

Novel Few-shots Learning Neural Network for Predicting Carbohydrate-active Enzyme (CAZyme) Affinity Towards Fructo-oligosaccharides

Shaoxun Liu

University of Southern California

Yi Kou (✉ yikou@usc.edu)

University of Southern California <https://orcid.org/0000-0002-7139-2043>

Lin Chen

University of Southern California

Research article

Keywords: Machine learning, few-shots learning, neural network, CAZyme, sugar binding affinity, resistant sugar, fructo-oligosaccharides, Poisson process

Posted Date: October 19th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-92730/v1>

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

[Read Full License](#)

1 **Title Page**

2 **Novel Few-shots Learning Neural Network for Predicting Carbohydrate-active**
3 **Enzyme (CAZyme) Affinity Towards Fructo-oligosaccharides**

4 Shaoxun Liu^{1, a}, Yi Kou^{1, a, *} and Lin Chen^{a, *}

5 ¹ Co-first Authorship

6 ^aMolecular and Computational Biology, Department of Biological Sciences, University
7 of Southern California, 1050 Childs Way, Los Angeles, CA 90089, USA

8 *Correspondence should be addressed to Y.K. (yikou@usc.edu) or L.C.
9 (linchen@usc.edu)

10

11 **Abstract**

12 **Background:** The enzymatic activity of the microbiome toward carbohydrates in the
13 human digestive system is of enormous health significance (Zou, Y., et al., 2019; Pinard,
14 D., et al., 2015). Predicting how carbohydrates through food intake may affect the
15 distribution and balance of gut microbiota remains a major challenge. Understanding
16 the enzyme-substrate specificity relationship of the carbohydrate-active enzyme
17 (CAZyme) encoded by the vast gut microbiome will be an important step to address
18 this question. In this study, we seek to establish an *in-silico* approach to studying the
19 enzyme-substrate binding interaction.

20 **Results:** We focused on the key carbohydrate-active enzyme (CAZyme) and
21 established a novel Poisson noise-based few-shots learning neural network (pFSLNN)

22 for predicting the binding affinity of indigestible carbohydrates. This approach
23 achieved higher accuracy than other classic FSLNNs, and we have also formulated new
24 algorithms for feature generation using only a few amino acid sequences. Sliding bin
25 regression is integrated with mRMR for feature selection.

26 **Conclusion:** The resulting pFSLNN is an efficient model to predict the binding affinity
27 between CAZyme and common oligosaccharides. This model can be potentially applied
28 to binding affinity prediction of other protein-ligand interactions based on limited
29 amino acid sequences.

30

31 **Key words**

32 Machine learning; few-shots learning; neural network; CAZyme; sugar binding affinity;
33 resistant sugar; fructo-oligosaccharides; Poisson process.

34

35 **Background**

36 In recent years, increased attention has been paid to the human c microbiome and its
37 health effect. Those microorganisms, mostly bacteria, inhabit the human
38 gastrointestinal tract and engage in a symbiotic relationship with their host
39 (Huttenhower, C. et al., 2012) (Conlon, M. A. and Anthony, B., 2015). The species of
40 microorganism presented in the human body varies among individuals. Up to now, over
41 2000 species of the human microbiome have been discovered, and this number is
42 predicted to increase as more human gut microbiome samples are collected (Almeida,

43 A. et al., 2019). Within those identified species, 997 species are proven to have
44 statistical significance in positively influencing human health, which is named probiotic
45 microbiome.

46

47 One conserved feature of the vast microbiome is the expression of carbohydrate-active
48 enzymes (CAZyme). CAZymes are enzymes that perform the synthesis, recognition,
49 and degradation (digestion) of carbohydrates. CAZymes are essential for the
50 microbiome to break down the complex carbohydrates from various food sources,
51 including plant cell walls and seaweeds (Huang, L. et al., 2018). Despite its prevalence,
52 the human genome only expresses approximately 17 CAZymes (Bhattacharya, T. et al.,
53 2015). As such, most of the food carbohydrates that are indigestible to human enzymes
54 also referred to as dietary fibers, are digested by CAZymes encoded in the gut
55 microbiome. Human gut microbiome CAZymes are highly diverse in sequences
56 (Huang, L. et al., 2018). Remarkably, CAZymes encoded by the gut microbiome
57 shows the adaptability to the carbohydrates accessible to the host (Jan-Hendrik
58 Hehemann, et al., PNAS, 109, 19786, 2012). These observations suggest that
59 understanding the enzyme-substrate specificity relationship of CAZyme of the gut
60 microbiome may provide a way to use specific carbohydrates (prebiotics) to modulate
61 population abundance and distribution of gut microbiota to promote probiotic effect.
62 Carbohydrate-binding modules (CBM) are non-catalytic modules of CAZymes that
63 facilitate substrate binding (Boraston, A.B. et al., 2004). This study will focus on

64 substrate binding by CBMs of CAZymes expressed by human microbiomes.

65

66 Oligosaccharides, consisting of 3-10 monosaccharides, are complex carbohydrates

67 found in a wide variety of biological systems. Oligosaccharides are abundant in

68 glycolipids and glycol proteins, where they play indispensable roles in cell recognition

69 and cell adhesion. One of the commonly seen oligosaccharides is fructo-

70 oligosaccharide (also called oligofructan), which is mainly found in fruits and

71 vegetables and has prebiotic activity as dietary fiber. Fructans can react with reactive

72 oxygen species, and this antioxidant activity can reduce potential inflammation

73 (Franco-Robles, E. and Mercedes G. L., 2015). While being a major component of the

74 human diet, Fructo-oligosaccharides cannot be digested by native human digestive

75 enzymes (Moise, A. and Maria R., 2017). They are solely digested by the human gut

76 microbiome by bacterial CAZymes (Franco-Robles, E. and Mercedes G. L., 2015).

77 However, the detailed enzyme-substrate recognition mechanisms between these

78 prebiotic oligosaccharides and microbiome encoded CAZymes remain largely

79 unexplored. To address this question, we seek to establish a high throughput and robust

80 computational approach that can be used to predict the carbohydrate substrate

81 preference by a given CAZyme of a specific bacterial species in the human gut

82 microbiome. As a first step toward this goal, here we studied the binding of four

83 model carbohydrate substrates to the active sites of CAZyme: 1-kestose, raffinose,

84 nystose, and stachyose. These are four fructo-oligosaccharides that are shown to be

85 digestible by the human microbiome (Hayakawa, K. et al., 1990), thus are used as
86 CAZyme binding substrates during protein docking to analyze their binding pattern.
87 These analyses will provide a structural basis for future exploration of the enzyme-
88 substrate specificity relationship of CAZymes in the human gut microbiome at the level
89 of big data.

90

91 To facilitate binding affinity prediction, we used protein models generated by I-Tasser
92 and binding affinity predicted by Molegro Virtual Docker (MVD). Structure simulation
93 and modeling by I-Tasser are carried out to predict the tertiary structure of a given
94 amino acid sequence (AAS). I-Tasser is by far one of the most accurate protein structure
95 prediction servers (MacCarthy, E. and Derrick, P., 2019) with more than 90% quality
96 prediction accuracy and 85.1% accuracy in assigned molecular functions (Roy, A. et al.,
97 2010). MVD is the software used for substrate-enzyme binding predictions between
98 selected CBMs and the four aforementioned oligosaccharides. This software is also
99 used in the docking analysis between chlorogenic acid and aldose reductase since it
100 provides a consistent and relatively accurate score for binding models with different
101 binding energy (Naeem, S. et al., 2013). MolDock algorithm used by this software
102 provides protein cavity and substrate binding location predictions with around 87%
103 accuracy and position deviation within 2Å (Thomsen, R. and Mikael H. Christensen.,
104 2006). MVD provides possible binding locations, binding energy scoring (rerank score),
105 as well as cavity related fragment sequences and substrate binding residuals. MVD is

106 by far the most optimum docking software considering accuracy, information output,
107 and runtime.

108

109 The simulation of binding interactions between the CAZyme CBM and substrate
110 oligosaccharide is relatively accurate but highly time-consuming. An average docking
111 process for one CAZyme on the I-Tasser server is approximately 30 hours. The advent
112 of machine learning provides an efficient approach for this time cost issue. However,
113 the difficulty for these enzymes and substrates to be in full data set simulation using
114 any machine learning models is its lack of ample sample space. Heavy data training is,
115 therefore, often impossible due to the scarcity of available enzyme/substrate structures.
116 To overcome this limitation, we sought to apply few-shots learning (FSL) ideas and
117 develop our version of protein sequence-based Poisson augmentation few-shots
118 learning network.

119

120 This study aims to establish a method of predicting enzyme-substrate (protein-ligand)
121 binding affinity across an unlimited number of proteins in a given sample ensemble,
122 based on a small sample (~50) of enzyme-substrate docking results. Few-shot learning
123 is used to generate a neural network that is capable of differentiating the distinctive
124 classes under various goals, for example, classifying pictures of different animal species
125 will small data set (Richard, Z. et al., 2017) (Li, Z. et al., 2017). This property is
126 especially important since the high variability of AAS in proteins allows testing samples

127 to consist of rather various sequences from training samples. A neural network
128 optimization algorithm of finding the loss of each round of neural network generation
129 is a core feature of few-shot learning (Garcia, V. and Joan B., 2017). The loss algorithm
130 applied is based on a prototypical neural network with adjustment of using accuracy
131 rate instead of Euclidean distance; since among the various neural networks available,
132 a prototypical neural network is the most reliable means of approach in this situation
133 thanks to its outstanding performance in the small sample space in practices (Pan, Y. et
134 al., 2019), which often outputs the prediction accuracy that has surpassed human
135 recognition (He, K. et al., 2016). By integrating with few-shot learning algorithms, the
136 prototypical neural network achieved an approximated 70% accuracy in 5-way 5-shot
137 image classification (Richard, Z. et al., 2017).

138

139 In this study, the whole set CAZyme CBMs of probiotic human microbiomes are
140 obtained from CAZy-database (Lombard, V. et al., 2014). The over 4000 proteins are
141 clustered based on K-nearest neighbors according to the primary structure. This study
142 provides the novel idea of selecting anchor protein as bases for feature generation,
143 including cavity site and protein binding site similarity calculated through fuzzy search
144 according to anchor protein binding site fragment sequences.

