

Supporting Information for Publication: ProtWave-VAE: Integrating Autoregressive Sampling with Latent-Based Inference for Data-Driven Protein Design

Niksa Praljak¹, Xinran Lian², Rama Ranganathan^{3,4}, and Andrew L. Ferguson⁴

¹Graduate Program in Biophysical Sciences, University of Chicago, Chicago IL, 60637,
USA

²Department of Chemistry, University of Chicago, Chicago, IL, 60637, USA

³Center for Physics of Evolving Systems and Department of Biochemistry and
Molecular Biology, University of Chicago, Chicago, IL, 60637, USA

⁴Pritzker School of Molecular Engineering, University of Chicago, Chicago, IL, 60637,
USA

Email: andrewferguson@uchicago.edu

Contents

1	Data preparation and preprocessing	S4
1.1	Protein family data collection and preprocessing	S4
1.2	Fitness prediction benchmarking data collection and preprocessing	S4
1.3	Chorismate mutase data collection and preprocessing	S5
1.4	SH3 protein data collection and preprocessing	S5
2	Model architecture and hyperparameterization	S5
2.1	Gated dilated encoder	S5
2.2	WaveNet decoder with latent conditioning	S7
2.3	Discriminate top model over the latent space for semi-supervision	S10
2.4	Model architecture and hyperparameter optimization for protein family task	S11
2.5	Model architecture and hyperparameter optimization for benchmark fitness task	S12
2.6	Model architecture and hyperparameter optimization for Chorismate Mutase task	S13
2.7	Model architecture and hyperparameter optimization for SH3 design task	S13
3	Supplementary Methods	S14
3.1	Random mutagenesis of SH3 sequences	S14
4	Supplementary Results	S15
4.1	Pfam task: ColabFold structure prediction analysis of design sequences	S15
4.2	SH3 design task: experimental results	S18
4.3	Comparison on generative performance for protein families	S19

List of Supplementary Figures

S1	Random mutagenesis control for elevating function.	S16
S2	Pfam task: infer biological representations and predict structure of design sequences. . .	S18
S3	Validation of the high-throughput select-seq assay.	S18

List of Supplementary Tables

S1	Comparison on generative performance on protein family tasks.	S19
S2	Comparison on generative performance on N-terminus prompting.	S19

1 Data preparation and preprocessing

1.1 Protein family data collection and preprocessing

The datasets used in this study were obtained from Rivoire et al. [13] or from the PFAM database (release 27.0). Accession codes PF00071, PF00186, and PF13354 were used for G proteins, DHFR, and class A β -lactamases, respectively. For each protein family, a reference sequence/structure was selected, namely rat trypsin (PDB 3TGI) for S1A proteases, human Ras (PDBs 5P21 and 4Q21) for G proteins, E. coli DHFR (PDB 1RX2) for DHFR, and E. coli TEM1 β -lactamase (PDB 1FQG).

The number of sequences in each dataset was as follows: DHFR (1759 sequences), G-protein (4974 sequences), β -lactamase (814 sequences), and S1A proteases (1344 sequences). To prepare the sequences for analysis, they were converted into one-hot encoded tensors. In this encoding scheme, each amino acid label corresponds to an index of the one-hot encoded vector along each position of the sequence. The size of the tensor was defined as the maximum length protein homolog in the family. Sequences with less than the maximum sequence length were padded with mask tokens, which have their own one-hot encoded index.

The DHFR, G-protein, beta-lactamase, and S1A proteases datasets were then converted into 2D tensors with sequence lengths of 150, 158, 199, and 205, respectively. The one-hot encoded vector used in the encoding scheme had a length of 21. The datasets can be found here: https://github.com/PrajakReps/ProtWaveVAE/tree/main/Pfam_analysis/data/protein_families.

1.2 Fitness prediction benchmarking data collection and preprocessing

The datasets used in this study include the Fitness Landscape Inference for Proteins (FLIP) and the Tasks Assessing Protein Embedding (TAPE). The FLIP datasets consist of two benchmarking tasks (AAV and GB1) and were obtained from FLIP benchmarks [4]. The TAPE datasets also include two benchmarking tasks and were obtained from TAPE benchmarks [10].

The AAV benchmarking task involves a mutational screening landscape of VP-1 AAV proteins [2, 21] (UniProt Accession P03135). Sequences were mutagenized along a 28-amino acid window from position 561 to 588 of VP-1, and the resulting variants were tested for fitness, with between 1 and 39 mutations. The VP-1 AAV benchmark task has seven dataset splits: (1) sampled, (2) sampled-designed, (3) design-sampled, (4) train on single mutants, (5) train on single and double mutants, (6) train on mutants with up to seven changes, and (7) train on low fitness, test on high.

The GB1 benchmarking task involves an exhaustive, combinatorial, and highly epistatic mutational landscape, ranging from variants with fewer mutations to predict activity of variants with more mutations [20]. The number of fitness measurements was 149,361 out of 160,000 possible combinations of mutations at four positions. The GB1 benchmark task also has seven dataset splits: (1) sampled, (2) train on single mutants, (3) train on single and double mutants, (4) train on single, double, and triple mutants, and (5) train on low fitness, test on high.

The TAPE benchmarking tasks include a green fluorescence protein (GFP) and stability measurements of candidate proteins. For the fluorescence regression predictions, a Deep Mutational Scanning (DMS) approach was used to characterize the local genotype-to-phenotype mapping of a single protein [16]. The training and validation sets include mutants with only Hamming distance 3 from the original protein, while the testing set includes mutants with Hamming distances 4-15. In contrast, the stability task [14] has a training and validation set containing proteins from four rounds of experimental data measuring candidate proteins, while the testing set consists of mutant variants of seventeen 1-Hamming distance neighborhoods. The FLIP datasets can be found here: <https://benchmark.protein.properties/>, while the TAPE datasets can be found here: <https://github.com/songlab-cal/tape#data>.

1.3 Chorismate mutase data collection and preprocessing

The dataset used in this study was obtained from Russ et al. [15]. For each protein family, four CM atomic structures (PDB entries 1ECM, 2D8E, 3NVT, 1YBZ) were selected as reference sequences/structures to build the multiple sequence alignment. The 1ECM PDB corresponds to the wild-type Chorismate mutase *E. coli*. The training dataset consists of natural homolog sequences, with 1130 sequences, while the testing dataset consists of synthetic natural homolog sequences, with 1618 sequences. Each protein sequence has a corresponding fitness value quantified by the normalized relative enrichment (r.e.).

To prepare the sequences for analysis, they were converted into one-hot encoded tensors, with each amino acid label corresponding to an index of the one-hot encoded vector along each position of the sequence. The size of the tensor was defined as the maximum length protein homolog in the family. Sequences with less than the maximum sequence length were padded with mask tokens, which have their own one-hot encoded index.

The CM dataset was then converted into a 2D tensor with a sequence length of 96. The one-hot encoded vector used in the encoding scheme had a length of 21. The dataset can be found here: https://github.com/PraljakReps/ProtWaveVAE/tree/main/Pfam_analysis/data/protein_families.

1.4 SH3 protein data collection and preprocessing

The datasets used in this study were obtained from [7], and consist of SRC homolog 3 (SH3) homologs. The dataset contains various natural paralogs, Sho1 orthologs, and synthetic domains with fitness measurement scores, quantified by the relative enrichment (r.e.). The dataset size is 17,218 sequences, which includes both natural and synthetic sequences.

To prepare the sequences for analysis, the input data tensor was set to a length of 82, corresponding to the maximum sequence length within the dataset. Shorter sequences were padded with mask tokens at the ends. Each sequence was then converted into 2D tensors using a one-hot encoded transformation, such that each amino acid and masked padded token corresponded to a one-hot encoded index.

The number of sequences with fitness measurements is 14,768. During stratified cross-validation, we set $k = 5$ and split the dataset into five partition bins, such that the training and validation set sizes are 80% and 20%, respectively. The dataset can be found here: https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/data.

