

# Amino Acid Composition drives Peptide Aggregation: Predicting Aggregation for Improved Synthesis

Bálint Tamás<sup>1†</sup>, Marvin Alberts<sup>1,2,3†</sup>, Teodoro Laino<sup>2,3</sup>, and Nina Hartrampf<sup>1\*</sup>

<sup>1</sup>University of Zürich, Department of Chemistry, Winterthurerstrasse 190, 8057 Zürich, Switzerland

<sup>2</sup>IBM Research Europe, Säumerstrasse 4, 8803 Rüschlikon, Switzerland

<sup>3</sup>National Center for Competence in Research-Catalysis (NCCR-Catalysis), Zürich, Switzerland

\*nina.hartrampf@chem.uzh.ch

## Abstract

Peptide aggregation is a long-standing challenge in chemical peptide synthesis, limiting its efficiency and reliability. Although data-driven methods have enhanced our understanding of many sequence-based phenomena, no comprehensive approach addresses so-called “non-random difficult couplings” (generally linked to aggregation) during solid-phase peptide synthesis. Here, we leverage existing peptide synthesis datasets, supplemented with newly acquired experimental data, to build a predictive model that deciphers the role of individual amino acids in triggering aggregation. First, we identified and experimentally validated composition-dependent aggregation as a stronger predictor than sequence-based patterns. This insight enabled the development of a composition vector representation, allowing insights into the aggregation propensities of individual amino acids. Applying an ensemble of trained models, we predict the aggregation properties of peptides and recommend optimized synthesis conditions. By elucidating each individual amino acid’s influence, this method holds the potential to accelerate synthesis optimization through existing data, offering a robust framework for understanding and controlling peptide aggregation.

---

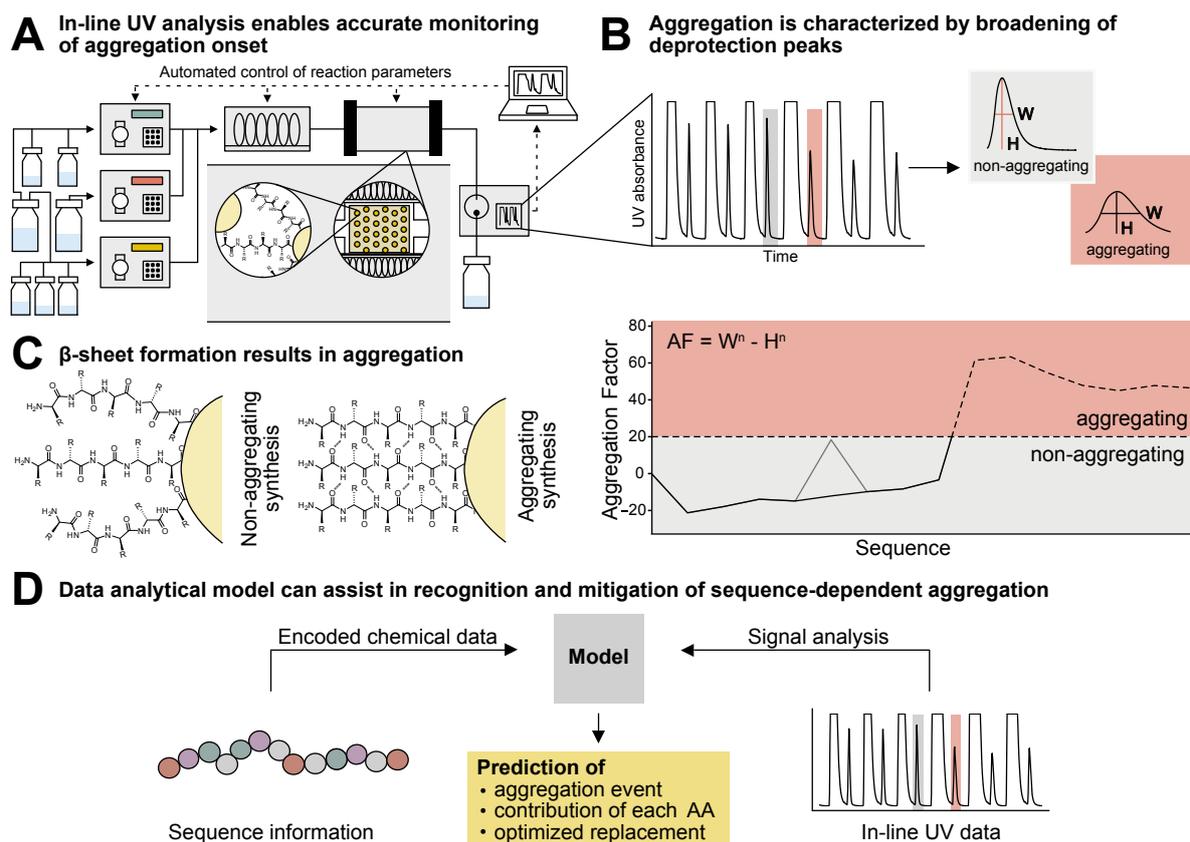
<sup>†</sup>Equal Contribution, Author order interchangeable

# 1 Introduction

Peptides and proteins play diverse biological roles, functioning as hormones, enzymes, and signalling molecules, which are critical for maintaining physiological processes. Their versatility and specificity have made them valuable therapeutic agents, driving innovations in the pharmaceutical industry. [1] Understanding their structures has been a long-standing challenge in biochemistry. [2,3] Despite key advances, human intuition alone has proven insufficient for a systematic understanding of the structure of proteins based on their primary sequence, leading to the widely known “Protein Folding Problem”. [4] With decades of accumulated data, computational methods have emerged as essential tools to predict the structure of proteins. [5–7] This evolution in methodology culminated in the development of AlphaFold and RoseTTAFold, effectively solving the problem of accurately predicting a protein’s structure from its sequence. [8,9]

While these developments have greatly enhanced our understanding of peptide and protein folding under physiological conditions, folding properties during solid-phase peptide synthesis (SPPS) remain comparatively unexplored. During SPPS, aggregation of resin- and linker-bound peptides often induces peptide folding, which can hinder synthetic efficiency and render certain sequences inaccessible. Aggregation is thought to originate from the undesired formation of  $\beta$ -sheet structures on the solid support. [10–13] This causes both truncations and deletions of the peptide sequence, often making it challenging, if not impossible, to isolate the desired peptide. Notably, even additional coupling or deprotection cycles and a large excess of amino acid do not lead to full conversion post-aggregation. Aggregation depends on several factors, such as synthesis temperature, loading of the solid support, and—most importantly—the peptide sequence and its amino acid side chain protecting groups. It has been shown that aggregation often occurs within 5–15 amino acids from the anchoring point to the resin. [14,15] Consequently, C-terminal amino acids exert the greatest influence on aggregation, with current literature suggesting that  $\beta$ -branched amino acids aggravate this effect. [14] Despite multiple attempts to understand the sequence-dependence of aggregation experimentally [16–18] and with advanced data analytics on UV data obtained from flow-SPPS [19,20], a robust method to predict aggregation and to propose an alternative synthesis strategy remains elusive.