145

146 Aim to establish an improved few-shots learning model, we bring in the data
147 augmentation through Poisson noise, since it represents the distribution of amino acid

148 in 1D. Previous research shows that the site substitution mutation of proteins can be
149 described by the Poisson-correction method (Sadygov, R. G., 2018). Especially when
150 the substitution rate is independent among sites, Poisson -correction can best describe
151 the scenario (Grishin, N. V., 1995). In this study, since the site-dependence of
152 substitution is unknown, site-independent substitution will be assumed. In addition, we
153 mapped the data into several higher dimensions. We have also compared the Poisson
154 data augmentation with Gaussian, random, and salt-and-pepper noises.

155

156 The major significance of this study is several folds: firstly, this study takes the first
157 step towards understanding enzymatic function at the scale of the gut microbiome,
158 which is a timely topic attracting much attention. Secondly, the study establishes a
159 generalized method pipeline for future similar few shots learning in biology and is the
160 first to try FSL and noise augmentation on proteins. Since the enzyme-substrate binding
161 predictions are based on primary structure instead of the tertiary structure used in most
162 other studies, the time of protein simulation can be reduced. Thirdly, the study sets up
163 the first example for future studies of protein-substrate interactions to be performed
164 with minimal data input and limited computational power with reasonable accuracy.

165

166 **Results**

167

168 **1** Prototypical neural network for Few-shot learning

169

170 This experiment aims to provide a method of prediction of protein-ligand interaction
171 based on a small amount of labeled data, since acquiring the labeled training set using
172 I-Tasser and Molgro is the most time-consuming.

173

174 Using a small training set, though timelier, provides less accurate prediction results
175 using traditional machine learning algorithms. To achieve better performance, we
176 adopted and modified the prototypical neural network for few-shot learning. This model
177 applies to our data set in two aspects. First is that most features of the data set resemble
178 distance from a specific anchor data, which renders each data point inherent distance to
179 a calculated prototype. This feature generating technique inherently implies the protein
180 evolution tree, where proteins with similar functions from similar organisms closely
181 resemble ligand binding site structure with each other. Second is that since the training
182 set is small, multiple epochs of neural network formation are best to run to exploit the
183 random selection of the starting point of linear regression so that the neural network
184 which has the best performance during cross-validation can be selected.

185

186 Applying prototypical neural network algorithms increases the F1 score of prediction
187 for 18%, comparing the next best machine learning model SVM (table 1). Data
188 augmentation techniques are also applied to the data set. Such a technique has been
189 used in image recognition in previous studies of artificial intelligence industries.

190 Previous studies on using data augmentation in FSL incorporated gaussian noise and
191 linear transformations, and the prediction accuracy on EMNIST data and Face
192 recognition reached accuracies of 80.25% and 58.46% using 25 samples (Antoniou, A.
193 et al., 2017). Poisson noise augmentation was applied to the data set, and in optimum
194 configurations, it increased the F1 score of prediction by 8.67%. The justification and
195 discussion were stated in the following section (table 1).

196

197 **2 Effect of Poisson augmentation**

198

199 The augmenting input sequence increases the sample size, which better supports the
200 neural net formation when experimental data scarce. The natural mutation of AAS
201 sequences is independent between each site of amino acid. The mutation rate of each
202 AA remains constant, disregarding the sequence. Assuming the probability of
203 occurrence of a specific event in a small interval of a sequence is equal to the
204 macroscopic intensity, such mutation rate can be described by a Poisson distribution,
205 where the value of lambda represents the mutation rate. This mutation rate consists of
206 both total mutation rate, the possibility that an amino acid site will mutate versus will
207 not mutate, and the amino acid-specific mutation rate, the probability of which amino
208 acid the site will mutate. The amino acid-specific mutation state was generated by
209 summarizing the occurrence frequency of each amino acid in the entire sample set. This
210 likely represents the relative abundance of each amino acid. Such a method was

211 compared with an evenly distributed model, in which the former model has better
212 prediction results. A range of different total mutation rates was tested, and 10% gave
213 the best result. The increase of the total mutation rate exhibits a possible trade-off
214 between overfitting and information perseverance.

215

216 **3 Feature importance**

217

218 Minimum Redundancy Maximum Relevance (mRMR) algorithm was used to calculate
219 the effectiveness and redundancy of the features. The mRMR score of each feature was
220 calculated, and the features were rearranged accordingly. Without loss of generality, a
221 sliding bin of 10-feature was used to slide over the rearranged features. Features in the
222 bin were the only input features for the FSL algorithm. F1 scores of those trails were
223 obtained and shown in figure 2.

224

225 Features ranking 1-40 shows a decreasing trend as expected. A second peak appeared
226 at feature group 51-60, suggesting that there are features, though redundant, are capable
227 of providing more substantial information. The redundancy of features can be explained
228 by either caused by combinations between features or the underlying scaling
229 mechanism of the mRMR algorithm. The features were further rearranged according to
230 the resulting F1 values to validate the claim above. Not to lose generality, a sliding bin
231 of size 20 was used to generate three groups, with 1-10 and 51-60 in group 1, 11-20 and

232 41-50 in group 2, and 21-30 and 61-70 in group 3. The resulting F1 values are shown
233 in figure 2.

234

235 51-60 contain all 9 Sugar binding Alignment scores and no Cavity Fragment Alignment
236 scores. This can be explained by the possible redundancy of Sugar Binding Alignment
237 as they are basically, as mentioned in the method section, Cavity Fragment Alignment
238 with a sugar-binding coefficient matrix applied to the AA exchange matrix. This
239 suggests that the modification of the sugar-binding coefficient does not have an
240 apparent effect on the prediction ability of cavity fragments.

241

242 The decreasing trend shown supports the hypothesis. This suggests that the most
243 important features in this pNN neural network come from group 1 of the rearranged
244 features. Note that the F1 score of intervals 1-10 and group 1 both exceeds the F1 score
245 of all features suggesting that contradicting features downplays the predicting
246 capability of the model.

247

248 Table 2 listed out the types of features that are selected from the 71-feature to the top
249 20 features generating the highest F1 value. High percentages of α -Helix prediction
250 score and Cavity Fragment Alignment score was selected. The identified important
251 features retain intrinsic structural significance. α -Helix prediction score consists of
252 individual residue count (SSSH_01) and long strand count (SSSH_02). A higher score

253 of either feature indicates the more abundant in α -helix (figure 4a, 4b). Cavity
254 alignment and Sugar binding alignment scores show the possibility of occurrence of a
255 similar cavity pattern between the sample protein and a given anchor protein (figure 4f).
256 The whole sequence of sugar-binding score shows gives the average affinity of residues
257 in the sample protein to a specific ligand (figure 4g).

258

259 Protein samples with a high correlation in values of important features are likely to
260 preserve similar structural identities and functionality. For example, P19 is a sample
261 protein that has resulted in a perfect score aligning with anchor protein P5 (CFAL_05,
262 SBAK_05), and such a relationship is validated by AAS alignment using super
263 algorithm resulting in an RMSD of 0.389 (figure 4e). Whole sequence alignment shows
264 the level of similarity in AAS between a sample and an anchor protein. P19, though
265 possesses a drastically different sequence compared with P5, has a relatively high whole
266 sequence alignment score (WSAL_05). Thus, a high structural similarity between P5
267 and P19 is identified (figure 4c, 4d). Subsequent research indicates that P5 and P19 are
268 glycogen debranching proteins in *B. glumae* and *A. veronii* (Lim, JaeYun, et al., 2009;
269 Yang, Honghui, et al., 1996). Since *B. glumae* and *A. veronii* are of different classes, it
270 is highly possible that P5 and P19 are enzymatic proteins that underwent convergent
271 evolution. The discovery of the Glxg proteins provides evidence that this pFSLNN has
272 the potential to identify proteins of similar function regardless of phylogenic origin.
273 Using the same important features for assessment, we also identified P10 and P52 to

274 contain similar residues in their binding cavity (figure 5). Those two proteins also come
275 from different bacteria species. Despite having drastically different AAS, the two
276 proteins showed similar key residues at their binding cavities. Such ability of this model
277 can serve to accomplish the goal of classifying the human microbiome basing on
278 enzymatic functionalities instead of 16s RNA, and be beneficial to the research of the
279 functionalities of probiotics.

280

281 **Discussion**

282

283 This study focuses on the novel FSL with Poisson augmentation on data sets. This idea
284 can be used in other fields such as genomic prediction, where datasets are few. The FSL
285 model was built upon techniques in the generation of feature matrices, which can be
286 applied to prediction models on interactions with unknown causal features but has
287 symbolic labeled subjects as anchors. Another important finding is that certain features,
288 including cavity fragment similarity and α -helix pattern, are important for the prediction
289 of binding affinity for resistant sugars. Moreover, the method of evaluating features by
290 sliding bin regression can be applied to other FSL learning models.

291

292 It should be noted that there are certain limitations for this study: a) Although it
293 facilitates the acquirement of virtual calculated binding affinity data when the dataset
294 is small and especially when the mass structure simulation is not an option, the

295 generalization of this specific Poisson augmented FSL pNN directly to others are
296 relatively not easy. b) It still remains as a question of whether certain features are
297 causally representative of a special kind enzyme as predicted and selected by this model,
298 and likewise, if an apparent unknown transferable feature set can be obtained among
299 similar protein/enzyme species. c) Feature encoding combining vector embedding
300 method and traditional ways in FSL needs to be further explored for the sake of
301 prediction confidence and power. d) There are certain other FSL modes and patterns
302 which could further enhance the prediction score if Poisson augmentation is added.
303 Other methods using deep learning networks (Thapa, N. et al., 2020) may achieve better
304 results for this kind of study, but in terms of time cost and coding vector embedding, it
305 may not be very well suited to FSL framework, especially with pNN, but is definitely
306 worthy of future investigation. e) Specifically, for CAZymes, more structural features
307 involving side chain interactions with certain sugar structure types can be further
308 explored. Depending on the interacting group characteristics, evaluation scores can be
309 reranked towards certain preferences such as H-bond/aromatic stacking, and the results
310 can be regionally optimized and cross-validated globally (McCartney, L. et al., 2004).
311 Still, this study is only a small step towards understanding the CAZyme features among
312 thousands of probiotic types. The intriguing world of probiotic bacteria and their
313 CAZyme relationships, together with the charming world of FSL modeling, is definitely
314 worth substantial future works to be devoted to.