2 Model architecture and hyperparameterization

2.1 Gated dilated encoder

The primary hyperparameters for the gated dilated convolutional encoder network $q_\phi(z|x)$ include encoder depth, input channel depth for the initial convolutional layer, hidden channel depth for the subsequent convolutional layers, kernel size for each convolutional layer, and the number of fully connected linear layers before transitioning the encoder outputs into the latent space. We increased the dilation window for the convolutional layer by 2^i , based on the encoder layer depth (where i is the depth index of the encoder), creating what is known as dilated convolution layers.

For each protein family task, the encoder depth was chosen to ensure that the original input sequence remained greater than length 0 after the convolutional layers reduced its length. Consequently, hyperparameter optimization for the encoder depth was not necessary, as it depended on the input sequence length `compute_max_enc_depth`. Following the approach used in WaveNet [9] and PixelCNN [19], we employed two sets of dilated convolutions: one for signal activations and another for gated activations.

Our nonlinear activation function was a linear gated activation function [5], which multiplied the outputs of the signal and gated convolution layers after applying a sigmoid function onto the gated

91 convolution outputs. This enabled the encoder to incorporate a memory gate similar to LSTMs and
92 GRUs [6, 3]. A batch normalization layer was applied after each gated activation operation. Once
93 the dilated convolution layers processed the hidden representations, signal and gate convolution layers
94 with kernel size 1 and channel depth 1 were applied to the outputs. Subsequently, a gated activation
95 operation and a batch normalization layer were added.

96 Next, we applied k fully connected layers, using leaky ReLU activation functions with an α hyper-
97 parameter of 0.1. The final outputs were then passed through a variational mean linear layer and a
98 variational variance module, which comprised a linear layer followed by a softplus activation function.
99 Ultimately, the final variational mean and variational variance outputs were combined using the repara-
100 rameterization trick (`reparam.trick`). This involved sampling Gaussian noise from a normal distribu-
101 tion ϵ , multiplying it with the square root of the variational variance, and adding the variational mean
102 to produce the latent vector. The forward pass for the encoder network is shown in pseudocode (Algo-
103 rithm 1). The source code and class object of the encoder network is in PyTorch and can be found in the
104 following link: [https://github.com/PraljakReps/ProtWaveVAE/blob/main/SH3_design_project/
105 source/model_components.py](https://github.com/PraljakReps/ProtWaveVAE/blob/main/SH3_design_project/source/model_components.py).

Algorithm 1 Gated dilated convolution encoder network $q_\phi(z|x)$

```
Input: {
x, encoder_depth, num_fc, sigmoid, lrelu, initial_conv_layer,
signal_dilated_conv_layers, gate_dilated_conv_layers,
batch_norm_layers, final_signal_conv_layers, final_gate_conv_layers,
fc_layers, mean_z_layers, var_z_layers
}
// *Description of inputs.*
// x --> input sequence data
// encoder_depth --> number of dilated convolution layers
// num_fc --> number of fully connected linear layers
// sigmoid --> sigmoidal activation function
// lrelu --> leaky ReLU activation function
// initial_conv_layer --> first convolution layer with kernel size 1
// signal_dilated_conv_layers --> list of layers along the signal network path,
// consisting of dilated convolutions with increasing dilations by factor of 2
// gate_dilated_conv_layers --> list of layers along the gated network path,
// consisting of dilated convolutions with increasing dilations by factor of 2
// batch_norm_layers --> list of batch normalization layers
// batch_norm_layers --> list of batch normalization layers
// final_signal_conv_layers --> single convolution layer with kernel size 1 and depth
// channel 1
// final_gate_conv_layers --> single convolution layer with kernel size 1 and depth
// channel 1
// final_gate_conv_layers --> single convolution layer with kernel size 1 and depth
// channel 1
// fc_layers --> fully connected linear layers
// mean_z_layers --> fully connected linear layers mapping encoder outputs to latent
// mean
// var_z_layers --> fully connected linear layers mapping encoder outputs to latent
// variance
1 batch_size, input_channels, max_protein_length = x.shape
2 encoder_depth = compute_max_enc_depth(
    max_protein_length=max_protein_length
) // Determine encoder depth based on protein sequence length
3 h = init_conv_layer(x)// 1x1 convolution operation
4 h = batch_norm[i](h)
5 for i ← 1 to encoder_depth
6 // signal and dilated convolution layer operation
7 h_signal=signal_dilated_conv_layers(dilation=2i-1)[i](h)
8 h_gate=gate_dilated_conv_layers(dilation=2i-1)[ii](h)
9 h = h_signal * sigmoid(h_gate)
10 h = batch_norm_layers[ii+1](h)
11 end for
// final single no-dilated convolution layer to map output depth channels to 1
12 h_signal = final_signal_conv_layer(h)
13 h_gate = final_gate_conv_layer(h)
// apply gated activation function
14 h_out = h_signal * sigmoid(h_gate)
15 h_out = final_gate_conv_layer(h_out)
16 for i ← 1 to num_fc
17 h_out = fc_layers[ii](h_out)
18 h_out = lrelu(h_out)
19 end for
20 z_mean, z_var = mean_z_layer(h_out), var_z_layer(h_out)
// infer latent variables using reparameterization trick
21 z = reparam_trick(z_mean, z_var) Output: z, z_mean, z_var
```

106 2.2 WaveNet decoder with latent conditioning

107 The WaveNet architecture is based on the work of van Oord et al. [9], utilizing dilated causal convo-
108 lution layers. Our approach differs in that we employ latent variables for conditional inference rather

109 than one-hot encoded class variables. In accordance with the class conditional WaveNet, we upsample
 110 our latent variables through a linear layer, enabling us to extend the latent vector to the sequence
 111 data dimension’s length. This is referred to global conditioning [9]. No activation function is applied
 112 during this linear upsampling for the latent variable, and the latent linear upscaling operation for
 113 global conditioning is the following:

$$Z = V_k^T z \quad (1)$$

114 where V_k is learnable linear project, z is the inferred latent low-dimensional variable, Z is the vector
 115 broadcast over the sequence position dimension. This operation is implemented by `latent_cond_net()`
 116 `nn.module`, and the source code is [https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd05](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L23)
 117 [4be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L23](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L23).

118 The WaveNet decoder has two primary components: the main module, which is composed of dilated
 119 causal convolutions, gated activation functions, residual connections, and skip connections; and the
 120 top head, which is inspired by van Oord et al.’s architecture [9] and features two convolution layers
 121 (non-dilated) with kernel size 1 and a ReLU activation function. The second component generates
 122 logits that is converted to the amino acid class distribution using a softmax function. The main
 123 hyperparameter optimization for the first component involves input channel depth, output channel
 124 depth, and the number of dilation rates (i.e., the WaveNet’s depth). The causal dilation convolutions’
 125 long-range reach is determined by the last hyperparameter. The top model’s hyperparameters include
 126 input channel depth, output channel depth, and hidden state depth. The output channel depth is set
 127 to the WaveNet’s input channel depth, which corresponds to the one-hot encoded vector’s length (20
 128 amino acids plus 1 padded token).

129 The WaveNet decoder processes the upsampled latent variable and input sequence. The first layer
 130 applies a convolution layer with kernel size 1 to the input sequence, creating an amino acid embedding
 131 layer $h = W_{k=1} x$, where $W_{k=1}$ is the 1×1 kernel with learnable weights. Subsequently, the causal
 132 dilated convolution is employed in two separate modules for signal and gated representations. At
 133 the same layer depth as the dilated convolutions, a vanilla convolution layer with kernel size 1 is
 134 also applied to the upsampled latent variable, mapping input depth 1 to output depth as per the
 135 hyperparameter C_{out} . The hidden representation outputs from the signal causal dilated convolution
 136 are then summed with the outputs from the signal latent convolution, and the outputs from the gated
 137 causal convolution are summed with those from the latent gated convolutions. A sigmoid function is
 138 applied as a gating function to the gated output representations, and these values are multiplied with
 139 the signal hidden representations. After completing this process for a given dilated causal convolution
 140 layer, a convolution layer corresponding to the skip module and another convolution layer for the
 141 residual module, both with kernel size 1, are applied. Thus, in terms of mathematical expression, the
 142 layer operations are the following:

$$h = (W_{a,k} * h + \tilde{W}_{a,k=1} * Z) \odot \sigma(W_{g,k} * h + \tilde{W}_{g,k=1} * Z) \quad (2)$$

143 where first term in the parentheses corresponds to linear activations a while the second term in the
 144 parentheses corresponds to the gated inputs to the sigmoidal function σ , acting as a gated function
 145 and outputting gated activations g . The learnable kernel weight for the sequence representations and
 146 latent variable representations are $W_{*,k}$ and $\tilde{W}_{*,k=1}$. The variables h and Z corresponds to the hidden
 147 representations of the input protein sequence and upsampled global latent conditioning.