In this study, we use machine learning on deprotection peak data collected from the in-line UV-Vis of an automated fast-flow peptide synthesiser (AFPS) [21]. This data directly correlates to the aggregation state of the peptide being synthesised on the resin (see Figure 1 B). We leverage the UV-Vis data to gain new insights into the factors contributing to peptide aggregation, including the influence of each individual amino acid. Through shuffling of peptide sequences, it was found that the composition of the peptide, rather than the specific sequence, largely determines the aggregation characteristics of a given peptide. We verify this claim through experimental results and ultimately demonstrate



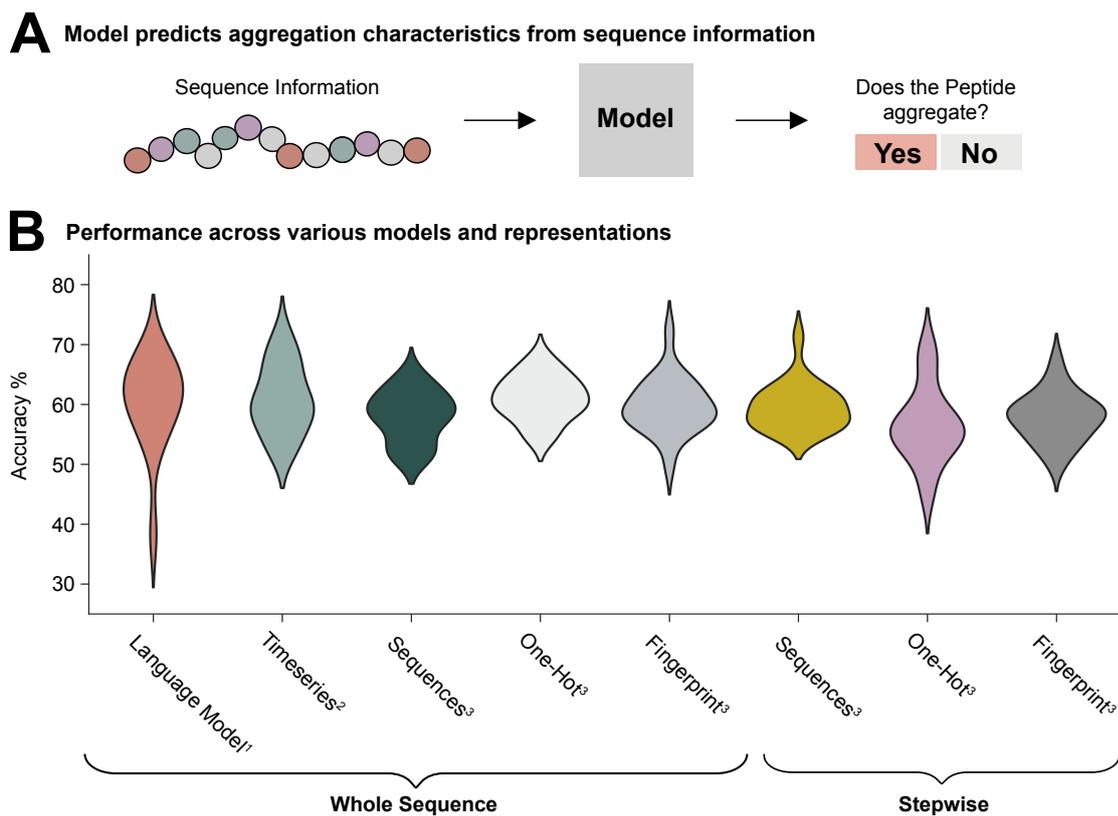
**Figure 1: Analytical data collected with an in-line UV module enables data-driven methods for synthesis analysis.** A) AFPS enables the precise monitoring of reaction kinetics, which corresponds to the aggregation of the sequences. B) Aggregation in the in-line UV traces is characterized as the broadening of the deprotection peak. Aggregation is quantified by an aggregation factor, calculated using the following formula:  $AF = W^n - H^n$ .  $W^n$ : half of the maximum height, normalized to the first peak,  $H^n$ : peak height normalized to the first peak. If  $AF > 20$ , the sequence is considered aggregating. C) Aggregation is driven by  $\beta$ -sheet formation between the growing peptide chains. D) In-line UV data collected during synthesis was leveraged to predict the occurrence of aggregation and the contribution of individual amino acids.

66 how these findings can be used to avoid aggregation.

## 67 2 Results and Discussion

### 68 Prediction of aggregation during SPPS is model- and 69 representation-independent.

70 Predicting peptide aggregation requires criteria to distinguish between aggregating and  
71 non-aggregating sequences. All data used in this study were collected on an AFPS  
72 platform equipped with an in-line UV-Vis detector monitoring coupling and deprotec-  
73 tion peaks during synthesis (Figure 1A). Deprotection peaks, which result from 9-  
74 fluorenylmethoxycarbonyl (Fmoc) removal, provide two crucial pieces of information:



**Figure 2: Prediction accuracy is independent of model or representation.** A) A variety of different models ranging from language models to classical machine learning models were trained to predict whether a given peptide sequence aggregates or not. B) Consistent prediction accuracy scores are observed across all models and representations regardless of the model, chemical representation, or if the sequence is fed stepwise or as a whole sequence.  
<sup>1</sup>: ESM 2.0, BERT. <sup>2</sup>: HIVE-COTE 2, WEASEL, TimeForest. <sup>3</sup>: XGBoost, Random Forest, KNN, Gaussian Processes.

75 their area indicates the coupling/deprotection efficiency, while their shape reflects the ag-  
 76 gregation state. [16, 20, 22, 23] Following Mohapatra et al. [19], we defined aggregation as  
 77 the deprotection peak broadening by more than 20% relative to the baseline. If any peak  
 78 during synthesis exceeds this threshold, we classify the entire sequence as aggregating. In  
 79 practice, this directly correlates to a decreased crude purity (Figure 1B).

80 Next, we used machine learning to predict the aggregation characteristics of a given pep-  
 81 tide using two datasets: One published by Mohapatra et al. [19] and one internal dataset.  
 82 Both were generated using similar AFPS platforms [21] and synthesis conditions, ensur-  
 83 ing minimal statistical deviation between the two. After curating and merging the two  
 84 datasets (see Methods 1), the combined dataset comprised 539 peptide sequences. Of  
 85 the total sequences, 420 were sourced from the Mohapatra dataset, with 48.8% showing

86 aggregation, and an additional 119 sequences from our internal dataset, where 53.8% ag-  
87 gregated. This resulted in a nearly balanced combined dataset, with 49.9% of sequences  
88 exhibiting aggregation. As aggregation typically occurs 5–15 amino acids from the an-  
89 choring point to the resin, all peptides longer than 20 amino acids were truncated and  
90 those shorter than five amino acids were discarded (see Supporting Information Figure 1  
91 for length distribution).

92 While extensive research has been conducted on identifying suitable statistical models  
93 and molecular representations for proteins, considerably less attention has been devoted  
94 to peptides. To address this gap, we explored a wide range of models and representations  
95 for peptide synthesis data. In all cases, we framed the problem as a binary classification  
96 task: Does a given peptide sequence aggregate or not? Our data was collected during  
97 synthesis, allowing for two distinct prediction approaches: either predicting the aggrega-  
98 tion characteristics of the final synthesized peptide directly or leveraging the step-by-step  
99 nature of the synthesis process. During synthesis, the peptide is elongated amino acid  
100 by amino acid, with information on whether the peptide has aggregated available at each  
101 synthesis step. We explored both approaches for the predictions (Figure 2): “Whole Se-  
102 quence” corresponds to predictions based on the final peptide sequence and “Stepwise”  
103 emphasises the step-by-step nature of the syntheses. For the step-by-step approach, all  
104 peptides are labelled as non-aggregating for the first few couplings. Once an aggregation  
105 event, i.e. broadening of the deprotection peak, occurs, all subsequent peptide couplings  
106 are labelled as aggregating.

107 To evaluate both approaches, we experimented with a range of models and represen-  
108 tations. One highly successful approach for proteins treats the amino acid sequence as  
109 text and leverages language models to predict protein properties. [24,25] Inspired by this  
110 approach, we fine-tuned a specialized protein language model (ESM2.0 [26]) as well as a  
111 generalist language model (BERT [27]) to classify whether a peptide aggregates or not.  
112 In addition to fine-tuning pretrained models, we also trained a BERT model from scratch  
113 (Figure 2, Language Models).