315

316 **Conclusion**

317 The study focused on the binding of 4 typical resistant sugars with key carbohydrate-
318 active enzymes (CAZyme) and established a novel Poisson noise-based few-shots
319 learning neural network (pFSLNN) for predicting the binding affinity of indigestible
320 carbohydrates. This approach achieved higher F1 scores than other classic
321 FSLNNs using Poisson noise augmentation, which has never been applied in the FSL
322 fields before. The Poisson augmentation is found to be optimal at a 10% noise level.
323 During the pFSLNN establishment, we have also formulated several new algorithms
324 for generating feature matrix depending on a few linear amino acid sequences, such as
325 sliding window fuzzy search and two-dimension threshold optimization. We have also
326 evaluated feature importance by novel sliding window method. Several discoveries
327 concerning the binding pattern of the resistant sugars have been made during the
328 pFSLNN prediction: 1) Different proteins share relatively similar binding cavities and
329 patterns concerning the same sugar substrate, with the same interaction residues and 3D
330 structures around the sugar. 2) The overall structures can be quite similar even across
331 different 16S classes with vastly distinctive sequences, which suggests that some key
332 residues and fragment parts far from the cavity are enough to reestablish the similar
333 same binding mode and the whole protein structures. These results suggest a new
334 binding function-based relationship between CAZymes and resistant sugars from the
335 structure perspective endowed by pFSLNN prediction.

336

337 **Methods**

338 **1 Data collection**

339 **1.1 Sample preparation**

340 A list of probiotic human microbiomes that contains 997 species was adopted from
341 previous research (Forster, Samuel C., et al., 2016)). Each species was searched on the
342 CAZy database (Alisdair B. et al., 2004) for its expressed CAZymes, and all CBM's
343 AAS were downloaded. This pool of CBM contains 3749 molecules. To acquire AAS
344 with distinct features, the pool of AAS was first grouped using K-mean cluster analysis.
345 Without loss of generality, $h = 500$ was selected as the cutoff line, and 10 groups were
346 generated. 6 AAS were randomly chosen from each group to make up the sample set of
347 60 AAS. Random selection after clustering can ensure AAS with different general
348 characteristics is evenly represented in the sample set.

349

350 The 60 AAS samples were uploaded to the I-Tasser server ([https://zhanglab.ccmb.
351 med.umich.edu/I-TASSER/](https://zhanglab.ccmb.med.umich.edu/I-TASSER/)) for protein structure modeling and simulation. Substrate
352 oligosaccharide models of 1-kestose (440080), raffinose (439242), nystose (166775),
353 and stachyose (439531) were obtained from PubChem databank (Berman, H.M. et al.,
354 2000). Molgro Virtual Docker was used to detect carbohydrate-binding cavities and
355 protein-ligand binding positions. The cavity was searched for each protein, and the
356 binding position search was performed within a 15 Å radius around the center of the
357 cavity after considering the sugar substrate sizes in this study. 10 binding simulations

358 were performed for each protein-ligand pair. The binding position with the lowest
359 Rerank score was recorded. For each AAS-oligosaccharide pair, ones with Rerank score
360 below -100 were labeled as 1, representing binding, and others were labeled as 0,
361 representing non-binding. Each AAS thus has four labels.

362

363 **1.2 Anchor AAS selection**

364 Assuming each AAS in different groups is distinct, one AAS from each group (10 in
365 total) were selected as anchor AAS. Those AAS were not used as testing samples in the
366 following few-shot learning process. For those 10 AAS, residues that are within 6Å
367 (Biro, J. C., 2006) of the cavity site were recorded as cavity related fragments with
368 connected residues in the same fragment. Fragments of less than three residues were
369 neglected. Sugar-binding fragments were also recorded based on the binding position
370 of each oligosaccharide. Those fragments were searched for in each AAS.

371

372 The key concept of the feature generation pipeline is to obtain the binding pattern of 10
373 anchor AAS. The higher similarity in secondary structure, cavity fragments, and sugar-
374 binding fragments between a tested AAS and an anchor AAS suggest a higher
375 possibility for the two proteins to have the same protein-ligand binding pattern. Anchor
376 AAS always has the maximum available score when compared to its secondary
377 structure, cavity fragment, and sugar-binding fragments; thus, they were taken as
378 feature standards by the prototypical neural network and remained in training set for

379 each round of learning.

380

381 **1.3 Feature value matrices preparation**

382 According to the secondary structure sequence returned from I-Tasser, the frequency of
383 each AA appearing as each general secondary structure type (Helix, Sheet, Coil) was
384 recorded. These data were used to predict secondary structure.

385

386 The AA exchange matrix was adopted from Lev Y. Yampolsky and Arlin Stoltzfus's
387 research on The Exchangeability of Amino Acids in Proteins (Yampolsky, L.Y. and
388 Arlin S., 2005), this matrix was used to assign similarity scores when performing fuzzy
389 search between the cavity fragments and AAS. The sugar-binding coefficient was then
390 applied to the AA exchange matrix to generate a sugar-binding exchange matrix.

391

392 **2 Neural network Preparation**

393 **2.1 Feature generation pipeline**

394 For each AAS sample, a total of 71 features (6 from secondary structure score, 10 from
395 binding cavity alignment, 10 from whole sequence alignment, 40 sugar-binding
396 fragment alignment, 4 from sugar-binding whole sequence alignment, 1 from sample
397 AAS length) were generated according to the AAS and the matrixes mentioned above:
398 6 features were generated for secondary structure score, including the estimated number
399 of promoting AA and estimated number of long consecutive representing each of the

400 three general secondary structure types. Those parameters of secondary structure give
401 hints to the overall shape of the protein, as more helix promoting AA with less helix
402 strand suggests helix strands being longer towards a rod shape. 10 features were
403 generated from cavity fragment alignment. Fuzzy search algorithm (Algorithm 2) was
404 applied to cavity fragments generated from anchor proteins on each sample AAS to
405 search for the longest succeeding fragment chain. A higher score against either anchor
406 protein implies a higher possibility for a similar cavity to form. 10 features of the whole
407 sequence alignment score implied the possibility of the whole sequence to present
408 similar interactions between the anchor protein and sample AAS. 40 features were
409 generated from the sugar-binding AA exchange matrix. The same algorithm was
410 applied, but the AA exchange matrix has been modified according to the frequency of
411 each AA binding with a respective oligosaccharide. 4 features of whole sequence sugar-
412 binding scores were generated using a fuzzy search algorithm with an interaction
413 coefficient modification to the AA exchange matrix. 1 feature of AAS length was added.
414

415 Using fragment and whole sequence similarity as a feature instead of direct and simple
416 AAS has three advantages. Firstly, the median length of sample AAS is approximately
417 530, introducing 530 features in building a neural network is unrealistically time-
418 consuming. Secondly, the properties of AA cannot be linearly represented due to 3D
419 intramolecular structures. This means that the feature matrix for each AA would be
420 indefinite and hard to be quantified in only 1 dimension. Thirdly, since the sample AAS

421 has a different length, a convolutional neural network that was to be applied would be
422 increasing its time cost. Since the aim of this study is to complete mass prediction in
423 the shortest time with only a limited sample size, applying a fixed number of features
424 that describe binding patterns would be optimum.

425

426 2.2 prototypical Neural Network (pNN) formation

427 For the 60 sample AAS, each AAS was 1:10 augmented by Poisson noise (detailed
428 description is in the data augmentation section below). A matrix of 660 AAS samples,
429 each with 71 features and 4 labels, was generated after the feature generation pipeline.
430 AAS samples augmented from the same AAS sample, including the original AAS, were
431 defined to possess the same root. The set of AAS samples was denoted D . The set of
432 110 anchor AAS were denoted D^{anchor} , where $D^{\text{anchor}} \in D$. 110 samples of 10 different
433 roots from $D - D^{\text{anchor}}$ was randomly selected as the training set, denoted as D^{train} , the
434 remaining is the testing set, denoted D^{test} . The ratio between the training and testing set
435 is 5:1.

436

437 For each epoch of FSL training, D^{train} was divided into supporting set S and query set
438 Q , where the number of samples in S and the number of samples in Q has a ratio of 4:1.
439 440 AAS samples of 40 different roots, including D^{Anchor} were used to compute the
440 prototype from S . The training algorithm of the Few-Shot learning model is the same
441 as Jake Snell, Kevin Swersky, and Richard Zemel's (Richard, Z. et al., 2017). D^{Anchor}

442 was always included in the supporting set as they provided the guidelines of each
443 feature. Anchor proteins contributed the most to the class generation as they have the
444 most distinct feature values. The neural network generated that performs the
445 classification is denoted NN. Q, consisting of the leftover 110 AAS, is used to compute
446 Loss-J of NN. Loss-J was modified to be the number of incorrect predictions of this
447 neural network on the validation set. As epochs proceed to 100, the NN with the least
448 loss-J was selected to be returned as the best neural network (bNN). The prediction
449 accuracy of bNN for D^{test} was recorded.

450

451 **3 Feature generation**

452 **3.1 Generation of secondary structure score**

453 -Feature structure

454 Each amino acid sequence (AAS) input returned 6 values that consist of a secondary
455 structure score, denoted as M_{ij} ($i \in [1,3], j \in [1,2]$). Those values included estimation
456 in the number of Amino acids (AA) promoting each classic secondary structure
457 (Helixes, β -Sheets, and random Coils) that were denoted as H, S, and C, as well as
458 the estimation of the number of long, consecutive strands (≥ 5) of AA promoting the
459 same secondary structure. To be convenient as a demo, we defined a fragment of AAS,
460 showing the consecutive occurrence of the same classic secondary structure over 5
461 times as a secondary structure strand.

