

# ESMS VAE: A Structure-Informed Variational Autoencoder for Protein Engineering

Danny Ahn<sup>1\*</sup>, Shihyun Moon<sup>2</sup>, Jooyoung Jung<sup>3</sup>, Minjae Lee<sup>4</sup>,  
Jeongsu Park<sup>5</sup>

<sup>1</sup>Daegu Science High School, 154 Dongdaegu-ro, Suseong-gu, Daegu,  
Republic of Korea.

\*Corresponding author(s). E-mail(s): [ahnd6474@gmail.com](mailto:ahnd6474@gmail.com);

## Abstract

Sequence-based protein Variational Autoencoders (VAEs) have definitive limits on generalization, reconstruction ability, and capturing meaningful information within the latent space. To overcome these limits, we introduce ESMS VAE, a novel VAE trained with a custom loss function designed to capture structural similarity. ESMS VAE achieved a 97.17% reconstruction rate on a test set randomly sampled from Uni Ref 50. This high reconstruction rate was maintained under small noise conditions and increased by approximately 1% with heavy noise. We randomly sampled and decoded latent vectors to test their ability to generate novel sequences. The decoded sequences showed a maximum identity of around 10% against sequences used in training. The latent space’s capability in downstream tasks was evaluated on the DMS dataset from protein gym and fluorescent proteins (FPs). In the case of DMS, this indicates  $\rho(\text{spearman})$  of 0.7779, effectively capturing mutational changes. In FP, it demonstrated excellent performance as an embedding for classifying FPs and non-FPs, with a 0.987 5-fold cross-validation accuracy, and for regressing excitation and emission wavelengths, achieving RMSE values of 2.7 nm and 3.8 nm, respectively, using Gaussian Process (GP) models. FPs were projected onto the latent space and were clustered using K-means. Consensus vectors were calculated and decoded, generating sequences that largely preserved their function. All trained models and code are available on GitHub ESMS VAE <https://github.com/Ahnd6474/ESMS-VAE>

**Keywords:** Variational Autoencoder, Protein Engineering, Deep Learning, Generative Models, Structural Biology, ESMS

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Protein VAE	2
1.2	ESM2 and ESMS	3
1.3	Our Contribution	3
<b>2</b>	<b>Methods</b>	<b>3</b>
2.1	Model Architecture and Loss Function	3
2.2	Training and Model Selection	4
<b>3</b>	<b>Results</b>	<b>6</b>
3.1	Model Reconstruction and Robustness	6
3.2	Mutation Effect Prediction on DMS dataset	7
3.3	Case Study: Fluorescent Protein (FP) Analysis and Generation	8
3.4	Ablation Study	11
3.5	Limitation: Thermo stability (Tm) Prediction	12
<b>4</b>	<b>Discussion</b>	<b>13</b>
<b>5</b>	<b>Conclusion</b>	<b>13</b>
	<b>References</b>	<b>13</b>

## 1 Introduction

A Variational Autoencoder (VAE) [1] is a generative model capable of generating similar but different data from the input. It is usually applied to images and is used to reconstruct and generate images [2]. In protein science, VAEs have been suggested to guide the exploration of the vast and empty protein space with their latent space [2]. A VAE is made up of an encoder and a decoder. The encoder compresses data and embeds it in a latent space, while the decoder reconstructs the original data from the latent space. The model is trained in a manner that minimizes information loss while compressing data and fits latent space variables to a normal distribution. However, VAEs are susceptible to issues such as posterior collapse, where the decoder ignores the latent space and generates the most common token, and KL vanishing, where an overly perfect fit of the latent space variables to a Gaussian distribution causes a loss of information in the latent space [3]. Both phenomena are associated with very low KL values, which is why a KL value of around 0.05 is often recommended.

### 1.1 Protein VAE

It has been demonstrated that the continuous latent space of a VAE can be applied to latent space interpolation and objective-driven amino acid optimization. Early models like Deep-sequence utilized a VAE that takes a Multiple Sequence Alignment (MSA) as input [4]. It was used to predict mutations, and its latent space could be used as an embedding for other machine-learning tasks. Subsequent protein VAEs continuously

evolved by taking a single sequence as input [5]; however, sequences contain smaller information compared to MSAs, and these models usually showed moderate reconstruction rates. A notable advance is ProT-VAE, which uses ProT5 embeddings and a transformer architecture to capture long-distance relations between amino acids [6]. ProT-VAE showed a nearly 100% reconstruction rate on fine-tuned protein families.

