta, and Siewert J Marrink. Polarizable water model for the coarse-grained martini force field. *PLoS Comput. Biol.*, 6(6):e1000810, 2010.
- [60] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011.
- [61] Rohit Batra, Henry Chan, Ganesh Kamath, Rampi Ramprasad, Mathew J Cherukara, and Subramanian KRS Sankaranarayanan. Screening of therapeutic agents for covid-19 using machine learning and ensemble docking studies. *J. Phys. Chem. Lett.*, 11(17):7058–7065, 2020.
- [62] Chiho Kim, Anand Chandrasekaran, Tran Doan Huan, Deya Das, and Rampi Ramprasad. Polymer genome: A data-powered polymer informatics platform for property predictions. *J. Phys. Chem. C*, 122(31):17575–17585, 2018.

Methods

Monte Carlo tree search

The AI-expert generates promising peptide sequences using the Monte Carlo tree search (MCTS), which utilizes a tree-structure based search to balance the exploration-vs-exploitation trade-off. MCTS builds a shallow tree of nodes, each containing a peptide sequence, that are inter-connected in a parent-child manner. A meaning to the structure of the tree is provided by ensuring that each child node contains a sequence that is only a minor perturbation of the sequence of the parent node. Thus, similar peptide sequences occur in a tree branch. The MCTS consists of four key stages; selection: based on a *tree policy* select the leaf node that has the highest current score; expansion: add a child node (with a slightly different sequences than the parent node) to the selected leaf after taking a possible action; simulation: from the selected node, perform Monte Carlo trials of possible actions using a *rollout policy* to estimate the associated expected reward; back-propagation: pass the rewards generated by the simulated episodes to update the scores of all the parent leaves encountered while moving up the tree. Here, we emphasize the distinction between score and reward; while former is computed using the full Eq. 1, the latter represents only the left term of this

equation. Starting from a random peptide sequence assigned to the root node, the MCTS iterates between these four stages as guided by the tree and the rollout policies. This results in continual growth of the search tree (expansion) in regions which have high scores either due to high rewards (exploitation) or due to their uniqueness (exploration). An advantage of the MCTS is that if the search gets trapped in a suboptimal point, it can quickly jump to other regions in the search space by growing other branches of the tree that have high exploration score.

The UCP tree policy (Eq. 1) used in this work is discussed in the main text, along with its reward and the uniqueness function. For the rollout policy, random perturbations were introduced to the peptide sequence of the selected node depending on its depth in the tree; the higher the level of depth the smaller were the perturbations. For example, in the case of pentapeptides, at depth 0 (or seed nodes) all 5 sequences were generated randomly. However, with each increase in the depth level, the number of sequences that were allowed to change decreased by one, for instance, at depth level 1 and 2, four and three sequences were changed randomly with respect to the selected node. Further, only one of the amino acid was randomly perturbed during rollout if the depth level was equal to or greater than the sequence length.

In this work, 10 seed nodes at depth 0 with random peptide sequences initiated the search. During each rollout, 10 Monte Carlo runs were conducted to obtain the expected reward of the selected node, results of which were back-propagated to all the parent nodes to update their scores. In the scenario when no RF model was used, all of the 10 Monte Carlo runs were performed on randomly generated sequences. In contrast, within the MCTS+RF scheme 5 of the 10 runs were performed on randomly generated sequences, while the remaining 5 were screened from a pool of 500 random sequences based on their reward as approximated by the RF model. The exploration constant, c , was set to 10. This value was chosen based on the optimization study on the dataset of tripeptides, as discussed in Supporting Information. The uniqueness function $f(\theta_j)$ was computed using the Dice similarity measure of the Morgan circular fingerprint of a peptide as implemented in the open-source rdkit library [51] with the radius parameter $m=3$ and using the feature based invariance [52] (useFeatures=True).

Molecular Dynamics Simulation

Peptide (tri or penta) coordinate files were created using VMD scripting tools [53] and converted to coarse-grained (CG) representation in the MARTINI force field (version 2.2) using the open-source script martinize.py [54]. Analogous to the previous study on tripeptides [24], the secondary structure input flag `--ss = EEE` was used. Since this choice was consistent for all peptide sequences studied, it is not expected to bias our search.