462

463 The function that generates secondary structure scores, denoted GenSS(), was based on
464 the secondary structure promotion matrix (sspM) and secondary structure promotion
465 threshold matrix (sspT). sspM recorded the frequency of each AA appearing in each
466 classic secondary structure after normalization according to the average and standard
467 deviation of the training set. The normalization of this matrix allowed the values to be
468 in a statistical range. sspT contained three sets of thresholds, which are the promotion
469 bar, demotion bar, and the tolerance number. The two bars characterized AA into
470 Promoting, indifferent, and demoting for each of the three secondary structure types.
471 The promotion bar is the lower bound of the sspM value of a given AA being
472 characterized to start or succeed a secondary structure strand. The demotion bar is the
473 higher bound of the sspM value of a given AA being characterized to prohibit or
474 terminate a secondary structure strand. AA with sspM value between the two thresholds
475 is considered indifferent to a classic secondary structure type. The tolerance number is
476 the minimum number of AA in a secondary structure strand to view the next indifferent
477 AA as a successor of the ongoing secondary structure strand.

478

479 -Generation process

480 The data input that generates those matrixes consisted the AAS training set, denoted
481 A^{trian}_{ij} ($i \in [1,50], j \in [1, \text{length}(A^{\text{trian}}_i)]$), where i and j mean the j^{th} AA from the i^{th} AAS
482 of the set, and the secondary structure sequence (SSS) returned by I-Tasser, denoted
483 S^{trian}_{ij} ($i \in [1,50], j \in [1, \text{length}(S^{\text{trian}}_i)]$).

484

485 The data collection for $m_{xy} \in \text{sspM}$ ($x \in [1,20]$, $y \in [1,3]$), where m_{xy} represents the x^{th}
486 AA, or the AA “x”, and the y^{th} classic secondary structure, is shown below. m_{xy} finds
487 the portion between an AA in the helix (sheet, coil) and the total amount of the AA type,
488 normalized by the average and standard deviation of all three classic secondary
489 structure types.

490 m_{xy}

$$491 = \frac{-\overline{m}_y + \frac{\sum_{i,j \in (A_{ij}^{\text{Train}}=x)} f(A_{ij}^{\text{Train}})}{\sum_{i,j \in (A_{ij}^{\text{Train}}=x)} 1}}{\sigma_{m_y}}$$

492 $f(A_{ij}^{\text{Train}})$

$$493 = \begin{cases} 1, A_{ij}^{\text{Train}} = S_{ij}^{\text{Train}} \\ 0, A_{ij}^{\text{Train}} \neq S_{ij}^{\text{Train}} \end{cases}$$

494 The threshold matrix $t_{xy} \in \text{sspT}$ ($x \in [1,3]$, $y \in [1,3]$), where t_{xy} represents the $x = 1$
495 (higher), $x = 2$ (lower), $x = 3$ (tolerance) threshold for the y^{th} classic secondary structure,
496 was optimized through linear regression of minimizing difference between $\text{GenSS}()$
497 output value of A^{train}_{ij} and S^{train}_{ij} . A demo run of $\text{GenSS}()$ using the values from a given
498 sspT was denoted by $\text{GenSS}^t()$. The sspT was optimized when the difference between
499 the estimated value and real value (running SSS in the same algorithms gives the real
500 value) is minimized.

$$501 \frac{d \sum \frac{|\text{GenSS}^t(A^{\text{train}}) - \text{GenSS}^t(S^{\text{train}})|}{\text{GenSS}^t(S^{\text{train}})}}{dt}$$

502 = 0

(3)

503 GenSS() has two parts: the first is counting the number of each secondary structure
504 promoting AA. For the c^{th} AAS, there are:

505 M_{i1}

506 $= \sum_{j \in A_{cj}^{\text{Train}}} f(A_{cj}^{\text{Train}})$

507 $f(A_{cj}^{\text{Train}})$

508 $= \begin{cases} 1, m_{A_{cj}^{\text{Train}} i} \geq t_{1i} \\ 0, m_{A_{cj}^{\text{Train}} i} < t_{1i} \end{cases}$

509 The generation of M_{i2} is further illustrated in Figure 8.

510

511 3.2 Generation of fragment binding and whole sequence binding score

512 -AA exchange fuzzy search

513 Fuzzy search is a searching algorithm based on the sliding-window idea with a penalty
514 of the difference applied to each distinctive element in the window (Vernica, R. and
515 Chen L., 2009). The advantage of the sliding-window algorithm against the Smith-
516 Waterman algorithm is that the Smith-Waterman algorithm aims to find local
517 alignments between the two strands, which neither must include the other, while the
518 sliding-window algorithm ensures to find consecutive and including alignments. In
519 addition, the Smith-Waterman algorithm aims to find the aligning strand while the aim
520 of the fuzzy search is to return the alignment score for each site.

521

522 The substitution matrix in this alignment was the AA exchange matrix, denoted aaEX_{ij}

523 ($i \in [1,20], j \in [1, 20]$), where i is the substituting AA, and j is the substituted AA. The
524 score for the same AA substitution is 1000, and the higher score indicated better
525 substitution efficiency. One AAS of sequence and fragment were inputs for one round
526 of fuzzy search. The fragment was being searched throughout the sequence. The
527 sequence is denoted S , where S_m is the m^{th} AA of the sequence. The fragment is denoted
528 F , where F_n is the n^{th} AA of the fragment. The Fuzzy search returns a vector of length
529 m , denoted R_m . Each value of R represented the alignment score between S and F at the
530 given position. The value was the average of aaEX values substituting each AA from
531 the sequence for the AA from the fragment.

532 $R_i =$

$$533 \frac{\sum_{j=1}^{\text{length}(F)} \text{aaEX}_{S_{i+j-1}F_j}}{\text{length}(F)}$$

534

535 -Longest increasing fragment

536 A vector of cavity fragments was obtained from each anchor protein. Each fragment
537 was labeled by its order in the anchor protein sequence. A new vector V of length m
538 was first filled with placeholder values. For each R_{xm} , which x indicates the x^{th} fragment
539 from the anchor, looping through R_x , each i ($i \in [1,m]$) that has R_{xi} larger than threshold
540 value T , V_i is labeled x .

$$541 V_i = f(R_{xi}) \quad f(R_{xi})$$

$$542 = \begin{cases} 0, & R_{xi} < T \\ x, & R_{xi} \geq T \end{cases}$$

543 A second vector W of the same length as V stored the fuzzy search alignment value of

544 each corresponding position. A longest increasing subsequence searching algorithm
 545 (Aldous, D. and Persi D., 1999) was applied to V with the weight of each position
 546 modified to its corresponding value in W. The returning cavity fragment alignment
 547 score was the total weight of the longest increasing subsequence in V that was divided
 548 by the total number of AA in all the fragments.

549

550 -Two-dimensional optimization of threshold T

551 Threshold T was a crucial parameter in the search for the longest cavity strand. This
 552 parameter was decided to increase the standard deviation of each column vector of
 553 feature matrix while maintaining the minimum difference of such standard deviation
 554 across the 10 features. We denote each column vector of the 10 cavity fragment
 555 alignment of all samples as F_i ($i \in [1,10]$), where i indicates the cavity fragment
 556 alignment vector with the i^{th} anchor protein. The standard deviation of F_i using t as
 557 threshold T is denoted sdF_i^t . The standard deviation of sdF_i^t for $i \in [1,10]$, is denoted
 558 $sd(sdF^t)$. t that fulfills equation 8 was chosen as T. T value is 608 in the experimental
 559 run.

$$560 \quad \frac{d \frac{\sum_{i=1}^{10} sdF_i^t}{sd(sdF^t)}}{dt}$$

561 = 0

562

563 -Whole sequence fuzzy search

564 For the whole sequence alignment score, the shorter sequence between the anchor

565 sequence and the sample sequence was viewed as a fragment. The same algorithm in
566 the previous section is applied with the substituting and substituted AA assigned
567 according to the compared length between the two AA. AA was substituted from the
568 sample AAS to the anchor AAS. The fuzzy search alignment value was returned as the
569 whole sequence alignment score.

570

571 -Sugar binding matrix

572 4 sugar-specific binding matrix was multiplied with a coefficient to aaEX to form sugar-
573 binding AA exchange matrix denote $sbEX_{mn}^i$ ($i \in [1,4]$, $m,n \in [1,20]$), where i
574 represents 1-kestose, nystose, raffinose, and stachyose, n , and m represent substituting
575 AA m with AA n . 4 vectors of AA-sugar interaction frequency, denoted sbV_j ($j \in [1,20]$),
576 where j represents AA. For $i \in [1,4]$, the AA-sugar affinity matrix was generated by
577 counting AA residuals that appear within 5\AA (Sharma, R. et al., 2008) of the sugar-
578 binding site. This method is similar to the method introduced in the work by Misaki
579 Banno when performing AA-sugar affinity prediction (Banno, M. et al., 2017). sbM_{mn}^i
580 was a matrix that contains the ratio between sbV_n and sbV_m . A larger ratio represents a
581 higher affinity of substituted AA. $sbEX^i$ was generated by applying a sbM^i filter,
582 multiplied by a factor F , on aaEX.

$$\begin{aligned}
583 \quad & sbEX_{mn}^i \\
584 \quad & = aaEX \\
585 \quad & \times (1 + F \\
586 \quad & \times (sbEX_{mn}^i \\
587 \quad & - 1)) \tag{9}
\end{aligned}$$

588 F was also optimized using equation 8, replacing T by F.

589

590 3.3 Whole sequence sugar-binding

591 4 whole sequence sugar-binding score, one for each oligosaccharide ligand, was

592 generated. AA-sugar interaction matrix was obtained using the same method as above.

593 The difference is that AA in the whole sequence was accounted for instead of AA

594 residuals that appear within 5Å of the sugar-binding site. And the average of AA-sugar

595 interaction score for all AA in the protein was calculated as the whole sequence sugar-

596 binding score.

597

598 3.4 Poisson augmentation

599 Poisson augmentation simulates the mutation of AAS to increase the sample size. The

600 usage of Poisson distribution relied on the assumptions 1) AA mutation chance is

601 independent of AA site; 2) the effect of minor mutations will not affect sugar-binding

602 efficiency. This augmentation process required an AA frequency matrix, denoted aaFM_i

603 ($i \in [1,20]$), where i represents AA, and a mutation chance at each site C . The

604 probability mass function of Poisson distribution is given by function (10). The value
605 of $k = 1$ and the value of λ was calculated by equation (11).

606 p

$$607 = \frac{\lambda^k e^{-\lambda}}{k!}$$

608 λ_i

609 $= C$

$$610 \times \frac{aaFM_i}{\sum aaFM}$$

611 20 λ values formed an accumulated interval sequence. After an AAS was input to the
612 Poisson augmentation function, a random number from 0 to 1 was generated for each
613 site of AA. The interval in which the random number falls into determines if the AA at
614 this given site would mutate and which AA it would mutate to.