114 Another common data type in machine learning are time series. A time series consists  
115 of a sequence of data points collected at regular time intervals. Following this definition,  
116 the stepwise synthesis of a peptide can be considered a time series with each addition of  
117 an amino acid corresponding to one time step. We trained three state-of-the-art time se-  
118 ries classification models on this problem, representing each amino acid with a numerical  
119 token and padding to accommodate varying sequence lengths (see Figure 2, Timeseries).

120 In addition to these models, we also explored the performance of classical machine  
121 learning models (e.g. Random Forest [28], XGBoost [29]) on three different represen-  
122 tations. These representations consist of a numerical token matching the approach for  
123 timeseries models, a one-hot encoding approach, and a fingerprint-based method inspired

124 by Mohapatra et. al (see Methods 5.1 for more detail). [19] All models were trained with  
125 five-fold cross-validation and the performance of each model was assessed using the accu-  
126 racy. Surprisingly, we observed similar performance across all representations, models, or  
127 hyperparameter configurations.

128 To further guide the model, we focused on labelling the most relevant segment of the  
129 sequence, leveraging the step-by-step nature of the synthesis process. We hypothesized  
130 that the sequence preceding the aggregation point, i.e. the amino acid coupling at which  
131 aggregation occurs, is the most informative to distinguish between aggregating and non-  
132 aggregating sequences. In contrast, the remaining peptide sequence beyond the aggre-  
133 gation point contains little to no meaningful information. Therefore, we systematically  
134 investigated how many amino acids before and after the aggregation point are ideal to  
135 label as aggregating: We evaluated ranges up to ten amino acids before and after the  
136 point of aggregation (see Supporting Information Figure 2). The model's performance  
137 remained consistent regardless of the modified hyperparameters, models, or representa-  
138 tion used, provided the sequences were sufficiently long to form secondary structures ( $\geq 6$   
139 amino acids). [30] This suggests that aggregation may be determined by factors other  
140 than peptide sequence or the models were unable to effectively capture the aggregation  
141 signal from the data.

## 142 **Amino acid composition, rather than the sequence itself, influences** 143 **aggregation.**

144 The consistent results across different models and representations prompted us to ques-  
145 tion the quality and consistency of our dataset. As a validation experiment, the models  
146 were trained on a shuffled version of the peptide sequence. Assuming aggregation is highly  
147 sequence-dependent, inconsistent performance with shuffled data would indicate that the  
148 models fail to capture a sequence-specific aggregation signal.

149 We trained XGBoost models using whole sequence representation on both the origi-  
150 nal and a randomly shuffled dataset. No significant difference in accuracy was observed  
151 ( $0.580 \pm 0.035\%$  for original sequences vs  $0.579 \pm 0.036\%$  for the shuffled sequences). This  
152 result was consistent across all tested representations and models (see Supporting In-  
153 formation Section 3). These findings challenge the widely accepted view of aggregation  
154 as a phenomenon that is highly dependent on peptide sequence. [15] To investigate this  
155 further, a simplified encoding method was developed, representing each sequence as a 20-  
156 dimensional vector corresponding to the normalized composition of amino acids. Using  
157 this minimal representation, the accuracy remains comparable ( $0.610 \pm 0.038\%$ ), reinforc-  
158 ing the notion that amino acid composition might outweigh sequence order in influencing  
159 aggregation (Figure 3A).

## A Computational testing of sequence vs composition dependence

### Dataframe with original sequences

```
FREAKAEGCDITIIIS
PNGGSTTLPLSPAPPASAGLKSHPPPEK
```

- Conventional, human readable
- Varying length
- Focus on the sequential nature

**Accuracy:**  
0.580±0.035



Sequence  
randomization

### Dataframe with shuffled sequences

```
ILARDEATIGCKFES
GAAGTPLHLKLPSPGPTPSSSESAPNPPPK
```

- Human-readable with little information
- Varying length
- Focus on the composition

**Accuracy:**  
0.577±0.033



Occurrence vectors  
calculation

[A%, C%, D%, E%, F%, ... W%, Y%]

### Dataframe with occurrence vectors

```
[0.13,0.06,0.06,0.13,0.06,0.06,0.0,0.19,...,0.06,0,0,0]
[0.11,0,0,0.04,0,0.11,0.04,0,0.07,0.07,...,0.07,0,0,0]
```

- Compatible with algorithms
- Fixed length
- Focus on the composition

**Accuracy:**  
0.595±0.019

## B Experimental confirmation of sequence shuffling

Sequence	Barstar [75–90]	hGH [176–191]	GLP-1	MYC [123–143]	NBDY [53–68]	GHRH	MYC [421–439]	PCP-4 [43–62]
Native	✓	✓	✓	✓	✗	✗	✗	✗
Shuffle 1	✓	✓	✓	✓	✗	✓	✗	✗
Shuffle 2	✓	✓	✓	✓	✓	✓	✗	✗
Shuffle 3	✓	✗	✓	✓	✗	✓	✗	✗
Shuffle 4	✓	✓	✓	✓	✗	✓	✗	✗
Shuffle 5	✓	✓	✓	✓	✗	✓	✗	✗

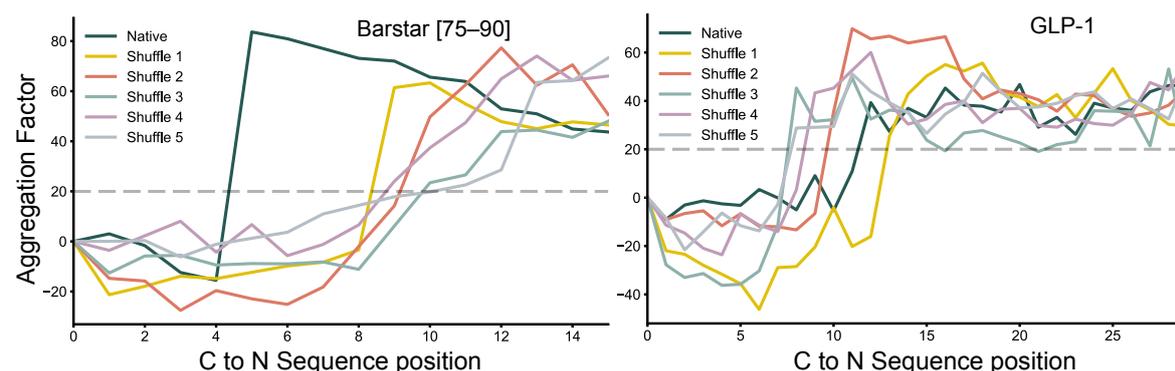


= Aggregating sequence



= Non-aggregating sequence

## C Point of aggregation is similar over multiple sequence randomizations



**Figure 3: Computational and experimental investigation of sequence shuffling on aggregation behaviour.** A) Training models on randomly shuffled sequences or only with a composition vector of the amino acids present in the peptide does not lead to a decrease in accuracy compared to training on the original sequences. B) To verify the computational results, four aggregating and four non-aggregating sequences were synthesized with five reproducible shuffles each. C) For aggregating test peptides, the point of aggregation remains consistent across the shuffled peptide sequences. Native UV-Vis traces for Barstar and GLP-1 were adapted from Tamás *et al.* [20] and Bürgisser *et al.* [13].