148 The skip output representations are accumulated (i.e., summed) after each dilated causal convo-
 149 lution layer, are implemented with a linear convolution with kernel size 1, and independently applied
 150 relative the to residual layer; and thus the mathematical operation for this layer are the following:

$$h_{skip}^l = V_{l,k=1}^{skip} * h \quad (3)$$

$$h_{cumm-skip} += h_{skip}^l \quad (4)$$

151 where $h_{cumm-skip}$ is the cumulative skip representations of all the skip output representations from
 152 the l dilated causal convolution layers applied to the input sequence. The output representations from
 153 the residual module’s convolution layers are added to the embedded input sequence, which is then fed
 154 back into the next dilated causal convolution layer with an increased dilation scale by 2^i , where i is
 155 the depth index of the dilated causal convolution layer. The mathematical operations are then the
 156 following:

$$h_{res} = V_{l,k=1}^{res} * h \quad (5)$$

$$h = h + h_{res} \quad (6)$$

157 The top model exclusively utilizes the accumulated skip connections as input and processes these
 158 hidden representations using a ReLU function, a convolution layer with kernel size 1, a second ReLU
 159 activation function, and finally a second convolution layer with kernel size 1. The top model generates
 160 logits that can be converted into probabilities using a softmax function. Thus, the final sequence
 161 probabilities using the top model is following:

$$\begin{aligned} h_{top-model}^1 &= Conv(ReLU(h_{cumm-skip})) \\ h_{top-model}^2 &= Conv(ReLU(h_{top-model}^1)) \end{aligned} \quad (7)$$

$$p(x|z) = softmax(h_{top-model}^2) \quad (8)$$

163 where $ReLU$ is elementwise nonlinear activation function and $Conv$ is a linear learnable affine trans-
 164 formation with bias and kernel size 1. The softmax layer converts the network logits into probabilities
 165 using the following function $softmax(x_i) = \frac{exp(x_i)}{\sum_j exp(x_j)}$

166 The source code for the class objects, written in PyTorch, can be found via this [https://github.com
 167 /PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_
 168 project/source/wavenet_decoder.py#L23](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L23), with the pseudocode provided below. The pseudocode for
 169 the upsampling latent network in PyTorch, comprising a solitary linear layer, is displayed below.

170 The subsequent PyTorch pseudocode represents the WaveNet decoder, encompassing numerous
 171 causal dilation convolutions, latent convolution layers, a skip module with convolution layers, and a
 172 residual module with convolution layers (see [https://github.com/PraljakReps/ProtWaveVAE/blob/
 173 01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py
 174 #L98](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L98)).

175 The final PyTorch pseudocode pertains to the TopHead of the WaveNet decoder, which produces the
 176 ultimate logits before utilizing a softmax function to transform them into amino acid class probabilities.
 177 This component features two straightforward convolution layers and two ReLU functions, while solely
 178 accepting cumulative skip outputs as input (see [https://github.com/PraljakReps/ProtWaveVAE
 179 /blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder
 180 .py#L222](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L222)).

181 Forward pass of the latent conditional WaveNet decoder in psuedocode (see [https://github.com
 182 /PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design
 183 _project/source/wavenet_decoder.py#L303](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/wavenet_decoder.py#L303)).

Algorithm 2 Forward pass for latent conditioned WaveNet decoder $p_{\theta}(x|z)$

Input: $x, z, \text{softmax}, \text{latent_cond_net}, \text{WaveNet_head}, \text{Top_head}$

// *Description of inputs.*

// x --> input sequence data// z --> inferred latent variable data// `latent_cond_net` --> upsampling linear layer for the latent variables// `WaveNet_head` --> WaveNet causal dilated convolution component// `Top_head` --> Top model component, consisting of convolution layers and outputting amino acid logits22 `z_upsampled = latent_cond_net(z)`23 `tot_cum_skip = 0 // Initialize the cumulative skip representations`24 `h, h_cum_skip = WaveNet_head(x, z_upsampled)`25 `tot_cum_skip += h_cum_skip // Cumulate WaveNet's skip representations`26 `logits = Top_head(tot_cum_skip) // final operation to compute amino acid logits`27 `p(x|z) = softmax(logits)`**Output:** $p(x|z)$

2.3 Discriminate top model over the latent space for semi-supervision

In our semi-supervised learning implementation, we utilize a top model that samples from the latent variables and predicts fitness regression values y and/or classifies sequences as functional or nonfunctional. For regression tasks, we employ the mean-squared error as our loss objective. Conversely, for classification tasks, we use the binary cross-entropy loss, $l(y, \tilde{y})$, where the loss value for a given sample sequence n is defined as $l_n = y_n * \log \tilde{y}_n + (1 - y_n) * \log(1 - \tilde{y}_n)$. Here, \tilde{y}_n represents the classification prediction probability, while y_n denotes the ground truth label (0 or 1).

The model architecture remains consistent, regardless of whether the task is regression or classification. It consists of hidden modules, with each module containing the following layers:

$$h = \text{Dropout}(\text{SiLU}(\text{LayerNorm}(\text{Linear}(z)))) \quad (9)$$

Dropout is a widely used regularization layer for neural networks [18]. SiLU represents a nonlinearity function, LayerNorm is a layer normalization layer [1], and Linear refers to a simple fully connected linear layer. The primary hyperparameters include `num_layers`, which defines the number of hidden modules to stack, and `hidden_width`, which sets the hidden width size of the linear layers. The hidden output representations h are then fed into a final linear layer for either regression or classification predictions. The classification path concludes with a sigmoidal function that maps the values to a binary probability, while the regression path does not require any final nonlinear layer since the outputs are continuous values. The hidden module's architecture is denoted as `TopModel_layer`, while the overall neural network architecture $p_{\omega}(y|z)$ is illustrated in object `Decoder_re`. The model and hidden module's source code can be found at this https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/model_components.py.

The hidden blocks found along the top model discriminator $p_{\omega}(y|z)$, which consist of linear layer followed by a layer normalization layer [1], then followed by a SiLU nonlinearity and dropout regularization layer [18]. The subsequent pytorch pseudocode represents these hidden block modules called `TopModule_layer`.

The top model discriminator's hidden blocks, $p_{\omega}(y|z)$, consist of a linear layer followed by a layer normalization layer [1], a SiLU nonlinearity, and a dropout regulation layer [18]. The subsequent PyTorch pseudocode represents these hidden block modules, referred to as `TopModule_layer`.

We present the entire neural model for the top discriminate model $p_{\omega}(y|z)$. This architecture encompasses both regression and classification tasks. However, the network can be easily adapted to cater to individual discriminate tasks. The subsequent PyTorch pseudocode represents these hidden block modules, denoted as `Decoder_re`. The source code can be found <https://github.com/PraljakReps/ProtWave>

215 [VAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/model_comp](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/model_components.py#L279)
216 [onents.py#L279](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/model_components.py#L279).

217 We demonstrate the forward pass of the latent discriminative top model $p_{\omega}(y|z)$ in pseudocode
218 (refer to the [https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca](https://github.com/PraljakReps/ProtWaveVAE/blob/01cecd054be7ca77540b23d19aeea5a3b9012d79/SH3_design_project/source/model_components.py#L279)).

