# Electronic Supplementary Information for "Data-Driven Models for Predicting Intrinsically Disordered Protein Polymer Physics Directly from Composition or Sequence"

Tzu-Hsuan Chao,<sup>†</sup> Shiv Rekhi,<sup>‡</sup> Jeetain Mittal,<sup>‡</sup> and Daniel P. Tabor<sup>\*,†</sup>

 †Department of Chemistry, Texas A&M University, College Station, TX 77843 USA
 ‡Department of Chemical Engineering, Texas A&M University, College Station, TX 77843 USA

E-mail: daniel\_tabor@tamu.edu

### Contents

$\mathbf{S1}$	Link to Repository Containing Code and Raw Data	$\mathbf{S3}$
S2	Simulation Snapshots	$\mathbf{S3}$
$\mathbf{S3}$	Illustrating Each Encoding Method on One Example Sequence	$\mathbf{S6}$
$\mathbf{S4}$	One-Hot Encoding Linear Regression Learning Curves	$\mathbf{S7}$
$\mathbf{S5}$	Learning Curves	$\mathbf{S8}$
$\mathbf{S6}$	2D Convolutional Neural Network	S12

$\mathbf{S7}$	Bag of Amino Acid Representation: Testing the Exponent in the Rep-	
	resentation	$\mathbf{S13}$
<b>S</b> 8	Count Encoding Extrapolation Test	S14
<b>S</b> 9	Prediction on Experimental Data	S15
<b>S10</b>	Reverse Sequences Testing	<b>S16</b>
$\mathbf{S11}$	Symmetric BAAI	S17
S12	Shuffled Test	<b>S18</b>
S13	Shuffle Test-Count Encoding	<b>S19</b>
$\mathbf{S14}$	Simulation Details	S20
$\mathbf{S15}$	Coarse-grained Model Parameters	S21
$\mathbf{S16}$	Best Parameters for Classical Regression Models	S22
$\mathbf{S17}$	Artificial Neural Network Architecture	S22
<b>S18</b>	Comparison of BAAI Representation to Augmented Fingerprints and	
	Performance on Alternative IDP Datasets	S24
<b>S19</b>	Analysis of Temperature-Dependent Model for Temperature Subsets	S28
Refe	erences	S31

## S1 Link to Repository Containing Code and Raw Data

The raw data, code, and results for this paper can be found at the following Github repository: https://github.com/Tabor-Research-Group/Predicting-Disordered-Polymer-Folding-Behavior-Directly-from-Sequences

### S2 Simulation Snapshots



Figure S1



Figure S2





Figure S3

## S3 Illustrating Each Encoding Method on One Example Sequence



Figure S4: Five encoding methods: categorical featurization (left) and physically-motivated featurization (right).

#### S4 One-Hot Encoding Linear Regression Learning Curves

Figures S5a and S5b contain the learning curves for the one-hot encoding linear regression models. The learning curve of one-hot encoding indicates that the model overfits for these sparse features, as the validation loss increases at the last three points as the training set size increases.



Figure S5: One-Hot Encoding Linear Regression Learning Curve

## S5 Learning Curves



Figure S6: Count Encoding Linear Ridge Regression Learning Curve



Figure S7: Count Encoding Kernel Ridge Regression Learning Curve



Figure S8: Count Encoding Support Vector Regression Learning Curve



Figure S9: Count Encoding Gaussian Process Regression Learning Curve



Figure S10: Ordinal Encoding Linear Ridge Regression Learning Curve



Figure S11: Ordinal Encoding Kernel Ridge Regression Learning Curve



Figure S12: Ordinal Encoding Support Vector Regression Learning Curve



Figure S13: Ordinal Encoding Gaussian Process Regression Learning Curve



Figure S14: One-Hot Encoding Kernel Ridge Regression Learning Curve



Figure S15: One-Hot Encoding Support Vector Regression Learning Curve



Figure S16: One-Hot Encoding Gaussian Process Regression Learning Curve

## S6 2D Convolutional Neural Network



Figure S17: Prediction using the Color Mapping representation and a2D convolutional neural network model.

## S7 Bag of Amino Acid Representation: Testing the Exponent in the Representation

Here, we tested the performance of a family of BAAI featurizations. The exponent that is used for each matrix element in the feature  $(\beta)$ :

$$BAAI = \sum_{i-j>1} (X_i - X_j)^{\beta}$$
(1)

was varied, similar to the approach employed in reference S1. For each of these formulations of the feature, we conducted the training and testing procedure as the model used in the main text. For these tests, we used support vector regression, as it was the best model for the baseline model. The first plot is the testing performance and the second plot is the extrapolation test using the same strategy as the main text but different exponent values used in each element.