Using the GROMACS code (version 5.1.2) [31, 32], 180 (300) zwitterionic pentapeptides (tripeptides) were randomly placed in a periodic cubic box of dimensions $13 \times 13 \times 13 \text{ nm}^3$ resulting in peptide concentration of 0.14 (0.23) mol/L in standard CG water. LJ interactions were shifted to zero in the range 0.9-1.2 nm, and electrostatic interactions in the range 0.0-1.2 nm for all simulations (no Particle Mesh Ewald method was used). A relative dielectric constant $\epsilon_r = 15$ was used in standard CG water simulations for screening of the electrostatic interactions, while 2.5 was used for simulations in polarizable water. To model the peptide structure, the simulations were conducted in a series. First, the box was energy minimized for 10000 steps or until forces on atoms converged to under 200 pN. Next, the minimized box was equilibrated at constant volume (NVT) for 15,000 steps of 6.125 fs, using v-rescale temperature coupling ($\tau_T = 0.1 \text{ ps}$) at around 303 K. Finally, the resulting structure was equilibrated for 2×10^6 steps of 6.125 fs using the Berendsen algorithms [55] to keep temperature ($\tau_T = 1.25 \text{ ps}$) and pressure ($\tau_P = 3 \text{ ps}$) around 303 K and 1bar, respectively. Bond lengths in aromatic side chains and the backbone-side chain bonds in I, V and Y were constrained using the LINCS algorithm [56]. The total simulation time for MCTS evaluation was 12.25 ns, which owing to the speed factor of the CG potentials [57, 58] tantamount to an ‘effective time’ of roughly 50 ns. For longer MD simulation on the screened top 100 pentapeptides, the water in the solvated energy-minimized box obtained after NVT simulation was converted to polarizable water (PW) [59] to better account for charge screening. This system was then energy-minimized again and run in the NPT ensemble for 8×10^6 steps, or roughly 200 ns effective time.

The GROMACS sasa tool was employed to compute AP values as the ratio of solvent accessible surface area (SASA) of the structures obtained at the start and the finish of the MD runs. The hydrophobicity of a peptide was computed using the Wimley–White whole-residue scale [21, 22], formulated as $\log P =$

$\sum_{a \in S} \Delta G_{\text{water-oct},a}$, where summation runs over all amino acids a in the peptide sequence S . The AP and logP values were normalized using the expression $x = (x - x_{\min}) / (x_{\max} - x_{\min})$, where x, x_{\min} and x_{\max} respectively denote the original peptide AP/logP value, and the associated possible minimum and the maximum values. For logP, the minimum and maximum values were computed by assuming all amino acids in the sequence to be either W (-2.09) or D (3.64). In the case of AP, the results on tripeptides provided the minimum and maximum values of 0.97 and 2.7, respectively.

Surrogate random forest model for peptide aggregation propensity

The random forest (RF) regression algorithm, as implemented in the scikit-learn [60], was used to learn the AP of peptides. RF is an ensemble of decision trees, that averages predictions from a large group of ‘weak models’ to overall result in a better prediction. The RF hyperparameter, i.e., the number of weak estimators, was set to 100 based on preliminary results using the dataset of tripeptides. As an input to the RF model, a 3-level hierarchical set of features, based on our past experience [61, 62], that capture different geometric and chemical information about the peptides at multiple length-scales (atomic, morphological) were considered. Further details on the model input features are provided in Supporting Information. The RF model was trained to minimize the mean absolute error (MAE). To estimate prediction errors on unseen data and showcase the improvement in the model performance with increasing training data, learning curves were generated by varying the sizes of the training and test sets. These results are included in Extended Data Figure 1. Statistically meaningful results were obtained by averaging over 10 different random test-train split.

Screening guidelines to human experts

Pentapeptide sequence design was requested from five experts. They were given minimal guidance in an effort to minimize design biasing. The guidelines were as follows: 1. Rationally design a self-assembling pentapeptide that is unmodified (i.e. select only from the 20 genetically encoded amino acids, no modified termini). 2. The peptide should assemble at neutral pH. 3. What morphology will the assembly yield? All five experts submitted multiple sequences, 29 in total. From this, we chose 11 based on the diversity of the sequences.

Peptide synthesis

Pentapeptide sequences were obtained from either the AI-expert or by the human experts as described in the main text. In an effort to minimize post-synthesis purification via lengthy HPLC methods, our solid-phase peptide synthesis (SPPS) methods were optimized to yield crude peptides with > 95% purity. SPPS of pentapeptides was carried out using Fmoc chemistry (CS Bio Co. automated peptide synthesizer, CS136XT). Preloaded Wang Resin (0.1mmol synthetic scale, Chem Impex) was used as the solid support. A solution of 20% piperidine (Sigma-Aldrich) in dimethylformamide (Fisher Chemical, Bioreagent Grade) was used as the deprotecting reagent with subsequent 5 and 20 minute deprotection times. Coupling was executed using ten fold equivalents of standard Fmoc protected amino acids (1 mmol, Chem Impex) and stoichiometric equivalents of diisopropylethylamine (DIEA, 1 mmol, Sigma Aldrich) and O-Benzotriazole-N,N,N',N'-tetramethyl-uronium-hexafluoro-phosphate (HBTU, 1mmol, Chem Impex) in DMF with a 90 minute coupling time. Final Fmoc deprotection was made following the same deprotection protocol listed above.