Algorithm 3 Forward pass for latent discriminative top model $p_{\omega}(y|z)$

Input: z , `Decoder_re`

`// *Description of inputs.*`

`// z --> inferred latent variable data`

`// Decoder_re --> discriminative regression and classification model`

218 `y_pred_R, y_pred_C = Decoder_re(z) // model predictions p(y|z)`

Output: `y_pred_R, y_pred_C`

219 2.4 Model architecture and hyperparameter optimization for protein fam- 220 ily task

221 The task of the protein family is to evaluate the ability of ProtWave-VAE to generate biologically
222 meaningful latent representations for a given protein family, without any supervision (i.e., unsuper-
223 vised). Additionally, since the model is generative, this task tests its ability to generate novel sequences
224 that are indistinguishable from the protein family in terms of tertiary structure, using known PDB
225 structure and ColabFold structure predictions. Therefore, the architecture of ProtWaveVAE for pro-
226 tein family tasks (specifically Gprotein, DHFR, S1A, and lactamase families) consists of the gated
227 dilated encoder $q_{\phi}(z|x)$ with the latent conditioned WaveNet decoder $p_{\theta}(x|z)$, since it is solely un-
228 supervised. The loss objective is the unsupervised loss, which includes the negative log-likelihood,
229 Kullback Leibler divergence between the posterior and prior latent distribution, and the max-mean
230 discrepancy between the aggregated posterior and prior latent distribution. Since each family has
231 different maximum sequence lengths, the encoder depth is determined based on the maximum possi-
232 ble depth, or the number of dilated convolutions that can be applied onto the input sequence before
233 compressing the sequence length below 0. The hyperparameters that were optimized are the latent
234 space dimension \mathbf{z} , channel depth of dilated convolution in the WaveNet decoder `whs`, channel depth
235 of the 1x1 convolution layers for WaveNet TopHead `hhs`, channel depth of the encoder convolution
236 layers `C.out`, number of fully connected layers along the encoder path `num_fc`, number of dilated
237 causal convolutions along the decoder path `ndr`, prefactor weight for the KL divergence loss term
238 `KL_weight`, prefactor weight for the negative log-likelihood term `NLL_weight`, and the prefactor weight
239 for the max-mean discrepancy term `MMD_weight`. These hyperparameters were optimized in order,
240 and the optimal values were determined by the optimal negative log-likelihood loss while maintain-
241 ing excellent max-mean discrepancy regularization loss on a random hold-out that consisted of 20%
242 of the original protein family dataset. The optimization was performed over the domains $\mathbf{z} \in [1, 20]$,
243 $\mathbf{whs} \in \{32, 64, 128, 256, 512\}$, $\mathbf{hhs} \in \{32, 64, 128, 256, 512\}$, $\mathbf{C.out} \in \{32, 64, 128, 256, 512\}$, $\mathbf{num_fc} \in [0, 4]$,
244 $\mathbf{ndr} \in [1, 10]$, $\mathbf{KL_weight} \in \{0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 0.99\}$, $\mathbf{NLL_weight} \in \{0.1, 0.5, 1.0, 5.0, 10.0, 50.0, 100.0\}$,
245 and $\mathbf{MMD_weight} \in \{1.0, 2.0, 5.0, 10.0, 20.0, 50.0, 100.0\}$. The training consisted of using an Adam opti-
246 mizer with a learning rate equal to 1×10^{-4} . The minibatch size was 256, and the number of epochs
247 was 500. The model was trained on unaligned sequences for each protein family.

248 The final optimized hyperparameter configurations for the G-protein, DHFR, S1A, and lactamase
249 families are presented below. For G-protein, the model configuration is $\mathbf{z}=6$, $\mathbf{whs}=128$, $\mathbf{hhs}=512$,
250 $\mathbf{C.out}=32$, $\mathbf{num_fc}=3$, $\mathbf{ndr}=5$, $\mathbf{KL_weight}=0.99$, $\mathbf{NLL_weight}=1.0$, $\mathbf{MMD_weight}=5$. For DHFR, the
251 model configuration is $\mathbf{z}=6$, $\mathbf{whs}=128$, $\mathbf{hhs}=256$, $\mathbf{C.out}=32$, $\mathbf{num_fc}=3$, $\mathbf{ndr}=5$, $\mathbf{KL_weight}=0.99$, $\mathbf{NLL_}$
252 $\mathbf{weight}=1.0$, $\mathbf{MMD_weight}=10$. For S1A, the model configuration is $\mathbf{z}=4$, $\mathbf{whs}=128$, $\mathbf{hhs}=512$, $\mathbf{C.out}=256$,
253 $\mathbf{num_fc}=3$, $\mathbf{ndr}=5$, $\alpha=0.99$, $\xi=10.0$, $\lambda=1$. For lactamase, the model configuration is $\mathbf{z}=3$, $\mathbf{whs}=128$,
254 $\mathbf{hhs}=512$, $\mathbf{C.out}=512$, $\mathbf{num_fc}=1$, $\mathbf{ndr}=7$, $\mathbf{KL_weight}=0.99$, $\mathbf{NLL_weight}=100.0$, $\mathbf{MMD_weight}=10$. These
255 models were k-fold cross-validated with $k=5$ and found that the average negative log-likelihood and

256 max-mean discrepancy were consistent across different folds and close to the training values (result
257 spreadsheets found [https://github.com/PraljakReps/ProtWaveVAE/tree/01cecd054be7ca77540b](https://github.com/PraljakReps/ProtWaveVAE/tree/01cecd054be7ca77540b23d19aeaa5a3b9012d79/Pfam_analysis/outputs/train_sess/pfam)
258 [23d19aeaa5a3b9012d79/Pfam_analysis/outputs/train_sess/pfam](https://github.com/PraljakReps/ProtWaveVAE/tree/01cecd054be7ca77540b23d19aeaa5a3b9012d79/Pfam_analysis/outputs/train_sess/pfam)), indicating that the models are
259 not overfitting on the training set. Next, after verifying that the model is not overfitting on the given
260 random training set split, the final model was trained on the entire dataset using the same number of
261 epochs (500). The work presented in protein inference and generative design here utilized these models
262 trained on the whole protein family dataset.

263 2.5 Model architecture and hyperparameter optimization for benchmark 264 fitness task

265 The fitness and function benchmarking tasks aim to assess the ability of ProtWave-VAE to predict
266 regression values for various extrapolative problems. By leveraging semi-supervision, we introduce a
267 second decoder head that incorporates a regression model, allowing us to sample latent variables z and
268 predict functional/fitness continuous values y . This task tests the model’s capacity to use the inferred
269 latent variables for more than simply generative design, but for extrapolative regression prediction,
270 and more importantly, compare it against state-of-the-art models, including large language models
271 (e.g., ESM-1 [12]). The architecture of ProtWave-VAE for the fitness benchmarking tasks consists of
272 a gated dilated encoder $q_\phi(z|x)$, a latent-conditioned WaveNet decoder $p_\theta(x|z)$, and a discriminative
273 multi-layer perceptron model $p_\omega(y|z)$. The decoder component that is implemented for regression
274 follows the architecture [Decoder_re](#).

275 The following hyperparameters were optimized for the ProtWave-VAE model: (1) latent space
276 embedding size `z_dim`, (2) number of kernel channels for the dilate convolutions along the encoder
277 path `C_out`, (3) number of fully connected layers at the end of the encoder path `num_fc`, (4) num-
278 ber of channels implemented in the dilated causal convolutions along the WaveNet decoder path `whs`,
279 (5) number of channels implemented in the 1x1 convolutions along the WaveNet TopHead decoder
280 path `hhs`, (6) number of dilated causal convolutions used `ndr`, (7) number of discriminative hidden
281 module layers `disc_num_layers`, (8) width of the hidden linear layers along the discriminative de-
282 coder `hidden_width`, (9) dropout probability for the discriminative decoder path `p`, (10) negative log-
283 likelihood prefactor weight `NLL_weight`, (11) Kullback-Leibler divergence prefactor weight `KL_weight`,
284 (12) max-mean discrepancy prefactor weight `MMD_weight`, and (13) discriminative likelihood prefactor
285 weight `gamma_weight`. These hyperparameters were optimized for each of the following four protein
286 families: AAV-capsid task (FLIP), GB1 task (FLIP), GFP task (TAPE), and stability task (TAPE).
287 The hyperparameters were optimized sequentially, with the optimal values chosen based on the optimal
288 spearman ρ correlation and mean square error loss while maintaining good negative log-likelihood loss
289 and excellent max-mean discrepancy regularization loss. For FLIP benchmarks, we hyperparameter
290 optimized over the random split train/valid split and used those hyperparameters for the remaining
291 benchmark splits. While for for the TAPE benchmarks, we simply hyperparameter optimized based
292 on the chosen validation set by the TAPE authors.