Figure S18



Figure S19

#### S8 Count Encoding Extrapolation Test

We performed extrapolation tests (as described in section 3.4 of the main text) on count encoding. For this feature, we only tested the extrapolation performance with a support vector regression model, as the count encoding only contains 20 features. The performance is similar to BAAI model on the IDP-10260 dataset.



Figure S20

### **S9** Prediction on Experimental Data

We apply both CE and BAAI models on 42 experimental data points. Three plots of correlation between prediction and experimental values are provided.



Figure S21: Prediction on Experimental Data

#### S10 Reverse Sequences Testing



Figure S22

As the underlying coarse-grained simulations do not depend on the amino acid sequence, we can test the performance of the trained models on predicting the properties of "reversed sequences." By virtue of its construction, the count encoding feature will have the same predictions, but the directional BAAI representation will have different representations for reversed sequences.

We tried 1) only reversing the training sequences or 2) only reversing the testing sequences. The resulting  $R^2$  (as a function of  $\beta$ ) shows that both cases give similar performance. This indicates that the BAAI feature can handle data that has properties that are symmetric upon reversal of the sequence while the generated matrix is transposed. This test was to ensure that this model, with this representation, doesn't "over-learn" from the intrinsically disordered proteins.

#### S11 Symmetric BAAI

Since the coarse-grained model that was used to generate the training data is "symmetric" (it would give identical results if the sequence were reversed), we built a symmetric representation, which encodes both forward (*e.g.*,  $C \rightarrow A$ ) and reverse (*e.g.*,  $C \leftarrow A$ ) interactions together in a single interaction. To make a non-redundant version of symmetric BAAI representation, we tested training the model with only 210 features (only the "upper triangle" of the BAAI matrix). The testing performance is similar to the original representation.



Figure S23: Performance ( $\mathbb{R}^2$  of the test set) of the symmetric BAAI representation (where are all interactions are symmetric, instead of proceeding along the chain) for the SVR models as a function of  $\beta$ .

#### S12 Shuffled Test

Here, we shuffle the  $R_g$  of a specific portion of training data and test using the correct  $R_g$ . The  $R^2$  values between shuffled and non-shuffled training data are plotted to see how robust the representation is to training set errors, which might be seen in larger, future datasets. The testing performance can still achieve over 0.75  $R^2$  when 40% of the training data is shuffled.



Figure S24: Performance of training and test set scores for BAAI representation (with SVR model) when the model is given "erroneous" random sequence data as a fraction of its training set.

#### **Count Encoding Shuffled Test** 0.8 Training score 0.6 **Testing score** <sub>2</sub>~ 0.4 0.2 0.0 0.2 0.4 1.0 0.0 0.6 0.8 Shuffled Ratio

### S13 Shuffle Test-Count Encoding



The shuffle test is also performed on count encoding. The training score goes down (as it should), while the testing score can still maintain an  $\mathbb{R}^2$  of 0.8 when 60% of the data is shuffled.

#### S14 Simulation Details

All the simulations were performed using LAMMPS. A Langevin thermostat with a friction coefficient of  $1 \text{ ps}^{-1}$  was used to for each simulation. Each simulation consisted of a 500 ns run. The first 100 ns was used for the system to reach equilibrium. Seven temperatures (270, 300, 330, 360, 390, 420, 450K) and 10 fs timestep were used for all sequences. Bonded, electrostatic and short-range pairwise interactions defined by amino acid hydropathy are included in the simulations. <sup>S2</sup> Bonded interactions are modeled using a harmonic potential with a spring constant of 10 kcal/mol Å<sup>2</sup> and a bond length of 3.8 Å. Electrostatic interactions are described using Debye-Huckel electrostatic screening:

$$E_{ij}(r) = \frac{q_i q_j}{4\pi\epsilon_0 Dr} \exp\left(-r/\kappa\right) \tag{2}$$

where  $\kappa$  is the Debye screening length and D is the dielectric constant of the solvent, 80 for water. The short-range pairwise interaction is described by Ashbaugh-Hatch functional form,

$$\Phi(r) = \begin{cases} \Phi_{LJ} + (1 - \lambda)\epsilon, & \text{if } r \le 2^{\frac{1}{6}}\sigma \\ \lambda \Phi_{LJ}, & \text{otherwise} \end{cases}$$
(3)

where  $\Phi_{LJ}$  is the standard Lennard-Jones potential. The lambda value is calculated by the average of the hydropathy values of two amino acids (Table S1).