160 To experimentally test whether composition, rather than sequence, determines aggrega-  
161 tion, we selected eight literature-known test peptides and synthesized five randomly shuf-  
162 fled variants of each peptide (Figure 3B). Barstar[75–90] [20], hGH[176–191]Y176F [13]  
163 (abbreviated as hGH), GLP-1 [13], and MYC[123–243] [13] were selected as aggregating  
164 sequences, and NBDY[53–68] [20], GHRH [21], MYC[421–439], and PCP-4[43–62] as non-  
165 aggregating sequences. The shuffled sequences were generated through a reproducible  
166 randomization process to avoid selection bias. The peptides, ranging from 16 to 28 amino  
167 acids in length, were experimentally evaluated for aggregation behaviour during AFPS.  
168 In alignment with the in silico results, 19 out of 20 of the shuffled aggregating peptides re-  
169 tained their aggregation characteristics, while 14 out of 20 of the shuffled non-aggregating  
170 sequences also maintained their non-aggregating character (Figure 3). The majority of  
171 peptides preserve their aggregation characteristics, regardless of amino acid order, as long  
172 as the overall composition remains unchanged. In addition, the aggregation point also  
173 remains similar for the shuffled sequences (Figure 3C). This suggests that factors beyond  
174 the sequence, i.e. amino acid composition, play a prominent role in determining peptide  
175 aggregation rather than sequence information alone.

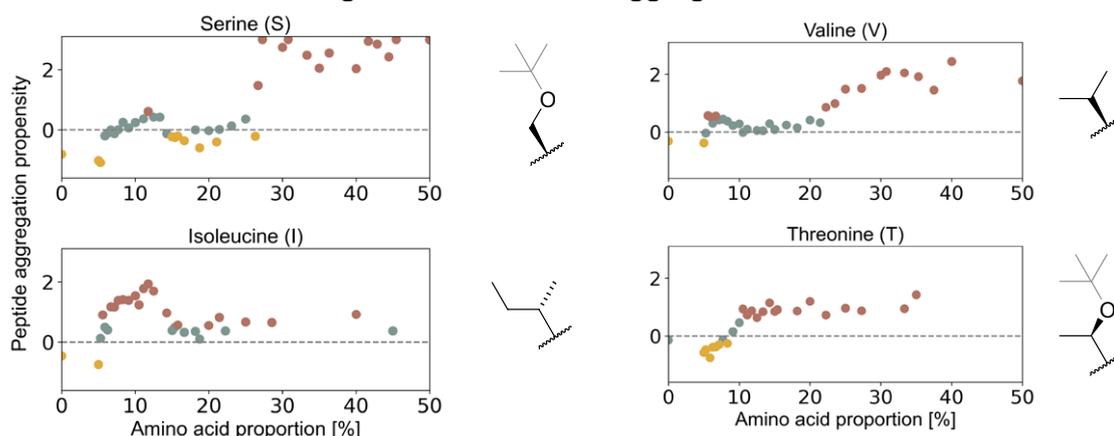
## 176 Interpretation of individual amino acid contribution to aggregation

177 To understand the impact of each individual amino acid on aggregation, we leveraged  
178 Shapley Additive Explanations (SHAP) values [31]. SHAP enables the quantification of  
179 the contribution of each amino acid to the aggregation propensity of a peptide sequence.  
180 In these experiments, the amino acid composition vector was used as the representation,  
181 establishing a direct link between the composition of amino acids and the model predic-  
182 tion.

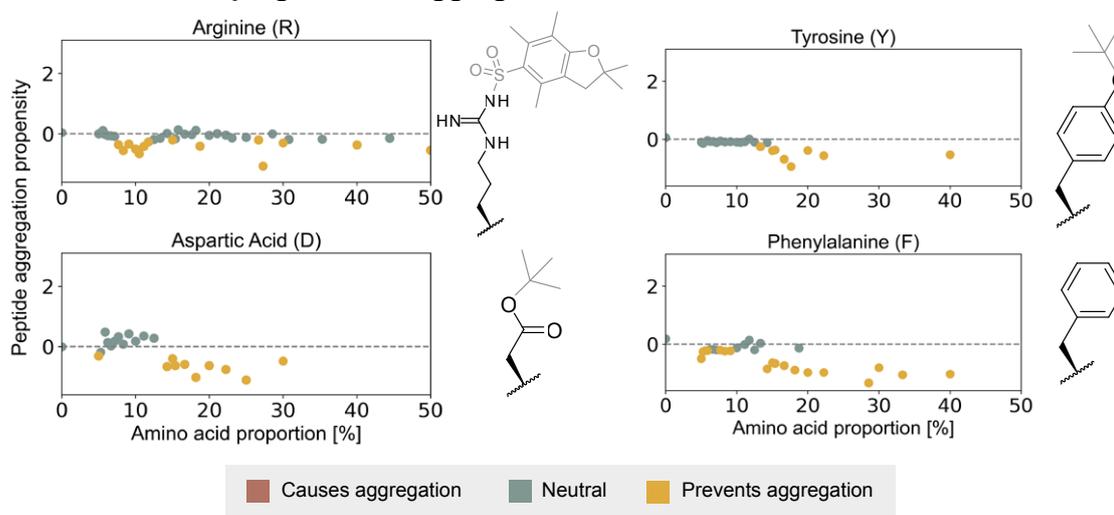
183 The analysis revealed distinct patterns in how different amino acids influence aggrega-  
184 tion (Figure 4). Amino acids such as Ser(*t*-Bu), Ile, Val, and Thr (*t*-Bu) were found to  
185 increase the likelihood of aggregation when present in higher proportions. Conversely, the  
186 presence of Phe, Asp(*t*-Bu), Tyr(*t*-Bu), and Arg(Pbf) tended to reduce aggregation. The  
187 remaining amino acids appeared to contribute neutrally, without a strong positive or neg-  
188 ative effect (see Supporting Information Section 5). While our analysis revealed peptide  
189 composition to be predominantly driving aggregation, other factors influence aggregation  
190 as well. To this end, we investigated the effect of dipeptide motifs on aggregation, with  
191 Gly–Ser and Leu–Leu contributing most to aggregation (see Supporting Information Sec-  
192 tion 6).

193 The aggregation-promoting amino acids generally have aliphatic, non-polar side chains,  
194 which seem to facilitate intermolecular interactions and packing between peptide strands.  
195 In contrast, amino acids that inhibit aggregation often have aromatic or polar side groups,  
196 which may increase spacing and disrupt aggregation-prone structures.