293 The optimization process covered the following domains: (1) `z_dim` $\in [1, 20]$, (2) `C_out` $\in \{32, 64, 128,$
294 $256, 512\}$, (3) `whs` $\in \{32, 64, 128, 256, 512\}$, (4) `hhs` $\in \{32, 64, 128, 256, 512\}$, (5) `num_fc` $\in [0, 5]$, (6) `ndr` \in
295 $\{2, 4, 6, 8, 10\}$, (7) `disc_num_layers` $\in 1, 2, 3, 4, 5$, (8) `hidden_width` $\in \{2, 5, 10, 20, 50, 100, 200, 500, 1000\}$,
296 (9) `p` $\in \{0.0, 0.1, 0.2, 0.3, 0.4, 0.5\}$, (10) `NLL_weight` $\in \{0.1, 0.5, 1.0, 5.0, 10.0, 50.0, 100.0\}$, (11) `KL_weight` \in
297 $\{0.8, 0.85, 0.9, 0.95, 0.99\}$, (12) `MMD_weight` $\in \{1.0, 2.0, 5.0, 10.0, 20.0, 50.0, 100.0, 200.0, 500.0, 1000.0\}$, and
298 (13) `gamma_weight` $\in \{0.1, 0.5, 1.0, 5.0, 10.0, 50.0, 100.0\}$. An Adam optimizer with a learning rate of
299 1×10^{-4} was used during training, and unaligned sequences were employed for all benchmarking tasks.
300 For the TAPE benchmark tasks, the training set was split into a train/validation set, and the previ-
301 ously mentioned metrics were optimized on the validation set. The AAV capsid task used an epoch of
302 500 and a batch size of 256, while the GB1 task used an epoch of 500 and a batch size of 512. For the
303 GFP task, a batch size of 256 and an epoch of 300 were used, while for the stability task, a batch size
304 of 256 and an epoch of 300 were used.

305 The final optimized hyperparameter configurations for the AAV capsid, GB1, GFP, and stability

306 tasks are presented below. For AAV capsid, the model configuration is `z_dim = 6, C_out = 128,`
307 `whs = 128, hhs = 256, num_fc = 3, ndr = 8, disc_num_layers = 1, hidden_layers = 50, p = 0.3,`
308 `NLL_weight = 100.0, KL_weight = 0.99, MMD_weight = 10.0,` and `gamma_weight = 10`. For GB1,
309 the model configuration is `z_dim = 3, C_out = 128, whs = 256, hhs = 512, num_fc = 3, ndr =`
310 `6, disc_num_layers = 5, hidden_layers = 20, p = 0.0, NLL_weight = 50.0, KL_weight = 0.99,`
311 `MMD_weight = 10.0,` and `gamma_weight = 10`. For GFP, the model configuration is `z_dim = 4, C_out =`
312 `512, whs = 256, hhs = 256, num_fc = 2, ndr = 8, disc_num_layers = 2, hidden_layers = 500,`
313 `p = 0.3, NLL_weight = 10.0, KL_weight = 0.99, MMD_weight = 20.0,` and `gamma_weight = 50`. For
314 stability, the model configuration is `z_dim = 11, C_out = 512, whs = 64, hhs = 128, num_fc = 0,`
315 `ndr = 10, disc_num_layers = 4, hidden_layers = 50, p = 0.3, NLL_weight = 0.50, KL_weight = 0.99,`
316 `MMD_weight = 100.0,` and `gamma_weight = 5.0`, while training configuration was epochs equal to 2000
317 and batch size equal to 512. After finding the optimal hyperparameter configuration for which task,
318 we used these hyperparameters on the train/test splits defined for the benchmarking tasks.

319 2.6 Model architecture and hyperparameter optimization for Chorismate 320 Mutase task

321 To assess if ProtWave-VAE can perform semi-supervised learning for reshaping the latent space con-
322 cerning fitness or function while retaining generative capacity, we compare an unsupervised ProtWave-
323 VAE to a similar architecture with an additional regression decoder that samples from z and predicts
324 fitness y . The Chorismate mutase dataset is an ideal protein dataset as it includes enzyme sequences
325 x from a specific protein family and fitness measurements y obtained by Russ et al. [15].

326 For the unsupervised learning architecture, we employed a model architecture similar to the one
327 described above and optimized over comparable hyperparameters based on the negative log-likelihood
328 on the hold-out set. In this instance, the training set consists of natural homologs, while the hold-out set
329 contains synthetic designs. The final hyperparameter optimization values are: `z_dim=4, C_out=256,`
330 `num_fc=0, wave_hidden_state=64, head_hidden_state=512, num_dil_rates=8, NLL_weight=1, KL_`
331 `weight= 0.95,` and `MMD_weight=10`. The hyperparameter optimization search results can be found
332 https://github.com/PraljakReps/ProtWaveVAE/tree/main/Pfam_analysis/outputs/hp_optimization/
333 CM. The optimizer is Adam with a learning rate of 1×10^{-4} , 300 epochs, and a batch size of 512. The
334 model was trained on unaligned sequences.

335 For semi-supervised learning, the gated dilated encoder and latent-conditioned WaveNet decoder
336 architecture utilize the previously described hyperparameter optimized unsupervised architecture. We
337 then introduce a discriminative decoder $p_\omega(y|z)$ with a depth of `disc_num_layers=2`, a linear layer
338 width of `hidden_width=10`, and `p=0.3`. Here, we only implement a regression decoder and omit a clas-
339 sification decoder from `Decoder_re`. Next, we add a mean-squared error loss between the relative en-
340 richment (r.e.) fitness prediction and ground truth to the unsupervised loss objective, applying a prefac-
341 tor weight of 1. The results presented in Chorismate mutase (CM) inference and protein sequence gen-
342 eration for unsupervised and semi-supervised architectures use the above hyperparameters. Pretrained
343 weights for these two models can be found [https://github.com/PraljakReps/ProtWaveVAE/tree/](https://github.com/PraljakReps/ProtWaveVAE/tree/main/Pfam_analysis/outputs/train_sess/pfam/CM)
344 [main/Pfam_analysis/outputs/train_sess/pfam/CM](https://github.com/PraljakReps/ProtWaveVAE/tree/main/Pfam_analysis/outputs/train_sess/pfam/CM).

345 2.7 Model architecture and hyperparameter optimization for SH3 design 346 task

347 To evaluate the ability of ProtWave-VAE to design functional sequences using semi-supervised learn-
348 ing and alignment-free sequence inference, we tasked the model with designing functional *in vivo* SH3
349 domains in *Saccharomyces cerevisiae*. The model consists of a gated dilated encoder $q_\phi(z|x)$, a latent-
350 conditioned WaveNet decoder $p_\theta(x|z)$, and a discriminative decoder that includes a regression and clas-
351 sification path $p_\omega(y|z)$. The regression path samples z and predicts continuous values y_{re} , which repre-
352 sent the relative enrichment (r.e.) scores that measure the functionality within the Sho1 osmosensing
353 assay using next-generation sequencing. The classification path is independent of the regression path

354 and samples z to predict class probabilities y_c using a final sigmoidal function layer. The two classes are
355 functional (sequences with r.e. ≥ 0.5) versus non-functional (sequences with r.e. < 0.5) sequences. The
356 loss objective for classification is binary cross-entropy, while the loss objective for regression is mean-
357 squared error. We performed hyperparameter optimization based on good negative log-likelihood,
358 mean-squared error, and binary cross-entropy values, while maintaining excellent max-mean discrep-
359 ancy regularization values. Based on the metrics and the procedure described in the previous sections
360 for hyperparameter optimization, the final hyperparameter optimization configuration is as follows: (1)
361 `z_dim=6`, (2) `C_out=128`, (3) `num_fc=2`, (4) `disc_num_layers=2`, (5) `hidden_width=10`, (6) `p=0.4`,
362 (7) `whs=256`, (8) `hhs=512`, (9) `num_dil_rates=12`, (10) `NLL_weight=1.0`, (11) `KL_weight=0.99`, (12)
363 `MMD_weight=10.0`, and (13) `gamma_weight=1.0`. The sequence data is alignment-free, and the optimizer
364 implemented is Adam with a learning rate of 1×10^{-4} . We used a batch size of 1024 and trained for 200
365 epochs. The results of the hyperparameter optimization can be found at [https://github.com/PraljakReps](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/hp_optim)
366 [/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/hp_optim](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/hp_optim).