Upon completion of the synthesis, the resin was transferred to a 20 mL scintillation vial equipped with a stir bar. The peptide side chains were deprotected and the crude peptide removed from the peptidyl resin with a standard trifluoroacetic acid solution (10 mL, 95% TFA, 2.5% triisopropylsilane, 2.5% water) and stirred for 3 h. If a cysteine residue was present in the sequence, the deprotection solutions was adjusted with ethane dithiol (EDT, Sigma Aldrich) (10 mL, 95% TFA, EDT, 2.5% triisopropylsilane, 2.5% water). The resulting solution was filtered (fritted peptide reaction vessel equipped with a side arm, Chem Glass) into a clean 20 mL glass vial. The crude peptide was precipitated out of solution via dropwise addition of the TFA solution into cold diethyl ether (90 mL). The suspension in diethyl ether was transferred into two centrifuge tubes (50 mL Falcon tubes). The precipitate was pelleted using a centrifuge. The off-white

to white precipitate was washed thrice with cold-diethyl ether yielding the crude material. Once dry, the material was reconstituted in water and then lyophilized to obtain a white powder.

Sample preparation

The lyophilized powder of each peptide was weighed and dissolved in MilliQ water ($R = 18.2 \text{ M}\Omega$) or deuterium oxide (D_2O , Sigma Aldrich) for solution infrared experiments. The pH was adjusted to 7 with ammonium hydroxide (1 M NH_4OH or 1 M ND_4OD prepared by diluting ammonium hydroxide in water or deuterium oxide). The sample was noted to either remain in solution, precipitate, or gel immediately after adjusting the pH and after 24 hrs. The samples were either used as prepared or diluted for further characterization.

Experimental measurements

LCMS was employed for not only peptide identification and purity analysis but also for quantifying the hydrophobicity as reported by retention time in a standardized method. An Agilent Technologies HPLC Workstation (Agilent 1260 Infinity equipped with an autosampling unit and multiwavelength detector) equipped with a c18 column (Jupiter Proteo $10 \times 250 \text{ mm}$, Phenomenex) was utilized. A linear purification method was employed using a polar mobile phase water (0.1% TFA) with a 4% (v/v) per minute increase of the non-polar mobile phase acetonitrile (0.1% TFA). The sample was prepared at a concentration of $300 \mu\text{g}/\text{mL}$ in water (0.1% TFA) with injection volumes of 0.9 mL. The retention time were recorded using the Agilent OpenChem software (Extended Data Table 2). Advion Expression CMS (ESI-MS) was employed to determine the correct mass of the isolated peptides Table S2.

Sample opacity was used as an indicator of aggregation or assembly. We added $100 \mu\text{L}$ of each sample to a 96 well plate and analyzed the absorption at 800 nm (OD800nm), Tecan Platereader, Magellan Software.

ExpScore is defined as the product of the normalized retention time (RT') and the normalized sample opacity (OD800nm'). RT' was normalized to the lowest (15 min) and highest, whole number retention time (22 min). OD800 nm was normalized to the value collected for water (0.04) and the highest value observed for the peptides (1.82). A complete table can be found in Extended Data Table 2.

A Thermofisher Nicolet fourier transform infrared (FTIR) spectrometer was used for analysis of the mid-infrared region, i.e., the amide region for peptides. Each spectrum is an average of 16 scans with a resolution of 4 cm^{-1} and background corrected for D_2O . Each sample ($10 \mu\text{L}$ of a 2wt% solution) was dropcast on a CaF_2 plate equipped with a 0.025 mm Teflon spacer in a solution infrared cell (Sigma Aldrich). A second CaF_2 plate was placed on top of the Teflon spacer and the assembly was sealed. Using the same spectrometer and settings as the solution FTIR except background corrected for CaF_2 only. Each sample was diluted tenfold to 0.2 wt% and $10 \mu\text{L}$ was dropcast onto a CaF_2 plate and dried ($\sim 30 \text{ min}$).