## A Amino acids with the largest contribution to aggregation



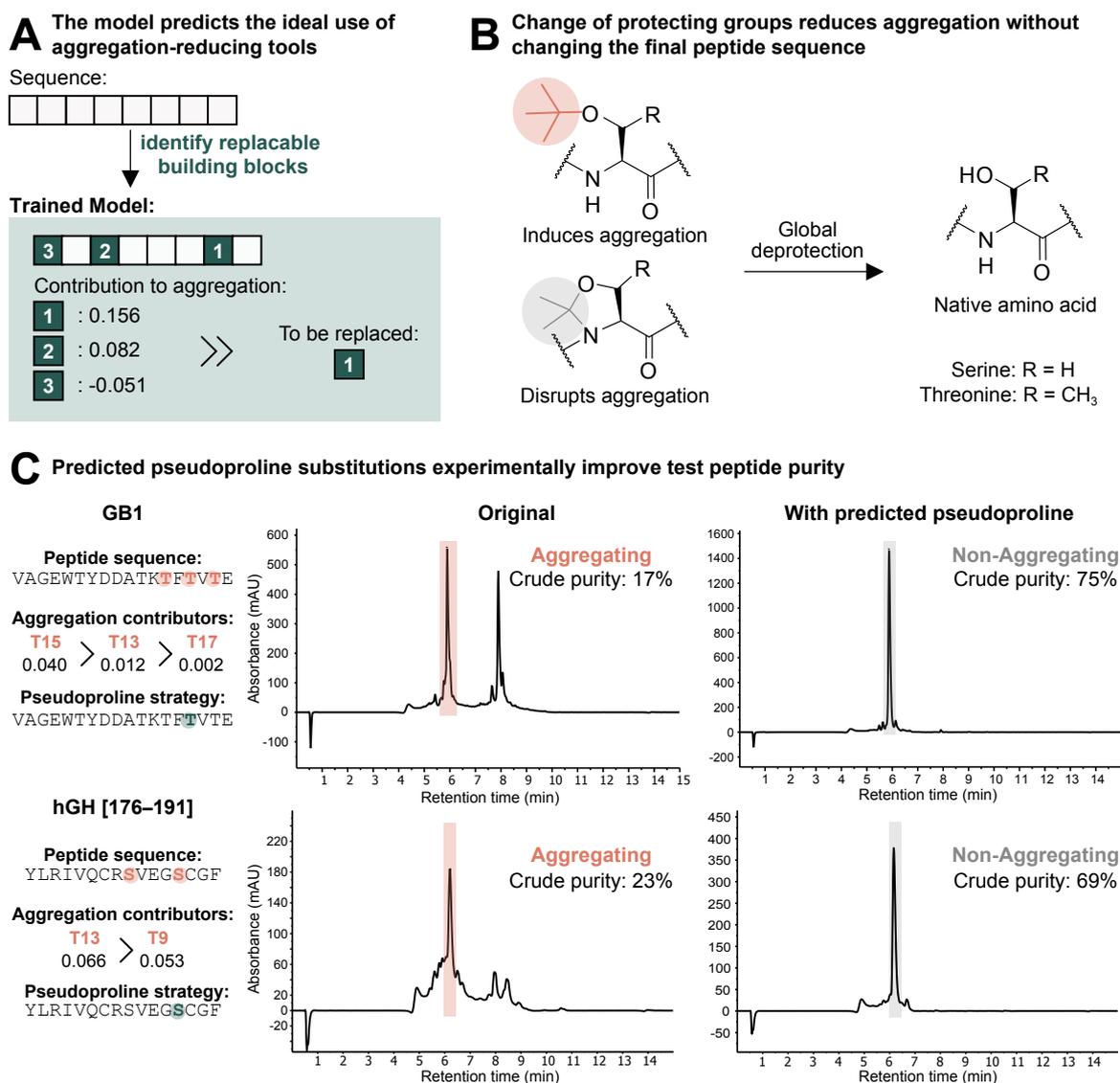
## B Amino acids helping to avoid aggregation



**Figure 4: Analysis of amino acids influencing the model's decision-making the most.** The X-axis represents the amino acid proportion in the sequences, with the Y-axis corresponding to the importance the model assigns to each data point. A positive value is associated with a higher likelihood of the model predicting aggregation and a negative one with a lower aggregation chance. A) Amino acids that contribute the most to aggregation: serine, valine, isoleucine, and threonine. B) Amino acids that contribute the least to aggregation: arginine, tyrosine, aspartic acid, and phenylalanine.

## 197 Trained models suggest conditions for improved solid-phase peptide 198 synthesis

199 The optimization of peptide synthesis can be a tedious process: As sequence-dependent  
200 events such as aggregation are difficult to predict, the usual workflow requires repetitive  
201 synthesis with the trial-and-error use of known aggregation-reducing tools. Our trained  
202 model not only enables the prediction of the aggregation propensity of a given peptide  
203 but also provides insights into how aggregation could be mitigated through strategic mod-  
204 ifications. By understanding the contributions of specific amino acids, we can predict the



**Figure 5: Leveraging the model for rational use of aggregation reduction tools to suggest improved synthesis conditions.** A) With the user input sequence and replaceable amino acids, the trained model ensemble predicts and scores the aggregation property of the sequence and predicts the contribution of the present amino acids in the early fragment of the peptide (position 2–12). This enables more effective introduction of aggregation-suppressing moieties. B) Serine and threonine, two *t*-Bu-protected amino acids with significant predicted contribution to aggregation, can also be introduced as pseudoprolines. The latter are established aggregation-reducing tools, which upon global deprotection yield the native amino acid. C) The potential of the model was tested in two known aggregating sequences: GB1 and hGH. The serines and threonines with the largest contribution were selected and replaced, resulting in a significant purity increase of 58% for the GB1 fragment and 46% for the hGH fragment.

205 most effective use of aggregation-reducing tools, such as different backbone and side chain  
 206 protecting groups. The algorithm we developed works as follows (Figure 5A): 100 models  
 207 were trained on varying splits of the data, forming an ensemble to avoid bias stemming  
 208 from the relatively small size of the dataset. The user inputs the peptide sequence and the  
 209 amino acids with available aggregation-reducing substitutions. The models then predict

210 whether the given sequence is likely to aggregate. If the sequence is predicted to be aggre-  
211 gating, the models analyse the key positions (2–12) to identify amino acids that could be  
212 substituted with their aggregation-reducing counterparts. These potential substitutions  
213 are then ranked in order of their relative contribution to aggregation, allowing the user to  
214 prioritize the most impactful changes. By substituting the highest-contributing residues,  
215 the synthesis process can be optimized to avoid aggregation issues.

216 To test this capability, we selected two aggregating sequences, hGH and GB1, and  
217 pseudoproline-protected amino acid building blocks as a widely used tool to mitigate ag-  
218 gregation [32]. The use of pseudoprolines is advantageous as they serve as an aggregation-  
219 disrupting equivalent of the two protected amino acids with the highest contribution,  
220 Ser(*t*-Bu) and Thr(*t*-Bu) (Figure 5B). For hGH, 74% of the models predicted aggrega-  
221 tion, whereas for GB1 this increased to 90%. Next, the contribution of Ser(*t*-Bu) and  
222 Thr(*t*-Bu) in the 2–12 amino acids from the resin (C-terminus) was assessed. In both cases  
223 structural motifs contributing to aggregation were identified: Ser(*t*-Bu) in position 13 for  
224 hGH and Thr(*t*-Bu) in position 15 for GB1. We synthesized both optimised sequences on  
225 the AFPS and, in both cases, we could confirm a reduction of aggregation via in-line UV  
226 signal and MS-MS. The incorporation of pseudoproline resulted in a crude purity increase  
227 from 23% to 69% for hGH and 17% to 75% for GB1. In summary, the developed algo-  
228 rithm can use the trained model to predict the aggregation property and suggest efficient  
229 incorporation of aggregation-reducing tools to increase synthetic efficiency.

### 230 3 Conclusions

231 In this study, machine learning was used as a discovery tool, uncovering a surprisingly  
232 strong composition-dependence of peptide aggregation. This finding was validated exper-  
233 imentally by synthesizing forty sequences (eight sequences, each shuffled five times). In  
234 the process, we developed a simple composition vector as a new peptide representation to  
235 investigate the aggregation character during SPPS. By leveraging the interpretability of  
236 this representation, we found that bulkier and more polar side chains or protecting groups  
237 have a tendency to reduce aggregation, while characteristically aliphatic side chains in-  
238 crease the likelihood of aggregation. In addition, we demonstrated the practical value  
239 of these findings by pinpointing the key amino acids contributing to aggregation in a  
240 given target peptide. By strategically introducing pseudoprolines at these positions, we  
241 observed a reduction in aggregation and an increase in the purity of two test sequences  
242 by 58% and 46%, respectively.

243 These findings question the understanding of aggregation as a mainly sequence-dependent  
244 event originating from intermolecular hydrogen bonding between backbones, resulting in  
245  $\beta$ -sheet structures. [10, 11, 13, 15, 30] For biological systems, it has been established that

246 amino acids with aliphatic side chains, such as valine or leucine, tend to be large con-  
247 tributors to  $\beta$ -sheet formation and aggregation. [33,34] Aromatic side chains also seem to  
248 have a major impact on the aggregation of native peptides and proteins under physiolog-  
249 ical conditions. [34] Our findings revealed that during SPPS amino acids with aliphatic  
250 side chains, such as valine or isoleucine, predominantly contribute to aggregation. Sim-  
251 ilarly, protecting groups that mimic these structures, such as *t*-Bu-protected serine or  
252 threonine, exhibit similar behavior during SPPS. In contrast to native peptides, amino  
253 acids with aromatic side chains or protecting groups, such as phenylalanine or tyrosine,  
254 tend to reduce aggregation occurrence. Furthermore, aggregation is widely considered  
255 sequence-dependent, yet our results indicate that during SPPS, amino acid composition  
256 is more influential. This discovery led to the development of the composition vector, a  
257 simplified representation of peptides allowing us to predict the onset of aggregation, while  
258 also recommending mitigation strategies.