615

616 The Poisson augmentation function was run on each sample AAS for 10 times to
617 generate a set of augmented AAS of the same root. Test trials of $C = 5\%$, 10% , and 20%
618 were performed, and $C = 10\%$ gave the best prediction results. Augmented AAS ran the
619 same feature generation pipeline.

620

621 **List of abbreviations**

622 CAZyme Carbohydrate-active Enzyme

623 FSL Few-Shots Learning

624 FSLNNs Few-Shots Learning Neural Network

625 pNN prototypical Neural Network

626 pFSLNN Poisson noise-based Few-Shots Learning Neural Network

627 mRMR minimum Redundancy Maximum Relevance

628 CBM Carbohydrate-Binding Modules

629 MVD Molegro Virtual Docker

630 AA Amino Acid

631 AAS Amino Acid Sequence

632 SVM Support Vector Machines

633 sspM secondary structure promotion Matrix

634 sspT secondary structure promotion Threshold matrix

635 SSS Secondary Structure Sequence

636 aaEX amino acid EXchange matrix

637

638 **Declarations**

639 **Ethics approval and consent to participate**

640 Not applicable

641

642 **Consent for publication**

643 Not applicable

644

645 **Availability of data and materials**

646 Data available on GitHub: https://github.com/ShaoxunLiu/CAZyme_FSL

647

648 **Competing interests**

649 The authors declare that they have no competing interests.

650

651 **Authors' contributions**

652 S. L and Y. K designed the model and implemented the study. S. L, Y. K and L. C drafted
653 the manuscript. Y. K and L. C supervised the whole project. All authors have read and
654 approved the manuscript.

655

656 **Funding**

657 No funding was obtained for this study

658

659 **Acknowledgements**

660 Great thanks to Professor Richard Zemel, Nicolas Terrapon, Alisdair Boraston, Dukka
661 KC, and Bruce Tidor for their kind and informative discussion.

662

663 **Reference**

664 Aldous, D. and Persi D. (1999) Longest increasing subsequences: from patience sorting
665 to the Baik-Deift-Johansson theorem. Bulletin of the American Mathematical Society
666 36.4: 413-432.

667 Alisdair B. et al. (2004) Carbohydrate-binding modules: fine-tuning polysaccharide
668 recognition. *Biochemical journal* 382.3: 769-781.

669 Almeida, A. et al. (2019) A new genomic blueprint of the human gut microbiota. *Nature*
670 568.7753: 499-504.

671 Antoniou, A. et al. (2017) Data augmentation generative adversarial networks. *arXiv*
672 preprint arXiv:1711.04340.

673 Berman, H.M. et al. (2000) The Protein Data Bank *Nucleic Acids Research*, 28: 235-
674 242.

675 Bhattacharya, T. et al. (2015) Global profiling of carbohydrate active enzymes in human
676 gut microbiome. *PloS one* 10.11: e0142038.

677 Biro, J. C. (2006) Amino acid size, charge, hydropathy indices and matrices for protein
678 structure analysis. *Theoretical Biology and Medical Modelling* 3.1: 15.

679 Bonk, B. M. et al. (2019) Machine Learning Identifies Chemical Characteristics That
680 Promote Enzyme Catalysis. *Journal of the American Chemical Society* 141.9: 4108-
681 4118.

682 Banno, M. et al. (2017) Development of a sugar-binding residue prediction system from
683 protein sequences using support vector machine. *Computational biology and chemistry*
684 66: 36-43.

685 Boraston, A.B. et al. (2004) Carbohydrate-binding modules: fine-tuning polysaccharide
686 recognition. *Biochemical journal* 382.3: 769-781.

687 Conlon, M. A. and Anthony, B. (2015) The impact of diet and lifestyle on gut

688 microbiota and human health. *Nutrients* 7.1: 17-44.

689 Forster, Samuel C., et al. (2016) "HPMCD: the database of human microbial
690 communities from metagenomic datasets and microbial reference genomes." *Nucleic
691 acids research* 44.D1: D604-D609.

692 Franco-Robles, E. and Mercedes G. L. (2015) Implication of fructans in health:
693 immunomodulatory and antioxidant mechanisms. *The Scientific World Journal* 2015.

694 Garcia, V. and Joan B. (2017) Few-shot learning with graph neural networks. *arXiv
695 preprint arXiv:1711.04043*.

696 Grishin, N. V. (1995) Estimation of the number of amino acid substitutions per site
697 when the substitution rate varies among sites. *Journal of molecular evolution* 41.5: 675-
698 679.

699 Han, X. et al. (2018) Fewrel: A large-scale supervised few-shot relation classification
700 dataset with state-of-the-art evaluation. *arXiv preprint arXiv:1810.10147*.

701 Hayakawa, K. et al. (1990) Effects of soybean oligosaccharides on human faecal flora.
702 *Microbial Ecology in Health and Disease* 3.6: 293-303.

703 He, K. et al. (2016) Deep residual learning for image recognition. *Proceedings of the
704 IEEE conference on computer vision and pattern recognition*. 770-778

705 Huang, L. et al. (2018) dbCAN-seq: a database of carbohydrate-active enzyme
706 (CAZyme) sequence and annotation. *Nucleic acids research* 46.D1: D516-D521.

707 Huttenhower, C. et al. (2012) Structure, function and diversity of the healthy human
708 microbiome. *Nature* 486.7402: 207.

709 Li, Z. et al. (2017) Meta-sgd: Learning to learn quickly for few-shot learning. arXiv
710 preprint arXiv:1707.09835.

711 Lim, JaeYun, et al. (2009) Complete genome sequence of Burkholderia glumae BGR1.
712 Journal of bacteriology 191.11: 3758-3759.

713 Lombard, V. et al. (2014) The carbohydrate-active enzymes database (CAZy) in 2013.
714 Nucleic acids research 42.D1: D490-D495.

715 MacCarthy, E. and Derrick, P. (2019) Advances in protein super-secondary structure
716 prediction and application to protein structure prediction. Protein Supersecondary
717 Structures. Humana Press, New York, NY. 15-45.

718 McCartney, L. et al. (2004) Glycoside hydrolase carbohydrate-binding modules as
719 molecular probes for the analysis of plant cell wall polymers. Analytical biochemistry
720 326.1: 49-54.

721 Moise, A. and Maria R. (2017) The Gut Microbiome: Exploring the Connection
722 Between Microbes, Diet, and Health. ABC-CLIO.

723 Montserrat, D. et al. (2017) Training object detection and recognition CNN models
724 using data augmentation. Electronic Imaging 2017.10: 27-36.

725 Naeem, S. et al. (2013) Docking studies of chlorogenic acid against aldose reductase
726 by using molgro virtual docker software. Journal of Applied Pharmaceutical Science
727 3.1:13.

728 Ndeh, D. et al. (2017) Complex pectin metabolism by gut bacteria reveals novel
729 catalytic functions. Nature 544.7648: 65-70.

730 Pan, Y. et al. (2019) Transferrable prototypical networks for unsupervised domain
731 adaptation. Proceedings of the IEEE Conference on Computer Vision and Pattern
732 Recognition. 2239-2247

733 Pinard, D., et al. (2015) "Comparative analysis of plant carbohydrate active enZymes
734 and their role in xylogenesis." BMC genomics 16.1: 1-13.

735 Richard, Z. et al. (2017) Prototypical networks for few-shot learning. Advances in
736 neural information processing systems. 4077-4087.

737 Roy, A. et al. (2010) I-TASSER: a unified platform for automated protein structure and
738 function prediction. Nature protocols 5.4: 725-738.

739 Sadygov, R. G. (2018) Poisson model to generate isotope distribution for biomolecules.
740 Journal of proteome research 17.1: 751-758.

741 Sharma, R. et al. (2008) The interaction of carbohydrates and amino acids with aromatic
742 systems studied by density functional and semi-empirical molecular orbital calculations
743 with dispersion corrections. Physical Chemistry Chemical Physics 10.19: 2767-2774.

744 Thapa, N. et al. (2020) DeepSuccinylSite: a deep learning based approach for protein
745 succinylation site prediction. BMC bioinformatics 21: 1-10.

746 Thomsen, R. and Mikael H. Christensen. (2006) MolDock: a new technique for high-
747 accuracy molecular docking. Journal of medicinal chemistry 49.11: 3315-3321.

748 Yang, Honghui, et al. (1996) Coordinate genetic regulation of glycogen catabolism and
749 biosynthesis in Escherichia coli via the CsrA gene product. Journal of bacteriology
750 178.4: 1012-1017.

751 Vernica, R. and Chen L. (2009) Efficient top-k algorithms for fuzzy search in string
752 collections. Proceedings of the First International Workshop on Keyword Search on
753 Structured Data. 9-14.

754 Yampolsky, L.Y. and Arlin S. (2005) The exchangeability of amino acids in proteins.
755 Genetics 170.4: 1459-1472.

756 Zou, Y., et al. (2019) "1,520 reference genomes from cultivated human gut bacteria
757 enable functional microbiome analyses." Nature biotechnology 37.2: 179-185.

758

759 **Figure legends**

760 Figure 1: Fructo-Oligosaccharides. Name, CAS number, 3D structure, and 2D structure
761 of the four fructo-oligosaccharides that are used as substrates in this study.

762 Figure 2: F1 scores of FSL models trained using each 10-feature groups arranged
763 according to mRMR scores.

764 Figure3: F1 score of FSL models inputting 20 features from three groups of mRMR
765 intervals.

766 Figure 4: Structural presentation of sample proteins. a) Structure of P9 with alpha
767 helixes in red. SSSH_01 = 12, SSSH_02 = 150. b) Structure of P6 with alpha helixes
768 in red, SSSH_01 = 1, SSSH_02 = 35. c) Structure of anchor protein P5 with alpha
769 helixes in blue and beta sheets in yellow; binding ligand1-kestose is shown in red. d)
770 Structure of protein P19 with alpha helixes in blue and beta sheets in yellow; binding
771 ligand1-kestose is shown in red. e) Structural alignment of cavity fragments of P5 (blue)

772 and P19 (gray) using align function in pymol with alignment RMSD = 0.389; residue
773 labels are shown in corresponding colors. f) Sample AAS alignment of P5 (GlgX [B.
774 glumae]) and P19 (GlgX [A. veronii]) around the two cavity fragments; aligned AA of
775 a given secondary structure or belongs to a cavity fragment is shown in the color scheme.
776 g) Sample AAS alignment of non-cavity parts of P5 and P19, with AA of a given
777 secondary structure shown in the color scheme.