Atomic force microscopy was obtained on a Bruker MultiMode 8 microscope using the Scanasyt mode. A silicon tip on a nitride lever was used (Scanasyt-Air Probe, Bruker). Each sample was diluted tenfold to 0.2 wt% and $100 \mu\text{L}$ was dropcast onto a freshly cleaved mica disk (top layer removed with scotch tape) affixed to a stainless steel disk. After 2 minutes the solution was removed by wicking away with filter paper. $10 \mu\text{m} \times 10 \mu\text{m}$ and $2 \mu\text{m} \times 2 \mu\text{m}$ images were collected at a scan rate of 1 Hz.

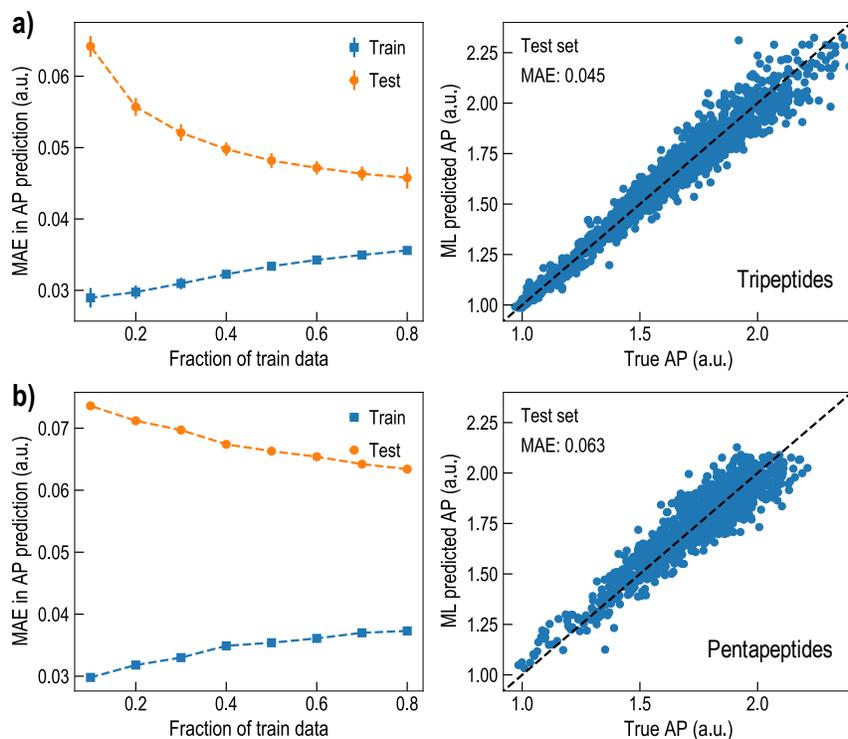
Extended Data

Extended Data Table 1: Top ranked tripeptides identified using the brute-force computational search on 8000 candidates. The score is based on the reward function r^{tri} . Abbreviations: AP, aggregation propensity; logP; hydrophobicity.

Rank	Peptide	AP	logP	Score
1	SYY	2.41	-0.96	0.214019
2	YKD	1.86	5.73	0.184754
3	DHK	1.83	6.55	0.184296
4	PSY	2.21	-0.11	0.184101
5	EYK	1.84	5.72	0.176396
6	TYT	2.19	-0.21	0.175317
7	EKW	1.89	4.34	0.174551
8	EWK	1.89	4.34	0.174551
9	SYH	2.18	-0.14	0.174447
10	YYS	2.27	-0.96	0.174427
11	DFK	1.87	4.73	0.173185
12	DKF	1.87	4.73	0.173185
13	FKD	1.87	4.73	0.173185
14	KYD	1.83	5.73	0.172508
15	KYF	2.12	0.38	0.170942
16	KYE	1.82	5.72	0.168379
17	KEY	1.82	5.72	0.168379
18	WKD	1.87	4.35	0.167202
19	DWK	1.87	4.35	0.167202
20	KYY	2.03	1.38	0.167073

Extended Data Table 2: Overall results for the synthesized pentapeptides. Computational (AP, logP) and experimental (LC(RT), OD800nm) measurements, along with the associated reward scores (r^{penta} , r^{tri}) and experimental score (ExpScore) are provided. β -sheet scale corrected r^{penta} and r^{tri} scores, respectively titled $r^{\text{penta}}_{\text{wB}}$ and $r^{\text{tri}}_{\text{wB}}$, are also included. Cases where aggregation (Agg.) was observed are marked 1 with a bold font.