367 We conducted stratified cross-validation instead of vanilla k-fold cross validation because the dataset
368 contained a larger number of nonfunctional sequences than functional sequences. For stratified cross-
369 validation, we trained five different train/validation configurations while monitoring the negative log-
370 likelihood, mean-squared error loss, binary cross entropy loss, and max-mean discrepancy regularization
371 loss, as well as classification performance metrics such as precision, recall, and F1 score on the validation
372 set. The results, which can be found [https://github.com/PraljakReps/ProtWaveVAE/tree/main](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/CV)
373 [/SH3_design_project/outputs/SH3_task/CV](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/CV), show that the values are consistent and good over the
374 five stratified train/validation splits. With the hyperparameter configuration and training configura-
375 tion described above, ProtWave-VAE was trained on the entire dataset before generating and designing
376 novel sequences for experimental testing. The pretrained weights can be found [https://github.com/Pr](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/final_model)
377 [aljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/final_model](https://github.com/PraljakReps/ProtWaveVAE/tree/main/SH3_design_project/outputs/SH3_task/final_model)

378 **3 Supplementary Methods**

379 **3.1 Random mutagenesis of SH3 sequences**

380 The goal of this study was to enhance the functionality of a weak binding natural ortholog (sub-
381 group III) and a weak binding natural hof1 paralog (subgroup IV) by using N-terminus plus latent
382 conditioning and inpainting the C-terminus missing region. To avoid elevating functionality based on
383 random chance, the generative model's performance was compared with that of random mutagenese-
384 sis. To achieve this, we mutated the same C-terminus region of the weak ortholog and weak paralog
385 randomly, while maintaining the same novelty as the designed inpainted sequences. The random mu-
386 tagenesis sequences were substituted with amino acids belonging to the reference sequence at random
387 sites along the C-terminus design region until we matched the minimum Levenshtein distance that cor-
388 responds to the generative designs. By doing so, we ensured that the random mutagenized sequences
389 were mutated along the same inpainted region as the designs, while retaining the same diversity and
390 novelty as the synthetic generative designs. The results showed that for weak binding ortholog sub-
391 group IV and weak binding paralog group V, the random mutagenized sequences had matching edit
392 distances to the design's minimum Levenshtein distance to the original weak binding natural paralog
393 or ortholog, respectively. The matching edit distances were obtained when the N-terminus condition-
394 ing was 25%, 50%, or 75%. Figure S1A and B show the corresponding results for the weak binding
395 ortholog and paralog, respectively.

396 4 Supplementary Results

397 4.1 Pfam task: ColabFold structure prediction analysis of design sequences

398 A crucial evaluation of the ProtWave-VAE involves determining the extent to which the model can
399 discern biologically relevant representations without receiving any annotations. We observed that
400 the latent space sorts training sequences into phylogenetic groups and functional subclasses for S1A
401 serine protease and beta-lactamase protein families (Figure S2A). The following assessment involves
402 confirming if ProtWave-VAE can sample from this specific biologically significant latent space and
403 produce new artificial sequences indistinguishable from the training set.

404 To evaluate the model’s generative capabilities, we created alignment-free sequences, predicted
405 associated tertiary structures utilizing ColabFold [8], and assessed if predictions mirrored the ter-
406 tiary structures of natural homologs. For DHFR, S1A proteases, and lactamase protein families, we
407 sampled 100 latent vectors z from an isotropic Gaussian distribution and employed the ProtWave-
408 VAE autoregressive decoder $p_{\theta}(x|z)$ to generate artificial sequences. Subsequently, we predicted the
409 tertiary structure of each artificial sequence and calculated the TMscores and heavy-atom root mean
410 squared distances (RMSDs). The anticipated tertiary structure with maximum, median, and minimum
411 TMscores for the DHFR, S1A protease, and lactamase protein families are displayed in Figure S2B.
412 We observed that the structure corresponding to the median TMscore accurately reflects the natural
413 homolog’s tertiary structure.

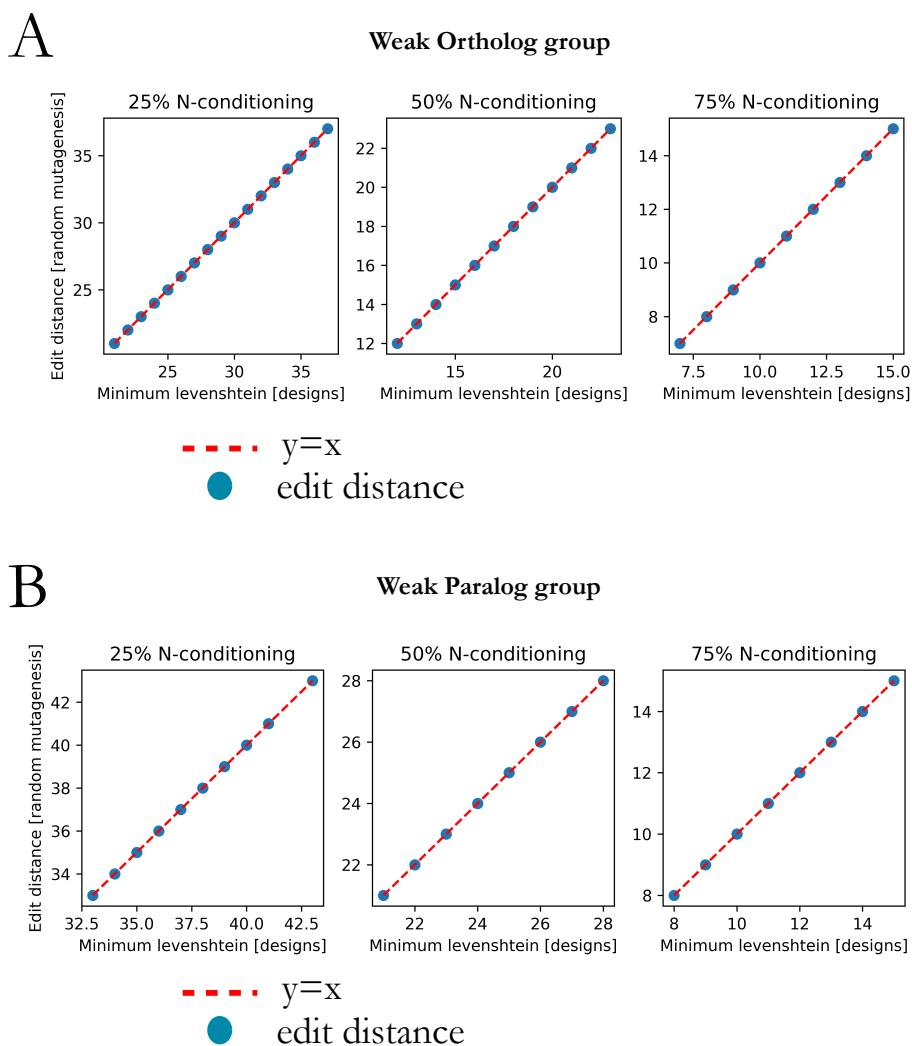


Figure S1: Random mutagenesis control for C-terminus diversification of protein designs with N-terminus plus latent conditioning. (A) The edit distance – number of random substitutions along the remaining C-terminus region of the reference weak binding ortholog – is shown on the y-axis for 25%, 50%, and 75% N-terminus conditioned subgroups. The x-axis shows the minimum Levenshtein distance between the designed sequences using N-terminus plus latent conditioning generative approach and weak binding natural ortholog sequence. (B) Similarly, the edit distance – number of random substitutions along the remaining C-terminus region of the reference weak binding paralog – is shown on the y-axis for 25%, 50%, and 75% N-terminus conditioned subgroups. The x-axis shows the minimum Levenshtein distance between the designed sequences using N-terminus plus latent conditioning generative approach and weak binding natural paralog sequence.