778 Figure 5: Predicted interaction of two CAZyme CBMs (P10 and P52) with 1-kestose.
779 P10 and P52 are highly correlated according to group 1 features, while having no
780 significant similarity in sequence. Similar amino acid residues are found in the radius
781 of the interaction of both P10 and P52, including SER, ALA, ARG, and LEU. This
782 suggests that the feature matrix may obtain inherent biological meaning. a) Interaction
783 model of P10 viewing from beta-D-fructofuranose. b) Interaction model of P10 viewing
784 from beta-D-fructofuranosyl residue. c) Interaction model of P10 viewing from alpha-
785 D-glucopyranosyl residue. d) Interaction model of P52 viewing from beta-D-
786 fructofuranose. e) Interaction model of P52 viewing from beta-D-fructofuranosyl
787 residue. f) Interaction model of P52 viewing from alpha-D-glucopyranosyl residue.

788 Figure 6: Flow chart for feature generation procedure. Arrows show derivation.

789 Figure 7: Flow chart of feature matrix generation and pFSLNN process.

790 Figure 8: This figure illustrates the bases of the algorithm estimating the number of
791 residuals of each general secondary structure type. Without loss of generality, α -Helix
792 (H) is shown as an example. A cursor scans over the input AAS to find the first H

793 promoting AA. A secondary structure strand of H begins. When the strand length is
 794 more than or equal to 5, the strand is considered a long strand. As the cursor proceeds,
 795 3 types of successor (promoting, demoting, indifferent) AA can result in 5 cases. Initial
 796 $sspTH = (-0.5, -1.6, 2)$.

797

798 **Tables**

799 Table 1: Accuracy and F1 score for experimented models

Method	Augmented	Accuracy	F1 score
Random Forest	N	67.70%	44.36%
SVM	N	72.50%	52.13%
FSL unaugmented	N	66.55%	70.13%
FSL Poisson aug (5%)	Y	67.00%	76.92%
FSL Poisson aug (10%)	Y	69.15%	78.67%
FSL Poisson aug (20%)	Y	66.55%	77.24%
FSL Poisson aug (30%)	Y	65.75%	75.38%
FSL Poisson aug (40%)	Y	66.15%	68.47%
FSL Salt & Pepper aug (10%)	Y	65.00%	67.73%
FSL Random mutation (10%)	Y	67.45%	73.22%
FSL Gaussian aug (10%)	Y	65.50%	65.87%

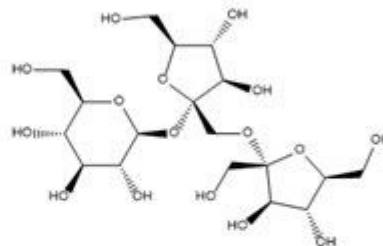
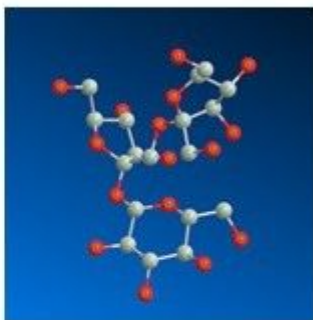
800 Table 2: Percentage of selected feature types in the top 20 features

Feature Type	Total	Top 20	Percent Selected
Cavity Fragment Alignment	10	4	40%
Whole sequence Sugar Binding	4	3	75%
Whole sequence Alignment	10	2	20%
α -Helix prediction	2	2	100%
Sugar binding Alignment	40	9	22.5%

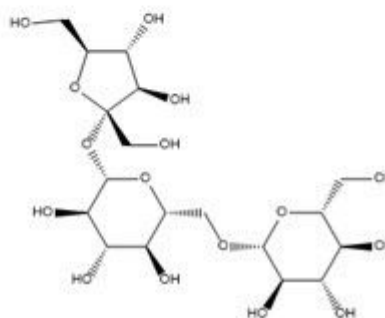
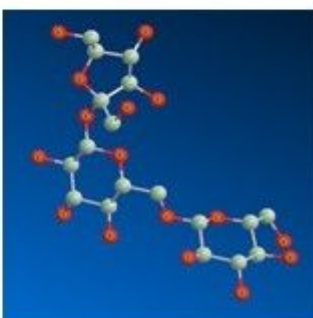
Figures

Resistant Oligossacharides

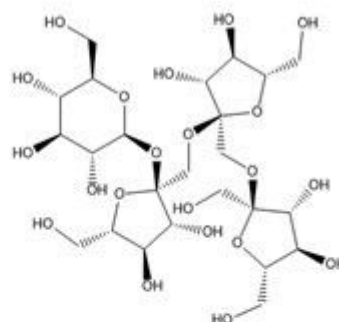
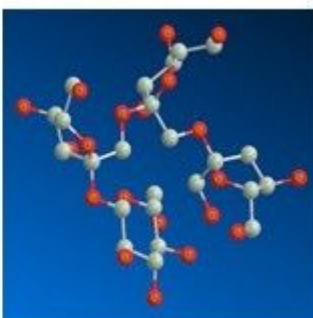
1-kestose
CAS:470-69-9



Raffinose
CAS: 512-69-6



Nystose
CAS: 13133-07-8



Stachyose:
CAS: 470-55-3

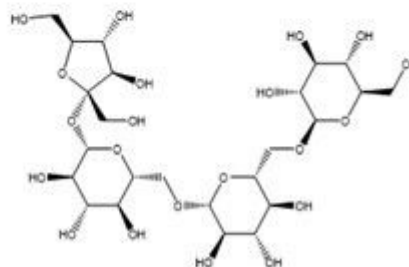
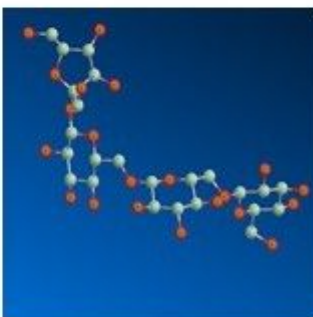


Figure 1

Fructo-Oligosaccharides. Name, CAS number, 3D structure, and 2D structure of the four fructo-oligosaccharides that are used as substrates in this study.

F1 Score of mRMR Intervals

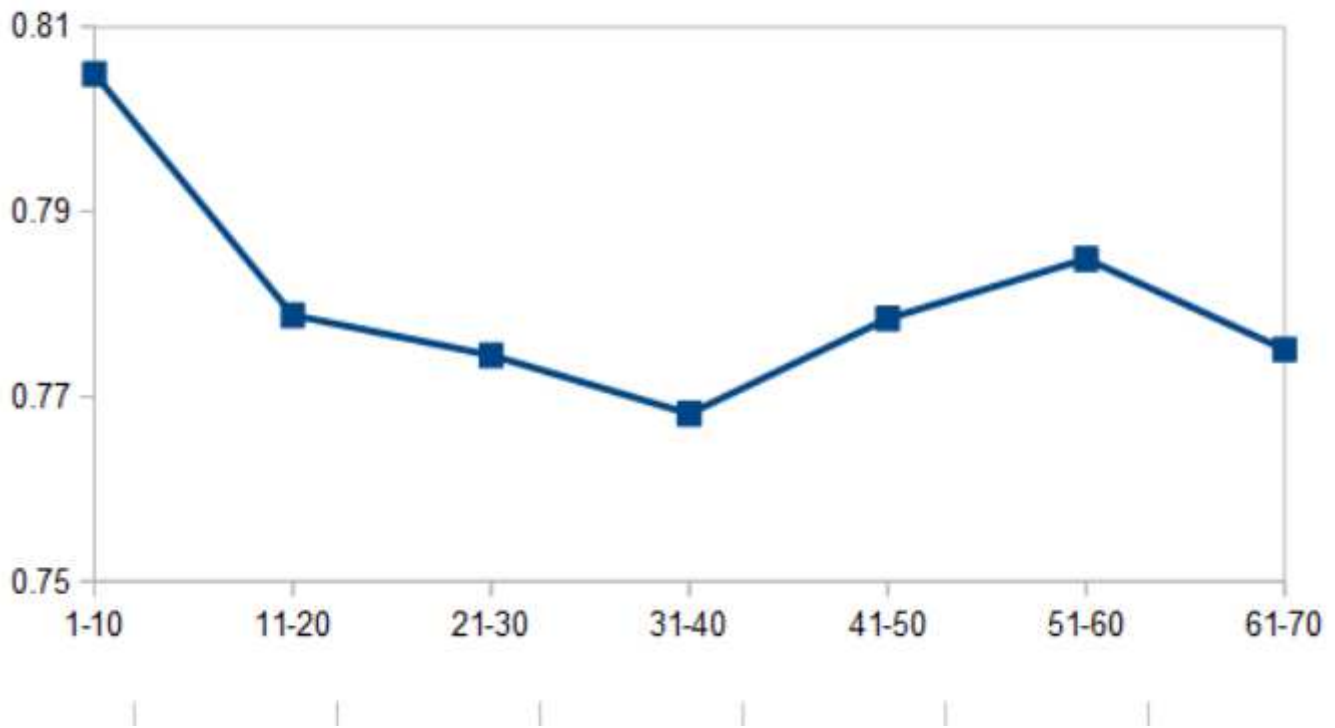


Figure 2

F1 scores of FSL models trained using each 10-feature groups arranged according to mRMR scores.