Peptide	Expert	Score	Agg.	LC(RT) min	OD800nm	ExpScore	AP	logP	r^{penta}	r^{tri}	$r^{\text{penta}}_{\text{wB}}$	$r^{\text{tri}}_{\text{wB}}$
VVVVV	Human	-	1	20.1	1.469	0.218	2.06	-2.30	0.212	0.113	0.212	0.113
FKFEF	Human	-	1	21.4	0.580	0.026	2.12	1.30	0.283	0.181	0.128	0.082
VKVEV	Human	-	1	18.2	0.192	0.046	1.84	5.05	0.186	0.137	0.126	0.093
VKVFF	Human	-	1	20.9	0.314	0.024	1.95	-1.54	0.179	0.100	0.124	0.069
KFFFE	Human	-	1	20.9	1.049	0.089	1.95	1.30	0.206	0.132	0.093	0.059
KFAFD	Human	-	1	19.5	1.320	0.257	1.86	3.52	0.185	0.129	0.066	0.046
RVSVD	Human	-	0	16.5	0.050	0.005	1.66	4.99	0.117	0.086	0.061	0.045
KKFDD	Human	-	0	16.2	0.053	0.006	1.64	11.17	0.130	0.113	0.033	0.029
VKVKV	Human	-	0	17.5	0.046	0.002	1.35	4.22	0.035	0.025	0.024	0.017
KVKVK	Human	-	0	15.4	0.047	0.004	1.19	7.48	0.013	0.010	0.007	0.005
DPDPD	Human	-	0	15.8	0.049	0.004	1.01	11.20	0.000	0.000	0.000	0.000
SYCGY	AI	r^{penta}	1	18.4	1.070	0.298	2.42	0.17	0.428	0.260	0.191	0.117
FFEKF	AI	r^{penta}	1	20.3	1.810	0.241	2.24	1.30	0.345	0.221	0.156	0.100
RWLDY	AI	r^{penta}	1	20.1	0.249	0.031	2.11	1.4	0.279	0.179	0.110	0.070
KWEFY	AI	r^{penta}	1	20.7	0.335	0.031	2.20	1.92	0.332	0.218	0.142	0.093
FKIDF	AI	r^{penta}	1	21.1	0.443	0.029	2.03	1.90	0.246	0.162	0.120	0.079
KWMDF	AI	r^{penta}	1	20.9	0.155	0.010	2.28	1.97	0.378	0.248	0.139	0.092
WKPY Y	AI	r^{penta}	0	19.6	0.044	0.001	2.11	-0.57	0.255	0.150	0.100	0.059
PPPHY	AI	r^{penta}	0	17.3	0.046	0.002	2.46	-0.18	0.444	0.266	0.087	0.052
PTPCY	AI	r^{penta}	0	18.5	0.052	0.004	2.54	-0.20	0.492	0.294	0.168	0.100
FFEKF	AI	r^{tri}	1	20.3	1.810	0.241	2.17	1.30	0.308	0.197	0.139	0.089
KFFDY	AI	r^{tri}	1	20.1	0.583	0.083	2.20	2.31	0.337	0.225	0.147	0.098
KDHFY	AI	r^{tri}	0	17.9	0.048	0.003	2.30	4.13	0.422	0.301	0.162	0.115
YTEYK	AI	r^{tri}	0	16.6	0.042	0.001	2.13	5.26	0.333	0.247	0.140	0.104
WKPY Y	AI	r^{tri}	0	19.6	0.044	0.001	2.29	-0.57	0.342	0.201	0.134	0.079
EPYYK	AI	r^{tri}	0	17.2	0.044	0.002	2.26	5.15	0.410	0.303	0.131	0.097
YDPKY	AI	r^{tri}	0	17.5	0.045	0.002	2.24	5.16	0.398	0.294	0.122	0.090
KDPYY	AI	r^{tri}	0	17.8	0.054	0.005	2.22	5.16	0.385	0.284	0.118	0.087
YEPYK	AI	r^{tri}	0	17.1	0.048	0.003	2.16	5.15	0.349	0.258	0.111	0.082



Extended Data Figure 1: Performance of the random forest (RF) model to predict the computed aggregation propensity (AP) in a) tripeptides and b) pentapeptides. In both cases improvement in the RF model performance with increasing size of training data (left panel) is shown, along with an example parity plot of the test data when it constitutes 20 % of the total dataset. In a) 10 statistical runs with a random split of test-train data (from 8000 total cases) were performed; error bars denote 3σ standard deviation. In b) the test-train split (from ~ 6600 total cases using r^{Penta}) was performed in a special manner to capture the progressive improvement of the RF model during the MCTS run. Since within the MCTS+RF scheme the training data was generated in an online fashion, the RF model training set consists of AP values evaluated in the early stages of the MCTS run while the test set contains AP values evaluated in the later stage of the run. Abbreviation: MAE, mean absolute error.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SI.pdf](#)
- [peptideMCTSNatchemcode.zip](#)