The scaling exponent  $\nu$  is derived using the formulation of Zheng *et al.*<sup>S1</sup>

$$R_g = \sqrt{\frac{\gamma \left(\gamma + 1\right)}{2 \left(\gamma + 2\nu\right) \left(\gamma + 2\nu + 1\right)}} b N^{\nu} \tag{4}$$

where  $\gamma = 1.1615$ , b = 0.55 and N is the number of peptide bonds.

## S15 Coarse-grained Model Parameters

Amino Acid	Charge	$\sigma(\text{\AA})$	$\lambda$
A	0	5.04	0.730
$\mathbf{C}$	0	5.48	0.595
D	-1	5.58	0.378
Ε	-1	5.92	0.459
F	0	6.36	1.000
G	0	4.50	0.649
Н	0.5	6.08	0.514
Ι	0	6.18	0.973
Κ	1	6.36	0.514
L	0	6.18	0.973
М	0	6.18	0.838
Ν	0	5.68	0.432
Р	0	5.56	1.000
Q	0	6.02	0.514
R	1	6.56	0.000
$\mathbf{S}$	0	5.18	0.595
Т	0	5.62	0.676
V	0	5.86	0.892
W	0	6.78	0.946
Υ	0	6.46	0.865

Table S1: CG model parameters

#### S16 Best Parameters for Classical Regression Models

#### **Count Encoding**:

Linear Ridge Regression: alpha: 0.6 Kernel Ridge Regression: alpha: 0.1, gamma: 1 Support Vector Regression: C: 100, gamma: 0.1 epsilon; 0.1 Gaussian Process Regression: alpha:0.05, length scale:1e5, length bounds:1e-20  $\sim$  1e20

#### **Ordinal Encoding**:

Linear Ridge Regression: alpha: 350 Kernel Ridge Regression: alpha: 0.6, gamma: 0.01 Support Vector Regression: C: 1, gamma: 0.1 epsilon; 0.1 Gaussian Process Regression: alpha:0.7, length scale:1e5, length bounds:1e-20~ 1e20

#### **One-Hot Encoding**:

Linear Ridge Regression: alpha: 120 Kernel Ridge Regression: alpha: 0.1, gamma: 0.001 Support Vector Regression: C: 10, gamma: 0.001 epsilon; 0.001 Gaussian Process Regression: alpha:1.5, length scale:1e5, length bounds:1e-20~ 1e20

#### **BAAI** representation:

Support Vector Regression: C:10 gamma:0.1 epsilon:0.1

#### S17 Artificial Neural Network Architecture

#### Feed Forward Neural Network(section 3.5)

Input layer: dimension:400, activation function: swish, kernel initializer: glorot uniform

One hidden layer: dimension:400, activation function: swish, kernel initializer: glorot uniform

Output layer: dimension:2, activation function: linear, kernel initializer: glorot uniform

#### Conv1D neural network(section 3.2)

One conv1d layer: filter:8, kernel size:3, stride:1, activation function: relu, kernel initializer: lecun normal

One hidden layer: dimension:10, no activation function, kernel initializer: lecun normal, dropout rate:0.2

output layer: dimension:1

#### Conv2D neural network(section 3.2)

One conv2d layer: filter:8, kernel size:(3,3), stride:(1,1), activation function: relu, kernel initializer: lecun normal

One hidden layer: dimension:10, no activation function, kernel initializer: lecun normal, dropout rate:0.2

output layer: dimension:1

## S18 Comparison of BAAI Representation to Augmented Fingerprints and Performance on Alternative IDP Datasets

Here, we test the performance of a set of the BAAI representations (specific details indicated in the figures) on the set of 2585 intrinsically disordered proteins given in Reference S3. We compare our results on both datasets to the augmented fingerprint approach described in the same reference, on both the randomly generated 10260 dataset in this paper and the 2585 dataset. The overall results are summarized in Table S2, with details shown in the text and figures that follow.