F1 Score of Rearranged Groups

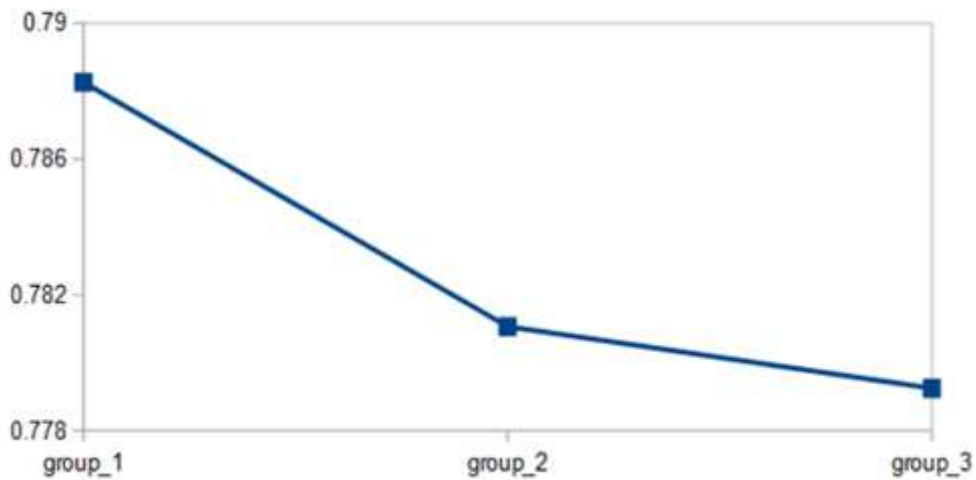


Figure 3

F1 score of FSL models inputting 20 features from three groups of mRMR intervals.

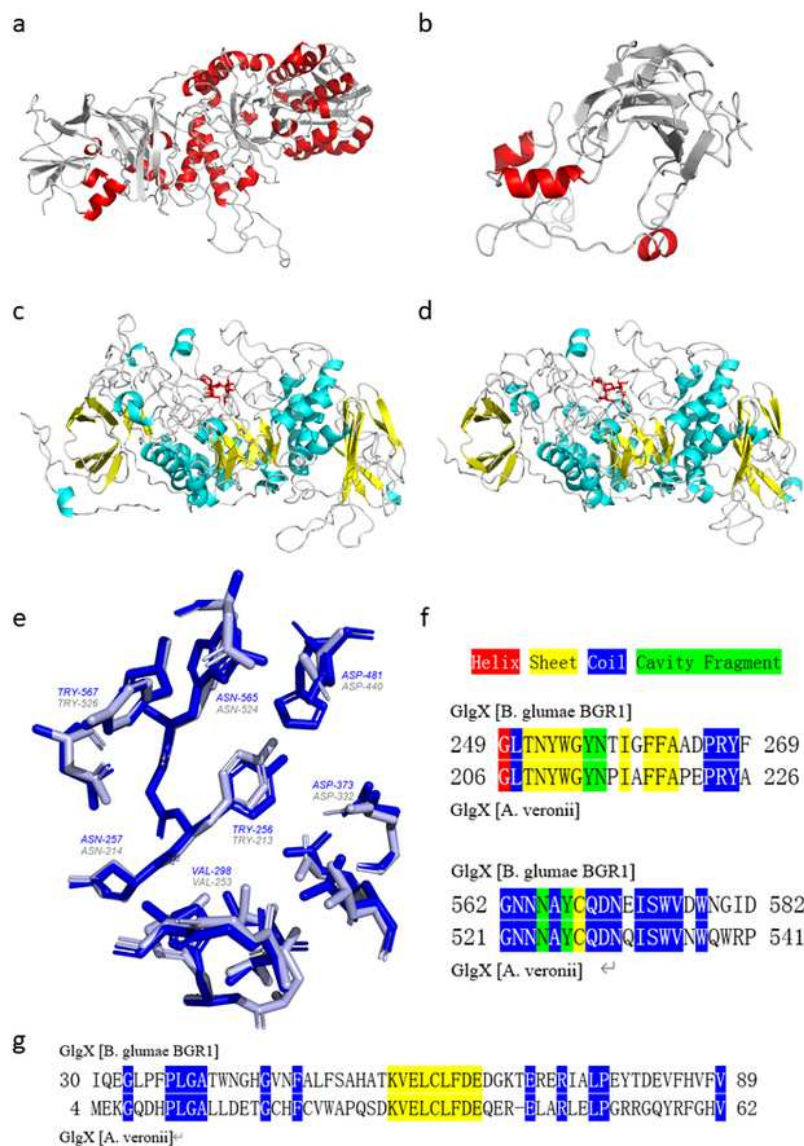


Figure 4

Structural presentation of sample proteins. a) Structure of P9 with alpha helices in red. SSSH_01 = 12, SSSH_02 = 150. b) Structure of P6 with alpha helices in red, SSSH_01 = 1, SSSH_02 = 35. c) Structure of anchor protein P5 with alpha helices in blue and beta sheets in yellow; binding ligand1-kestose is shown in red. d) Structure of protein P19 with alpha helices in blue and beta sheets in yellow; binding ligand1-kestose is shown in red. e) Structural alignment of cavity fragments of P5 (blue) and P19 (gray) using

align function in pymol with alignment RMSD = 0.389; residue labels are shown in corresponding colors. f) Sample AAS alignment of P5 (GlgX [*B. glumae*]) and P19 (GlgX [*A. veronii*]) around the two cavity fragments; aligned AA of a given secondary structure or belongs to a cavity fragment is shown in the color scheme. g) Sample AAS alignment of non-cavity parts of P5 and P19, with AA of a given secondary structure shown in the color scheme.

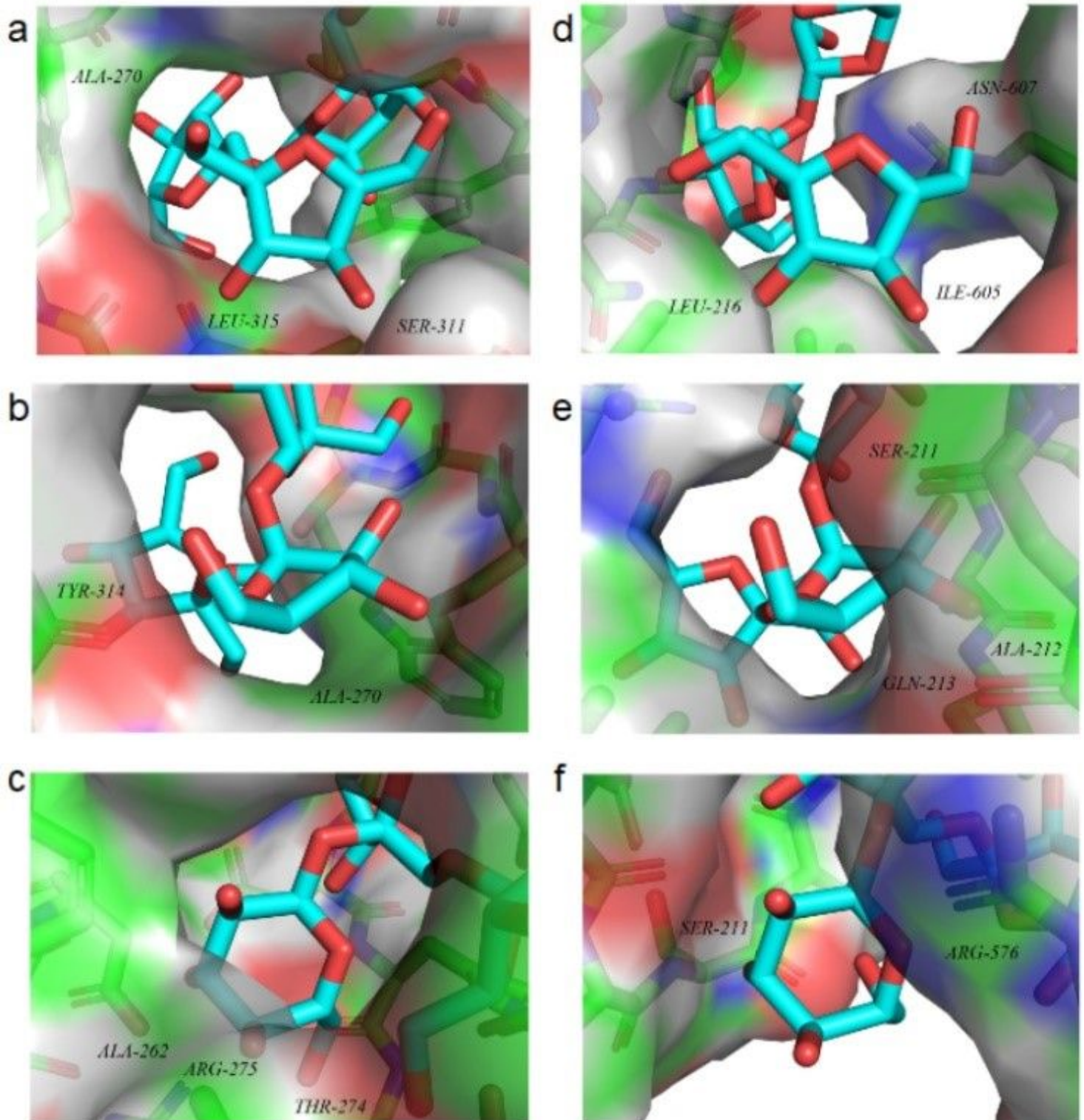


Figure 5

Predicted interaction of two CAZyme CBMs (P10 and P52) with 1-kestose. P10 and P52 are highly correlated according to group 1 features, while having no significant similarity in sequence. Similar amino acid residues are found in the radius of the interaction of both P10 and P52, including SER, ALA, ARG, and LEU. This suggests that the feature matrix may obtain inherent biological meaning. a) Interaction model of P10 viewing from beta-D-fructofuranose. b) Interaction model of P10 viewing from beta-D-fructofuranosyl residue. c) Interaction model of P10 viewing from alpha-D-glucopyranosyl residue. d) Interaction model of P52 viewing from beta-D-fructofuranose. e) Interaction model of P52 viewing from beta-D-fructofuranosyl residue. f) Interaction model of P52 viewing from alpha-D-glucopyranosyl residue.