259 Our machine learning driven approach revealed previously undetected patterns in pep-  
260 tide aggregation. The strong correlation between peptide composition and aggregation  
261 emerged only through the use of computational analysis, highlighting how machine learn-  
262 ing can discover complex relationships in chemical systems. This work demonstrates that  
263 machine learning's value in chemistry extends beyond its common applications in prop-  
264 erty prediction and molecular generation: It serves as a powerful discovery tool that can  
265 challenge established paradigms and uncover hidden patterns in molecular data.

## 266 4 Methods

### 267 4.1 Computational Methods

#### 268 4.1.1 Dataset Curation

269 The data used in this study consists of the UV-traces gathered during the SPPS of various  
270 peptides. We used the dataset published by Mohapatra *et. al.* containing 769 unique  
271 syntheses in addition to an internal dataset of 167 unique syntheses. Both datasets were  
272 combined, and all syntheses containing non-canonical amino acids, steps not performed  
273 on an AFPS (e.g. batch synthesis of a pre-chain), and synthesis of peptides with fewer  
274 than five amino acids were removed. As aggregation was reported to primarily occur  
275 between amino acids 5 and 15, only the synthesis steps of the first 20 amino acids were  
276 considered. [15] In addition, we filtered all duplicated sequences from the dataset. This  
277 reduced the size of the combined dataset to 539 unique syntheses. We defined aggre-  
278 gation as a broadening of the deprotection peak in excess of 20% compared to the first  
279 deprotection peak. During the synthesis, the addition of histidine and cysteine requires  
280 changes in the temperature of the reactor causing a broadening of the deprotection peak.

281 Following Tamas et. al. [20] we ignored these peaks and interpolated with the previous  
282 and subsequent peaks for all histidine and cysteine additions.

### 283 4.1.2 Data Processing

284 We used the following processing strategies for the peptide sequences:

285 *Step-by-Step:* Since SPPS builds the peptide sequence one amino acid at a time and  
286 aggregation information is available for each synthesis step, the problem can be framed  
287 as predicting whether a peptide sequence has aggregated at a given synthesis step. In  
288 theory, this approach has multiple advantages. It exposes the model to a considerably  
289 larger amount of training data (a total of 7.000 synthesis steps in the dataset) and enables  
290 the practitioner to not only predict whether a peptide will aggregate, but also pinpoint  
291 where aggregation occurs. In total, this approach yielded 7.000 training samples.

292 *Whole Peptide:* In this approach, we only considered the full peptide sequence and  
293 labeled it as aggregating or not aggregating. This yields 500 training samples.

### 294 4.1.3 Peptide Representation

295 *Text:* In this approach we leveraged pretrained Transformer models to predict whether a  
296 peptide aggregates or not. The peptide sequence is used as is and fed into the tokenizer  
297 of the Transformer model. ESM and BERT models were used.

298 *Sequence:* This representation converts a peptide into a vector by mapping each amino  
299 acid to a value between 1 and 20. We padded each sequence to the maximum sequence  
300 length (in this case 20) and fed this vector into the models.

301 *One Hot Encoding:* This approach works similarly to sequence representation. Instead  
302 of mapping each amino acid to a numerical value, we one-hot encoded each amino acid  
303 and concatenated the vectors. In addition, we pad the resultant vector to match the  
304 maximum sequence length.

305 *Fingerprint:* This approach is inspired by Mohapatra et. al. Here we used a Morgan  
306 Fingerprint [35] with a radius of three and a bit size of 128 to represent each amino acid.  
307 We concatenated the fingerprint for each amino acid and padded the vector with zero to  
308 a uniform length regardless of the sequence size.

309 *Composition Vector:* For a given peptide we constructed a normalized vector where  
310 each index corresponds to a specific amino acid. This vector is built as follows: Assign  
311 a fixed index to each of the 20 standard amino acids, creating a 20-dimensional vector  
312 followed by counting the number of occurrences of each amino acid and populating the  
313 corresponding vector indices. This vector is normalized by dividing by the total number  
314 of amino acids, ensuring that the vector represents the proportional composition of the  
315 peptide independent of its length.

#### 316 4.1.4 Models

317 All models were trained with five-fold cross-validation and a fixed seed.

318 *Fine-tuning ESM 2.0 and BERT:* For ESM 2.0 and BERT, the implementations pro-  
319 vided on Huggingface were used. The problem is phrased as a sequence classification  
320 task for a given peptide sequence. The entire model is fine-tuned. We used a standard  
321 Huggingface trainer with a learning rate of 2.5e-5, a batch size of 16 and a weight decay  
322 of 0.01. Adam is used as an optimiser with  $\beta_1$  of 0.9 and  $\beta_2$  of 0.99. We trained each  
323 model for 15 epochs and evaluated the model with the best validation loss. For ESM  
324 2.0 we evaluated the sizes varying from 8M, 35M, 350M to 650M whereas for BERT we  
325 evaluated the base and large checkpoints. For ESM 2.0 we only used pretrained models  
326 whereas for BERT we both fine-tuned a pretrained model and trained a model for each  
327 size from scratch.

328 All time series models were used as implemented in the SKTIME library [36] using the  
329 default parameters.

330 *HIVE COTE V2:* We used the implementation as provided by SKTIME with 500  
331 estimators and a time limit of 10 minutes. [37]

332 *WEASEL:* Weasel is used with Anova and bi-grams using “information-gain” as the  
333 binning strategy. [38]

334 *Time Forest:* The time series forest classifier is used with a minimum interval of three  
335 and 200 estimators. [39]

336 *XGBoost:* We used the implementation in the XGBoost library [29] with the default  
337 settings.

338 *Scikit-learn Models:*

339 All scikit-learn models are used with the default hyperparameters. We evaluated the  
340 Random Forest-, Gaussian Processes-, and KNN-Classifier. [40]

#### 341 4.1.5 Explainability

342 We used the Shap library [31] to explain the predictions of the models. Specifically, we  
343 leveraged the TreeExplainer and we trained a total of 50 models on random splits of the  
344 data to avoid noise in the explanations.

## 345 4.2 Experimental

### 346 4.2.1 Reagents and solvents

347 Fmoc- and side chain-protected L-amino acids (Fmoc-Ala-OH, Fmoc-Arg(Pbf)-OH, Fmoc-  
348 Asn(Trt)-OH, Fmoc-Asp(*O**t*-Bu)-OH, Fmoc-Cys(Trt)-OH, Fmoc-Gln(Trt)-OH, Fmoc-Glu(*O**t*-  
349 Bu)-OH, Fmoc-Gly-OH, Fmoc-His(Trt)-OH, Fmoc-Ile-OH, Fmoc-Leu-OH, Fmoc-Lys(Boc)-  
350 OH, Fmoc-Met-OH, Fmoc-Phe-OH, Fmoc-Pro-OH, Fmoc-Ser(*t*Bu)-OH, Fmoc-Thr(*t*-Bu)-