Table S2: Test set  $\mathbf{R}^2$  for BAAI representations and Augmented Fingerprint Representations (SVR Models)

	Count Encoding	Augmented Fingerprint	BAAI	$\begin{array}{c} \mathbf{BAAI} \\ + \mathbf{CL} \end{array}$	Symmetric BAAI	$\begin{array}{c} {\rm Symmetric} \\ {\rm BAAI} + {\rm CL} \end{array}$	
IDP-10260 (This work)	0.942	0.961	0.936	0.937	0.937	0.937	
IDP-2585 set $^{S3}$	0.964	0.983	0.953	0.965	0.948	0.969	$(\beta = -0.5)$
			0.967	0.967	0.974	0.974	$(\beta = -2.0)$

Although the BAAI + SVR model showed some ability to extrapolate to longer sequence lengths in the 10260 dataset, the 2585 dataset contains much longer sequences (similar in length to the sequences shown in Figure S21). In addition, the 2585 simulation dataset was obtained with the temperature-corrected coarse-grained model, whereas we used a model where the potential is not temperature-dependent for the 10260 dataset. Similar to our results in Section S9, this means that the model was trained on a just the IDP-10260 set does not extrapolate well to higher values of  $R_g$ . However, we are able to improve the BAAI + SVR model with retraining, which we show both in Table S2 above and the figures below.



Figure S26: Direct prediction or  $R_g$  based on BAAI representation (SVR model) trained on our data set. The R<sup>2</sup>, MSE, and RMSE are indicated on the inset of the plot, indicating that retraining is necessary.



Figure S27: Retraining on the IDP-2585 set. A test on the exponent  $\beta$  was performed to evaluate the sensitivity of the model.

Though the BAAI models do not reach the same accuracy as the augmented fingerprint models, we expect that some of this behavior could be explained by the size of the dataset. Since all BAAI representations have at least 210 features, they may need a larger volume of training data to show a payoff in improved performance, though the performance is relatively good, especially for the symmetric BAAI with  $\beta = -2.0$ , where only a few points could make difference in the model's topline R<sup>2</sup> performance. In addition, as mentioned in the main text, certain sequences may need to be explicitly included to challenge the models more than randomly generated (or DisProt-sourced) sequences.



Figure S28: Extrapolation test using different  $\beta$  values in BAAI on the 2585 dataset. The R<sup>2</sup> values and percent error are plotted to examine the extrapolation ability of  $\beta$  values of -0.5 or -2.0.



Figure S29: Performance of BAAI and BAAI + Chain Length using different  $\beta$  values.



Figure S30: Performance of symmetric BAAI and symmetric BAAI + Chain Length (CL) using different  $\beta$  values.



Figure S31: Performance of count encoding and augmented fingerprint

## S19 Analysis of Temperature-Dependent Model for Temperature Subsets

The following figures show the results for the models presented in Figure 5 of the main text but broken down by performance in each temperature. We present the data in four different ways due to the changing and length scales for the  $R_g$  of the polymer as a function of temperature. First, we present the  $\nu$  and  $R_g$  scatter plots as shown in Figure 5 in the text. Then, we show the  $R^2$ , MAE, RMSE, and average percent error predictions as a function of temperature for both of these metrics. Overall, the models are relatively consistent across different temperatures. However, because  $R_g$  increases as a function of temperature, MAE and RMSE will increase in parallel. In the intermediate temperatures, we also see an increase in the MAE and  $\nu$  (which, although it should be length-independent, has previously been noted to be length-dependent for this dataset in the text). Finally, as all of the IDPs become, on average more "random walk" like at higher temperatures (note the clustering of the data above 360 K on  $\nu$  in Figure S32), it becomes more difficult to differentiate the polymers from each other, and thus, R<sup>2</sup> decreases, even though the percent error in the prediction remains the same.



Figure S32: Performance of temperature-incorporated model for each temperature for  $\nu$ .



Figure S33: Performance of temperature-incorporated model for each temperature for  $\nu$ .



Figure S34: Performance of temperature-incorporated model for each temperature for  $R_g$ 



Figure S35: Performance of temperature-incorporated model for each temperature for  $R_g$ 

#### References

- [S1] Zheng, W.; Dignon, G.; Brown, M.; Kim, Y. C.; Mittal, J. Hydropathy patterning complements charge patterning to describe conformational preferences of disordered proteins. J. Phys. Chem. Lett. 2020, 11, 3408–3415.
- [S2] Dignon, G. L.; Zheng, W.; Kim, Y. C.; Best, R. B.; Mittal, J. Sequence determinants of protein phase behavior from a coarse-grained model. *PLOS Comput. Biol.* 2018, 14, e1005941.
- [S3] Patel, R. A.; Borca, C. H.; Webb, M. A. Featurization strategies for polymer sequence or composition design by machine learning. *Mol. Syst. Des. Eng.* 2022, 7, 661–676.