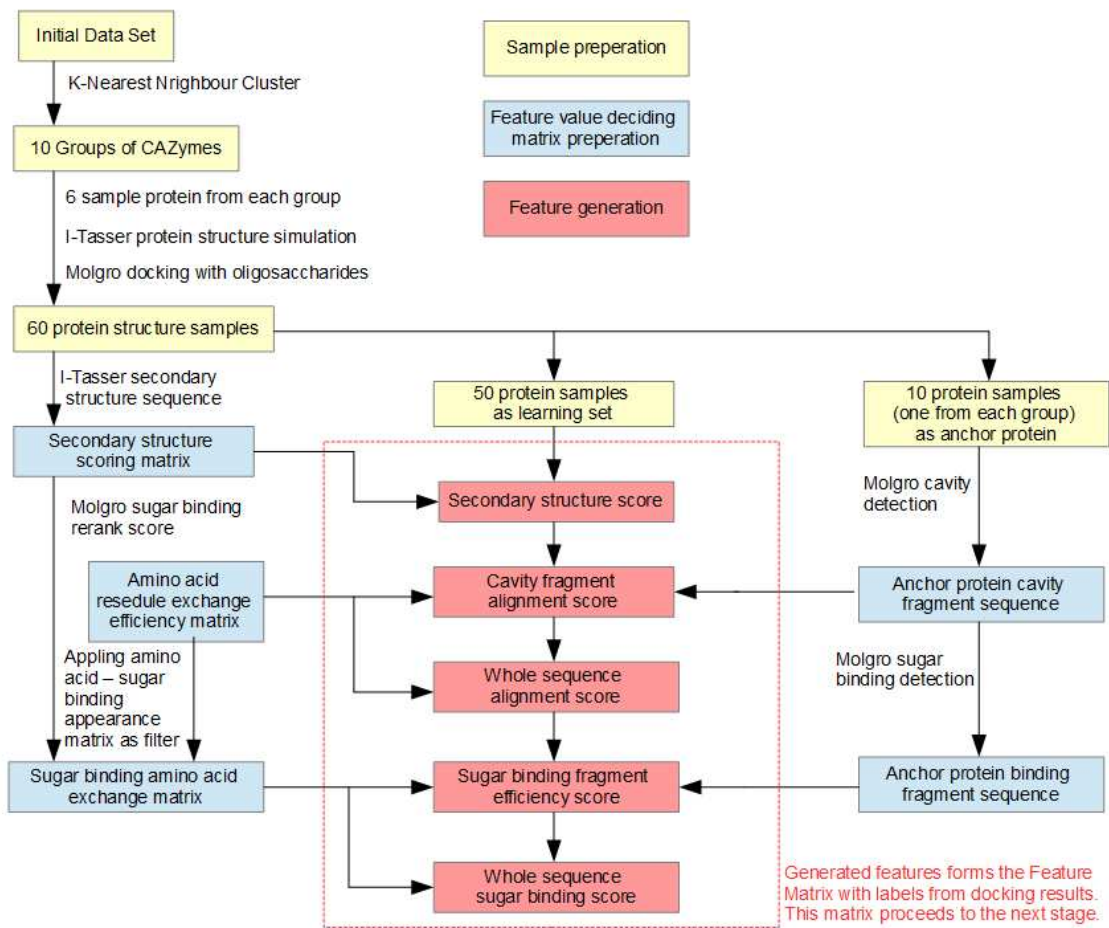


Figure 6

Flow chart for feature generation procedure. Arrows show derivation.

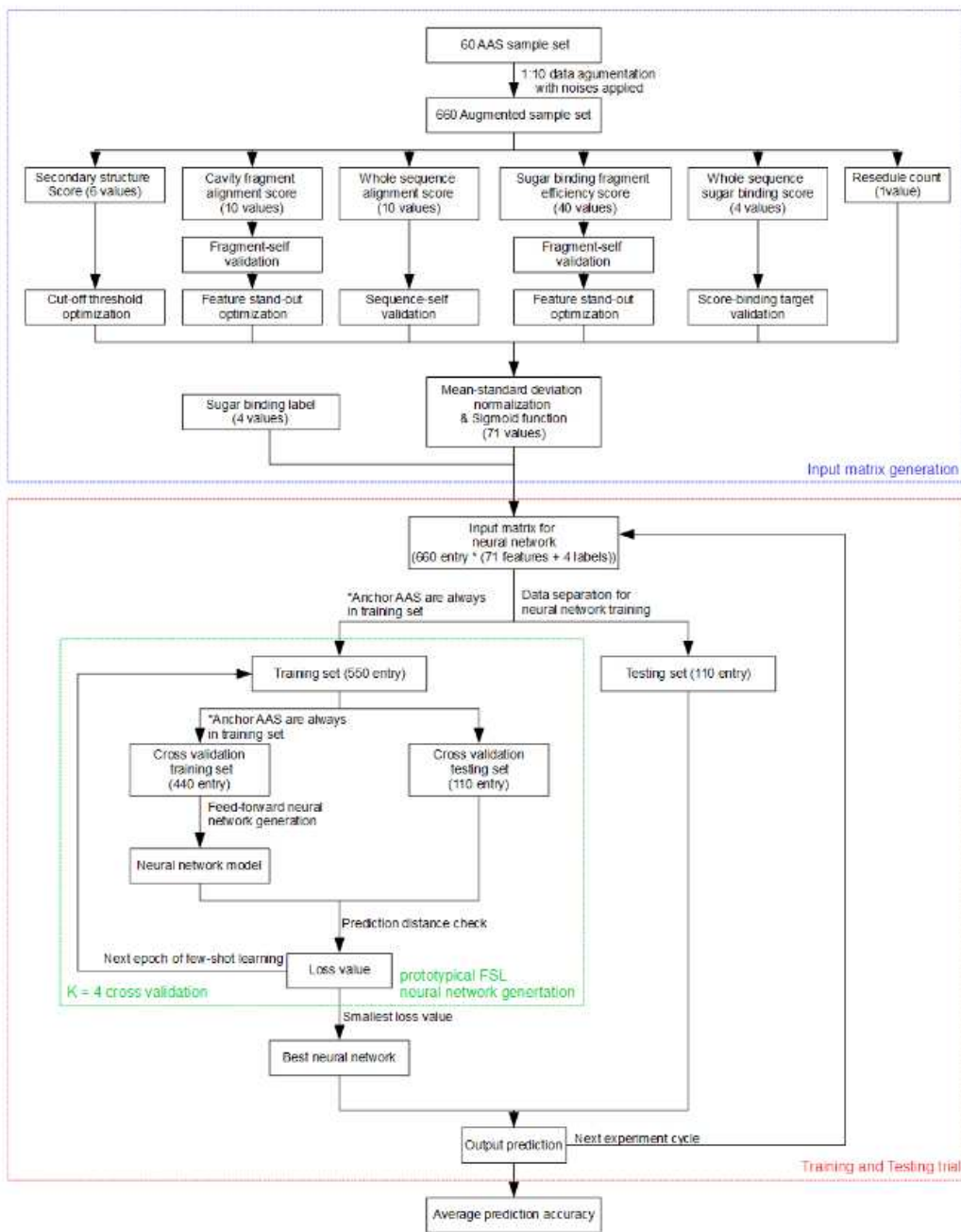


Figure 7

Flow chart of feature matrix generation and pFSLNN process

Algorithm 1

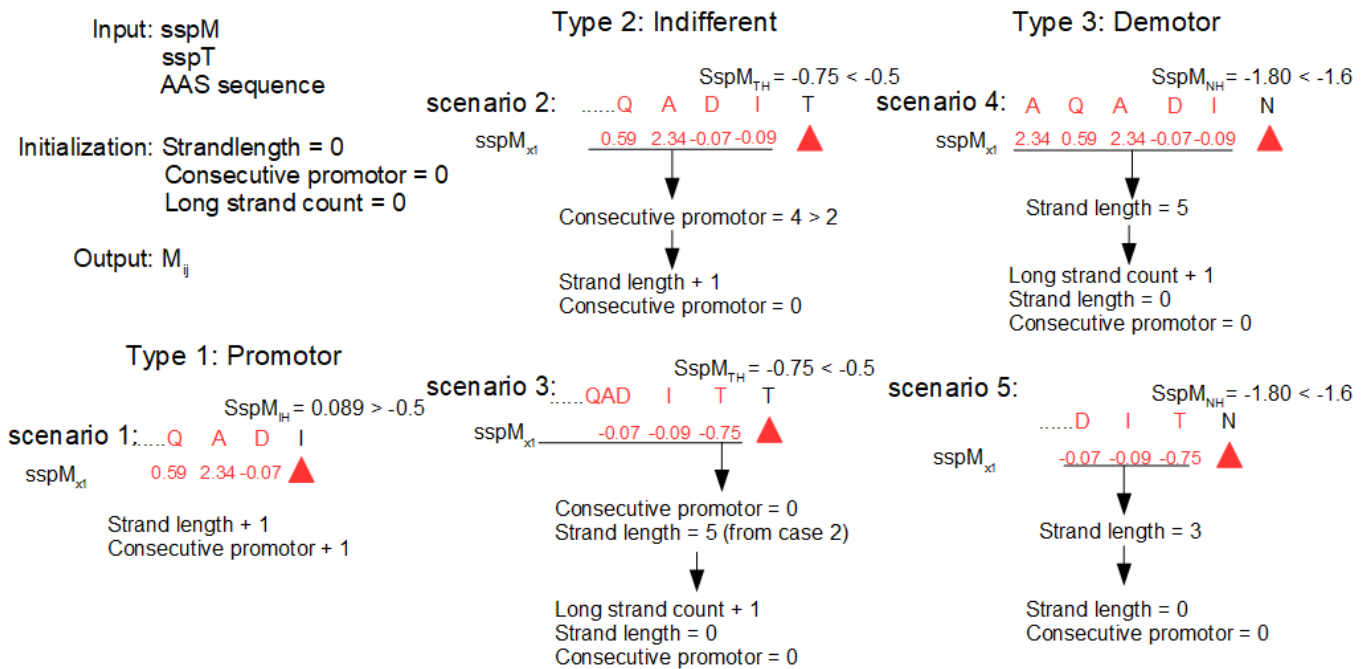


Figure 8

This figure illustrates the bases of the algorithm estimating the number of residuals of each general secondary structure type. Without loss of generality, α -Helix (H) is shown as an example. A cursor scans over the input AAS to find the first H promoting AA. A secondary structure strand of H begins. When the strand length is more than or equal to 5, the strand is considered a long strand. As the cursor proceeds, 3 types of successor (promoting, demoting, indifferent) AA can result in 5 cases. Initial $ssp_{TH} = (-0.5, -1.6, 2)$.

Supplementary Files

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

- [Supplimental1.bmp](#)
- [Supplimental2.png](#)