351 OH, Fmoc-Trp(Boc)-OH, Fmoc-Tyr(*t*-Bu)-OH, Fmoc-Val-OH) and *N*'-tetramethyluronium  
352 hexafluorophosphate (HATU) were purchased from Bachem; O-(7-azabenzotriazol-1-yl)-  
353 *N,N,N*' and (7-azabenzotriazol-1-yloxy)tripyrrolidinophosphonium hexafluorophosphate  
354 (PyAOP) were purchased from Advanced ChemTech; *N,N*-diisopropylethylamine (*i*-Pr<sub>2</sub>NEt,  
355 DIPEA, 99.5%) was purchased from Sigma-Aldrich; trifluoroacetic acid (TFA, for HPLC,  
356  $\geq 99.0\%$ ), triisopropylsilane (TIPS, 98%) and 3,6-dioxa-1,8-octane-dithiol (DODT, 95%)  
357 were purchased from Sigma-Aldrich. *N,N*-Dimethylformamide (DMF) was purchased  
358 from the from VWR International GmbH; dichloromethane (DCM,  $\geq 99.8\%$ ) was pur-  
359 chased from Fisher Scientific Ltd.; diethyl ether was purchased from Honeywell Riedel-de  
360 Haën; acetonitrile (MeCN, for HPLC gradient grade,  $\geq 99.9\%$ ) was purchased from Sigma-  
361 Aldrich. NovaPEG Rink Amide resin (0.41 or 0.20 mmol/g loading) was purchased from  
362 the Novabiochem-line from Sigma-Aldrich Canada Ltd. Piperidine ( $>99\%$ , for synthe-  
363 sis) was purchased from Carl Roth GmbH. Formic acid (reagent grade,  $>95\%$ ) and Al-  
364 draAmine trapping agent added to DMF were purchased from Sigma-Aldrich Canada  
365 Ltd.

#### 366 4.2.2 Automated flow-based peptide synthesis (AFPS)

367 Peptides were synthesized on an automated flow system built in the Hartrampf lab, which  
368 is similar to the published AFPS system. [21] Capitalized letters refer to L-amino acids.  
369 For all synthesis (referred to as standard AFPS protocol) the following settings were  
370 used for peptide synthesis: flow rate = 20 mL/min for coupling and deprotection steps,  
371 temperature = 90 °C (loop) for all canonical amino acids, except histidine and cysteine  
372 which were coupled at room temperature and 90 °C (reactor). The standard synthetic  
373 cycle involves a first step of prewashing the resin at 90 °C for 60 s at 40 mL/min. During  
374 the coupling step, three HPLC pumps are used: a 50 mL/min pump head pumps the  
375 activating agent, a second 50 mL/min pump head pumps the amino acid, and a 5.0  
376 mL/min pump head pumps *i*-Pr<sub>2</sub>NEt (neat). The 50 mL/min pump head pumps delivered  
377 0.398679 mL of liquid per pump stroke, the 5.0 mL/min pump head pumps  $3.9239 \times 10^{-2}$   
378 mL of liquid per pump stroke.

379 All peptides were prepared by AFPS on NovaPEG Rink Amide resin (0.41 mmol/g)  
380 and standard Fmoc/*t*-Bu protected amino acids (0.40 M in DMF) were coupled using  
381 HATU (0.38 M in DMF) or PyAOP (0.38 M in DMF) with DIPEA (neat, 3.0 mL/min)  
382 at a total flow rate of 20 mL/min. For amino acids D, E, F, G, I, K, L, M, P, S, W,  
383 and Y, a total volume of 6.4 mL of the “coupling solution” (i.e. amino acid (0.20 M),  
384 HATU or PyAOP (0.19 M), and DIPEA in DMF) was applied for each coupling. For  
385 amino acids A, C, H, N, Q, R, S, T, and V, a total of 10.4 mL of “coupling solution”  
386 was applied for each coupling. All amino acids except C and H were preheated at 90 °C  
387 during the activation step with HATU or PyAOP, whereas C and H were preactivated  
388 with PyAOP at room temperature. Removal of the N $\alpha$ -Fmoc group was achieved using

389 20% piperidine with 1% formic acid in DMF at a flow rate of 20 mL/min and a total  
390 volume of 6.4 mL at 90 °C. Between each coupling and deprotection step, the resin was  
391 washed with DMF (32 mL) at 90 °C with a flow rate of 40 mL/min. After completion of  
392 the peptide sequence, the resins were manually washed with DCM (3 × 5 mL) and dried  
393 under reduced pressure.

## 394 **5 Data and Models availability**

395 The code for generating the data and training the models is freely available on GitHub:  
396 <https://github.com/rxn4chemistry/AI4Aggregation> and the data is available on Zen-  
397 odo: <https://zenodo.org/records/14824562>

## 398 **6 Competing Interests**

399 All authors declare no competing interests.

## 400 **7 Acknowledgments**

401 Financial support for this project was provided by the Swiss National Science Foundation  
402 (N.H.; project grant no. 200021\_200865) and the University of Zurich. B.T. is additionally  
403 supported by a Candoc grant from the University of Zurich. This publication was created  
404 as part of NCCR Catalysis (M.A. and T.L.; grant number 180544), a National Centre of  
405 Competence in Research funded by the Swiss National Science Foundation. We would  
406 like to thank Dr. A. Jeandin and C. E. Grigglesome for helpful discussions during the  
407 preparation of this manuscript.

## 408 **References**

- 409 [1] Lei Wang, Nanxi Wang, Wenping Zhang, Xurui Cheng, Zhibin Yan, Gang Shao,  
410 Xi Wang, Rui Wang, and Caiyun Fu. Therapeutic peptides: Current applications  
411 and future directions. *Signal Transduction and Targeted Therapy*, 7(1):1–27, 2022.
- 412 [2] M. F. Perutz, M. G. Rossmann, A. F. Cullis, H. Muirhead, G. Will, and A. C. North.  
413 Structure of haemoglobin: A three-dimensional Fourier synthesis at 5.5 Å. resolution,  
414 obtained by X-ray analysis. *Nature*, 185(4711):416–422, 1960.
- 415 [3] Letícia M. F. Bertoline, Angélica N. Lima, Jose E. Krieger, and Samantha K. Teix-  
416 eira. Before and after AlphaFold2: An overview of protein structure prediction.  
417 *Frontiers in Bioinformatics*, 3, 2023.

- 418 [4] Ken A. Dill, S. Banu Ozkan, M. Scott Shell, and Thomas R. Weikl. The Protein  
419 Folding Problem. *Annual Review of Biophysics*, 37:289–316, 2008.
- 420 [5] David E. Shaw, Paul Maragakis, Kresten Lindorff-Larsen, Stefano Piana, Ron O.  
421 Dror, Michael P. Eastwood, Joseph A. Bank, John M. Jumper, John K. Salmon,  
422 Yibing Shan, and Willy Wriggers. Atomic-Level Characterization of the Structural  
423 Dynamics of Proteins. *Science*, 330(6002):341–346, 2010.
- 424 [6] Marcin J. Skwark, Daniele Raimondi, Mirco Michel, and Arne Elofsson. Improved  
425 Contact Predictions Using the Recognition of Protein Like Contact Patterns. *PLOS*  
426 *Computational Biology*, 10(11):e1003889, 2014.
- 427 [7] Kim T. Simons, Rich Bonneau, Ingo Ruczinski, and David Baker. Ab initio protein  
428 structure prediction of CASP III targets using ROSETTA. *Proteins: Structure,*  
429 *Function, and Bioinformatics*, 37(S3):171–176, 1999.
- 430 [8] Carol A. Rohl, Charlie E. M. Strauss, Kira M. S. Misura, and David Baker. Protein  
431 structure prediction using Rosetta. *Methods in Enzymology*, 383:66–93, 2004.
- 432 [9] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,  
433 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna  
434 Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard,  
435 Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas  
436 Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski,  
437 Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein,  
438 David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli,  
439 and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold.  
440 *Nature*, 596(7873):583–589, 2021.
- 441 [10] T. Miyazawa and E. R. Blout. The Infrared Spectra of Polypeptides in Various  
442 Conformations: Amide I and II Bands. *Journal of the American Chemical Society*,  
443 83(3):712–719, 1961.
- 444 [11] Mitsuaki Narita, Toshihiko Ogura, Kazuhiro Sato, and Shinya Honda. Design of  
445 the Synthetic Route for Peptides and Proteins Based on the Solubility Prediction  
446 Method. I. Synthesis and Solubility Properties of Human Proinsulin C-Peptide Frag-  
447 ments. *Bulletin of the Chemical Society of Japan*, 59(8):2433–2438, 1986.
- 448 [12] Marta Paradís-Bas, Judit Tulla-Puche, and Fernando Albericio. The road to the  
449 synthesis of “difficult peptides”. *Chemical Society Reviews*, 45(3):631–654, 2016.
- 450 [13] Héloïse Bürgisser, Elyse T. Williams, Aliénor Jeandin, Robin Lescure, Adhvitha Pre-  
451 manand, Songlin Wang, and Nina Hartrampf. A Versatile “Synthesis Tag” (SynTag)