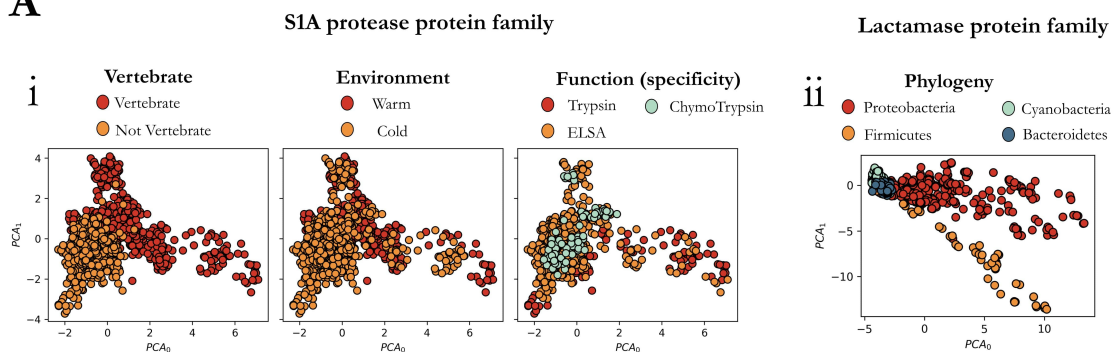
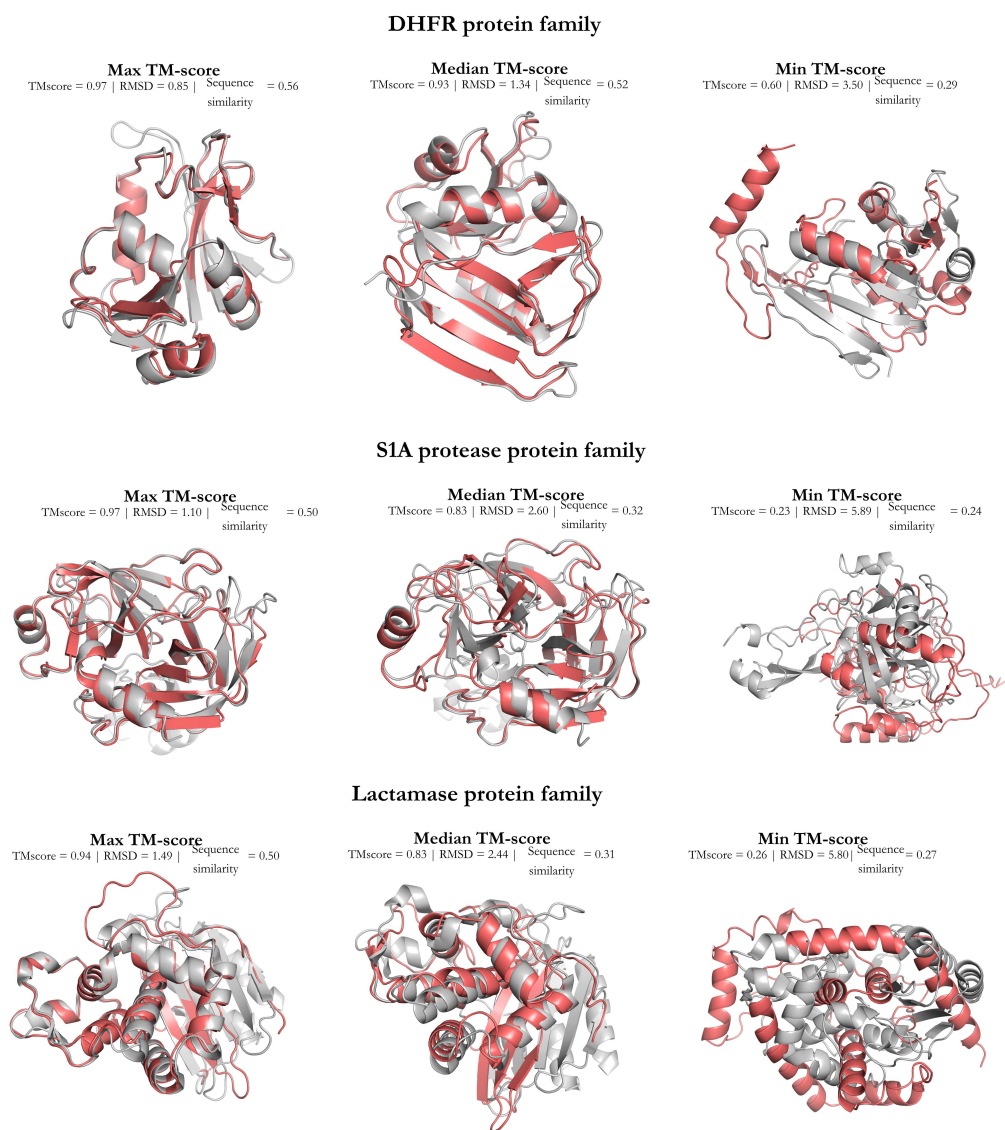
A**B**

Figure S2: ProtWave-VAE infers meaningful biological representations on alignment-free protein families. (A) Principal component analysis (PCA) projections of the inferred latent spaces for the S1A proteases and beta-lactamase families are presented. The unsupervised model of the S1A protease family disentangles homologs within the inferred latent space based on vertebrate, environmental conditions, and functional specificity. For the beta-lactamase family, the model disentangles homologs in the inferred space in terms of the phylogeny. (B) To test the generative capacity of ProtWave-VAE, we randomly sampled 100 latent vectors z for each protein family from a normal distribution $\mathcal{N}(0, I)$, corresponding to the latent prior. Using a computational structure prediction workflow (ColabFold + TMalign), we predicted each structure of the sample sequences and compared the predicted structure against a natural homolog that defines the corresponding protein family. We retrieved TMscores and root-mean-square distance (RMSD) scores. The structure predictions of ProtWave-VAE novel design sequences (red) for DHFR, S1A protease, and lactamase are visualized with the alignment of maximum, median, and minimum TM-score synthetic sequences against the natural reference homolog structure (grey).

414 4.2 SH3 design task: experimental results

415 The relative enrichment (r.e.) score provides a quantitative measurement of the degree to which our
 416 designed SH3 domains are functional *in vivo* and capable of activating a homeostatic osmoprotective
 417 response. The assay shows good reproducibility in independent trials ($R^2 = 0.94$, $\rho_{pearson}=0.97$,
 418 $n=1002$, $p < 10^{-307}$, Figure S3).

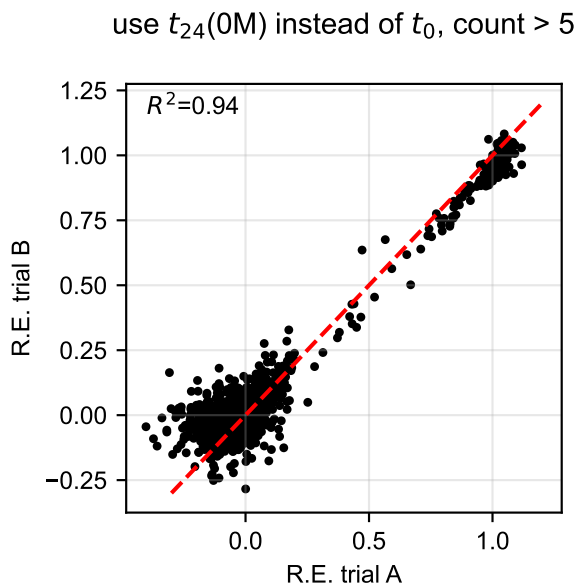


Figure S3: Validation of the high-throughput select-seq assay. A scatter plot of the enrichment score relative to wild-type *r.e.* for two independent ($n = 1002$ including standard curve mutants) trials of the select-seq assay under the same experimental conditions. The position of the wild-type Sho1 sequence is normalized to be at (1, 1) and of the null allele is at (0, 0). The red dashed line indicates the identity trace. The data shows that the select-seq assay shows good reproducibility between independent runs ($R^2 = 0.94$; $\rho_{Pearson} = 0.97$, $n = 1002$, $p < 10^{-307}$)

419 **4.3 Comparison on generative performance for protein families**

Table S1: This table provides a comprehensive overview of the scores associated with the four protein families as predicted by the three generative models: WaveNet (Shin et al. [17]), ProteinGAN (Repecka et al. [11]), and ProtWave-VAE (ours). It offers values for RMSD, TMscore, and sequence similarity (seq. sim.) in terms of their maximum, median, and minimum values. In the table, bold text is used to highlight the best scores for each metric and protein family combination.

Metric	WaveNet				ProteinGAN				ProtWave-VAE			
	G-protein	DHFR	Lactamase	SIA proteases	G-protein	DHFR	Lactamase	SIA proteases	G-protein	DHFR	Lactamase	SIA proteases
RMSD (max) ↓	1.902	1.466	1.9	1.662	5.824	5.754	6.028	6.328	3.386	3.49	6	6.17
RMSD (median) ↓	1.412	1.033	1.256	1.058	1.601	1.874	1.402	3.075	1.633	1.31	2.644	2.322
RMSD (min) ↓	1.093	0.806	1.038	0.796	1.245	1.102	1.152	1.432	0.926	0.277	1.494	1.1
TMscore (max) ↑	0.955	0.967	0.968	0.98	0.945	0.946	0.962	0.943	0.967	0.967	0.941	0.967
TMscore (median) ↑	0.933	0.954	0.96	0.969	0.917	0.877	0.951	0.725	0.911	0.933	0.836	0.861
TMscore (min) ↑	0.903	0.923	0.905	0.931	0.198	0.223	0.229	0.245	0.549	0.6	0.256	0.227
Seq. sim. (max) ↓	0.627	0.687	0.714	0.659	0.557	0.587	0.668	0.473	0.658	0.564	0.503	0.537
Seq. sim. (median) ↓	0.563	0.6	0.678	0.595	0.5	0.457	0.623	0.37	0.473	0.48	0.333	0.335
Seq. sim. (min) ↓	0.474	0.507	0.512	0.5	0.234	0.245	0.213	0.243	0.266	0.278	0.2	0.226

Table S2: This table provides a comprehensive overview of the scores associated with the N-terminus prompting for the Chorismate mutase family as predicted by the ProtWave-VAE (ours) and WaveNet decoder (Shin et al. [17]). Bold values indicates the **best values** between WaveNet decoder versus ProtWave-VAE.

Metric	No N-terminus Prompt		N-terminus Prompt	
	WaveNet	ProtWave-VAE	WaveNet	ProtWave-VAE
RMSD (max) ↓	3.660	4.192	0.772	2.454
RMSD (median) ↓	1.407	2.063	0.561	1.142
RMSD (min) ↓	0.626	0.876	0.468	0.572
TMscore (max) ↑	0.966	0.920	0.952	0.943
TMscore (median) ↑	0.871	0.746	0.945	0.888
TMscore (min) ↑	0.444	0.318	0.929	0.748
Seq. sim (max) ↓	0.882	0.463	0.936	0.734
Seq. sim (median) ↓	0.532	0.340	0.883	0.617
Seq. sim (min) ↓	0.404	0.229	0.819	0.489

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] Drew H Bryant, Ali Bashir, Sam Sinai, Nina K Jain, Pierce J Ogden, Patrick F Riley, George M Church, Lucy J Colwell, and Eric D Kelsic. Deep diversification of an aav capsid protein by machine learning. *Nature Biotechnology*, 39(6):691–696, 2021.
- [3] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [4] Christian Dallago, Jody Mou, Kadina E Johnston, Bruce J Wittmann, Nicholas Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K Yang. Flip: Benchmark tasks in fitness landscape inference for proteins. *bioRxiv*, pages 2021–11, 2021.
- [5] Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. In *International Conference on Machine Learning*, pages 933–941. PMLR, 2017.
- [6] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [7] Xinran Lian, Nikša Praljak, Subu Subramanian, Sarah Wasinger, Rama Ranganathan, and Andrew L Ferguson. Deep learning-enabled design of synthetic orthologs of a signaling protein. *bioRxiv*, pages 2022–12, 2022.
- [8] Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin Steinegger. Colabfold: making protein folding accessible to all. *Nature Methods*, 19(6):679–682, 2022.
- [9] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [10] Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in Neural Information Processing Systems*, 32, 2019.
- [11] Donatas Repecka, Vykintas Jauniskis, Laurynas Karpus, Elzbieta Rembeza, Irmantas Rokaitis, Jan Zrimec, Simona Poviloniene, Audrius Laurynenas, Sandra Viknander, Wissam Abujawa, et al. Expanding functional protein sequence spaces using generative adversarial networks. *Nature Machine Intelligence*, 3(4):324–333, 2021.
- [12] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- [13] Olivier Rivoire, Kimberly A Reynolds, and Rama Ranganathan. Evolution-based functional decomposition of proteins. *PLoS Computational Biology*, 12(6):e1004817, 2016.
- [14] Gabriel J Rocklin, Tamuka M Chidyausiku, Inna Goreshnik, Alex Ford, Scott Houliston, Alexander Lemak, Lauren Carter, Rashmi Ravichandran, Vikram K Mulligan, Aaron Chevalier, Cheryl H Arrowsmith, and David Baker. Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science*, 357(6347):168–175, 2017.

- 462 [15] William P Russ, Matteo Figliuzzi, Christian Stocker, Pierre Barrat-Charlaix, Michael Socolich, Pe-
463 ter Kast, Donald Hilvert, Remi Monasson, Simona Cocco, Martin Weigt, and Rama Ranganathan.
464 An evolution-based model for designing chorisimate mutase enzymes. *Science*, 369(6502):440–445,
465 2020.
- 466 [16] Karen S Sarkisyan, Dmitry A Bolotin, Margarita V Meer, Dinara R Usmanova, Alexander S
467 Mishin, George V Sharonov, Dmitry N Ivankov, Nina G Bozhanova, Mikhail S Baranov, Onu-
468 ralp Soylemez, Natalya S. Bogatyreva, Peter K. Vlasov, Evgeny S. Egorov, Maria D. Logacheva,
469 Alexey S. Kondrashov, Dmitry M. Chudakov, Ekaterina V. Putintseva, Ilgar Z. Mamedov, Dan S.
470 Tawfik, Konstantin A. Lukyanov, and Fyodor A. Kondrashov. Local fitness landscape of the green
471 fluorescent protein. *Nature*, 533(7603):397–401, 2016.
- 472 [17] Jung-Eun Shin, Adam J Riesselman, Aaron W Kollasch, Conor McMahon, Elana Simon, Chris
473 Sander, Aashish Manglik, Andrew C Kruse, and Debora S Marks. Protein design and variant
474 prediction using autoregressive generative models. *Nature Communications*, 12(1):2403, 2021.
- 475 [18] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
476 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*
477 *learning research*, 15(1):1929–1958, 2014.
- 478 [19] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, and Ko-
479 ray Kavukcuoglu. Conditional image generation with pixelcnn decoders. *Advances in Neural*
480 *Information Processing Systems*, 29, 2016.
- 481 [20] Nicholas C Wu, Lei Dai, C Anders Olson, James O Lloyd-Smith, and Ren Sun. Adaptation in
482 protein fitness landscapes is facilitated by indirect paths. *Elife*, 5:e16965, 2016.
- 483 [21] Ran Zhang, Lin Cao, Mengtian Cui, Zixian Sun, Mingxu Hu, Rouxuan Zhang, William Stuart,
484 Xiaochu Zhao, Zirui Yang, Xueming Li, Yuna Sun, Shentao Li, Wei Ding, Zhiyong Lou, and
485 Zihe Rao. Adeno-associated virus 2 bound to its cellular receptor aavr. *Nature Microbiology*,
486 4(4):675–682, 2019.