- 452 for the Chemical Synthesis of Aggregating Peptides and Proteins. *Journal of the*  
453 *American Chemical Society*, 146(50):34887–34899, 2024.
- 454 [14] R. C. de L. Milton, Saskia C. F. Milton, and Paul A. Adams. Prediction of difficult  
455 sequences in solid-phase peptide synthesis. *Journal of the American Chemical Society*,  
456 112(16):6039–6046, 1990.
- 457 [15] Stephen B. H. Kent. Chemical Synthesis of Peptides and Proteins. *Annual Review*  
458 *of Biochemistry*, 57:957–989, 1988.
- 459 [16] Joy Bedford, Carolyn Hyde, T. Johnson, Wen Jun, D. Owen, M. Quibell, and R.c.  
460 Sheppard. Amino acid structure and “difficult sequences” in solid phase peptide  
461 synthesis. *International Journal of Peptide and Protein Research*, 40(3-4):300–307,  
462 1992.
- 463 [17] P. Y. Chou and G. D. Fasman. Conformational parameters for amino acids in he-  
464 lical, beta-sheet, and random coil regions calculated from proteins. *Biochemistry*,  
465 13(2):211–222, 1974.
- 466 [18] P. Y. Chou and G. D. Fasman. Prediction of protein conformation. *Biochemistry*,  
467 13(2):222–245, 1974.
- 468 [19] Somesh Mohapatra, Nina Hartrampf, Mackenzie Poskus, Andrei Loas, Rafael Gómez-  
469 Bombarelli, and Bradley L. Pentelute. Deep Learning for Prediction and Optimiza-  
470 tion of Fast-Flow Peptide Synthesis. *ACS Central Science*, 6(12):2277–2286, 2020.
- 471 [20] Bálint Tamás, Pietro Luigi Willi, Héloïse Bürgisser, and Nina Hartrampf. A robust  
472 data analytical method to investigate sequence dependence in flow-based peptide  
473 synthesis. *Reaction Chemistry & Engineering*, 9(4):825–832, 2024.
- 474 [21] N. Hartrampf, A. Saebi, M. Poskus, Z. P. Gates, A. J. Callahan, A. E. Cowfer,  
475 S. Hanna, S. Antilla, C. K. Schissel, A. J. Quartararo, X. Ye, A. J. Mijalis, M. D.  
476 Simon, A. Loas, S. Liu, C. Jessen, T. E. Nielsen, and B. L. Pentelute. Synthesis of  
477 proteins by automated flow chemistry. *Science*, 368(6494):980–987, 2020.
- 478 [22] Monica Dettin, Stefano Pegoraro, Paolo Rovero, Silvio Bicciato, Andrea Bagno, and  
479 Carlo Di Bello. SPPS of difficult sequences. *The Journal of Peptide Research*,  
480 49(1):103–111, 1997.
- 481 [23] Alexander J. Mijalis, Dale A. Thomas, Mark D. Simon, Andrea Adamo, Ryan Beau-  
482 mont, Klavs F. Jensen, and Bradley L. Pentelute. A fully automated flow-based  
483 approach for accelerated peptide synthesis. *Nature Chemical Biology*, 13(5):464–466,  
484 2017.

- 485 [24] Dan Ofer, Nadav Brandes, and Michal Linial. The language of proteins: NLP,  
486 machine learning & protein sequences. *Computational and Structural Biotechnology*  
487 *Journal*, 19:1750–1758, 2021.
- 488 [25] Jia-Ying Chen, Jing-Fu Wang, Yue Hu, Xin-Hui Li, Yu-Rong Qian, and Chao-Lin  
489 Song. Evaluating the advancements in protein language models for encoding strate-  
490 gies in protein function prediction: a comprehensive review. *Frontiers in Bioengi-*  
491 *neering and Biotechnology*, 13, 2025.
- 492 [26] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu,  
493 Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos San-  
494 tos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander  
495 Rives. Evolutionary-scale prediction of atomic-level protein structure with a lan-  
496 guage model. *Science*, 379(6637):1123–1130, 2023.
- 497 [27] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT:  
498 Pre-training of Deep Bidirectional Transformers for Language Understanding, 2019.  
499 arXiv:1810.04805.
- 500 [28] Leo Breiman. Random Forests. *Machine Learning*, 45(1):5–32, 2001.
- 501 [29] Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System, 2016.  
502 arXiv:1603.02754.
- 503 [30] V. N. Rajasekharan Pillai and Manfred Mutter. Conformational studies of  
504 poly(oxyethylene)-bound peptides and protein sequences. *Accounts of Chemical Re-*  
505 *search*, 14(4):122–130, 1981.
- 506 [31] Scott Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predic-  
507 tions, 2017. arXiv:1705.07874.
- 508 [32] Thomas Haack and Manfred Mutter. Serine derived oxazolidines as secondary struc-  
509 ture disrupting, solubilizing building blocks in peptide synthesis. *Tetrahedron Letters*,  
510 33(12):1589–1592, 1992.
- 511 [33] James S. Nowick and Shabana Insaf. The Propensities of Amino Acids To Form  
512 Parallel  $\beta$ -Sheets. *Journal of the American Chemical Society*, 119(45):10903–10908,  
513 1997.
- 514 [34] Jiaqi Wang, Zihan Liu, Shuang Zhao, Tengyan Xu, Stan Z. Li, and Wenbin Li.  
515 Aggregation Rules of Short Peptides. *JACS Au*, (9):3567 – 3580, 2024.
- 516 [35] H. L. Morgan. The Generation of a Unique Machine Description for Chemical  
517 Structures-A Technique Developed at Chemical Abstracts Service. *Journal of Chem-*  
518 *ical Documentation*, 5(2):107–113, 1965.

- 519 [36] Markus Löning, Anthony Bagnall, Sajaysurya Ganesh, Viktor Kazakov, Jason Lines,  
520 and Franz J. Király. *sktime: A Unified Interface for Machine Learning with Time*  
521 *Series*, 2019. arXiv:1909.07872.
- 522 [37] Matthew Middlehurst, James Large, Michael Flynn, Jason Lines, Aaron Bostrom,  
523 and Anthony Bagnall. *HIVE-COTE 2.0: a new meta ensemble for time series clas-*  
524 *sification*, 2021. arXiv:2104.07551.
- 525 [38] Patrick Schäfer and Ulf Leser. *WEASEL 2.0 – A Random Dilated Dictionary Trans-*  
526 *form for Fast, Accurate and Memory Constrained Time Series Classification*, 2023.  
527 arXiv:2301.10194.
- 528 [39] Houtao Deng, George Runger, Eugene Tuv, and Martyanov Vladimir. *A Time Series*  
529 *Forest for Classification and Feature Extraction*, 2013. arXiv:1302.2277.
- 530 [40] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand  
531 Thirion, Olivier Grisel, Mathieu Blondel, Andreas Müller, Joel Nothman, Gilles  
532 Louppe, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexan-  
533 dre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard  
534 Duchesnay. *Scikit-learn: Machine Learning in Python*, 2018. arXiv:1201.0490.