## 1.2 ESM2 and ESMS

ESM-2 is a large language model released by Meta AI, trained with masking, and is capable of capturing long and short-distance interactions within a sequence [7]. ESM-2 was used as an embedding for ESMFold, a 3D structure prediction model known to have similar capabilities to AlphaFold 2 [8, 9]. It is commonly used as an embedding for various tasks like secondary structure prediction and thermostability prediction. Although its ability to capture structural information is undisputed, it is heavy. Needing a lighter model, Knowledge Distillation [10] was applied to ESM embeddings. A new small transformer model, ESMS, was trained to mimic ESM2’s embedding direction and size. Direction and size were each calculated with cosine similarity and RMSE. Cosine similarity and RMSE on the test set were 0.9647 and 1.2998.

## 1.3 Our Contribution

With ProT-VAE, it has been shown that transformer architectures are capable of capturing hidden information in a sequence. Also, rather than training big transformer models, using a pre-trained model as an embedding was demonstrated to be effective. However, only training with a sequence has clear limits on generalization, function, and versatility. It is clear that structural information must be introduced to train a fully functional VAE. ESMS VAE, in order to overcome such limits, introduced a custom loss function that explicitly forces the VAE to learn and represent structural information in the latent vector space.

# 2 Methods

## 2.1 Model Architecture and Loss Function

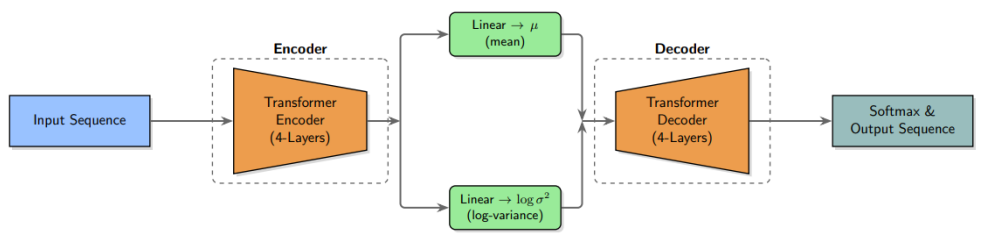
To capture structural information, the ESMS model was used to calculate the loss function. ESMS, the student model of ESM2 650M, learned embeddings that capture structural information and were used to calculate the difference in structure between the original (origin) and reconstructed (recon) sequence. This way, the latent vector space was trained to capture structural information. The full ESM2 model could not be used because of limits in computational resources, especially due to memory usage; thus, as a substitute, ESMS was chosen. ESMS is suitable for this task because it was trained using masking and shows how likely an amino acid is to take place in a certain position. If A and B were both likely to be in a certain position, A and B would have similar embeddings. This is not only because they play a similar role and can be replaced, but also because they both cannot be replaced with certain other amino acids in that position. Thus, structural loss will not penalize amino acids that are likely

to be substituted but will penalize those that are not. This also contributes to preventing posterior collapse, which usually occurs by learning the most common amino acid and relying on it. If such an event were to occur, structural loss would penalize it highly and force the VAE to rely on the latent space. This structural loss is similar to perceptual loss in image VAEs. Both concepts try to overcome the information bottleneck from pixel and sequence level matching and try to use deep features that represent perceptual and structural information [11]. Image VAE was able to achieve detail preservation: blur reduction, semantic consistency, Robustness against noise, and improved perceptual quality. In ESMS VAE, we can expect conservation of structural information and motifs, robustness against mutation and noise, and improvement in the quality of generated sequences. The composite loss function is defined as:

$$L = \lambda(L_{\text{MSE}} + L_{\text{COS}}) + \alpha \cdot L_{\text{CE}} + \beta \cdot L_{\text{KL}} \quad (1)$$

where  $L_{\text{COS}} = 1 - \text{COS}(\text{ESMS}(\text{origin}), \text{ESMS}(\text{recon}))$   
and  $L_{\text{MSE}} = \text{MSE}(\text{ESMS}(\text{origin}), \text{ESMS}(\text{recon}))$

The weights are set as  $\lambda = 5$ ,  $\alpha = 30$  for epochs  $< 100$  (then 0.1), and  $\beta = 0$  for epochs  $< 100$  (then 0.1). These hyperparameters were chosen so that VAE concentrates on reducing CE first and reduces CE even when it is small.  $\beta$  Unlike other VAEs is rather unstable, but allows structural and KL loss to have similar losses at beginning of training.



**Fig. 1:** The architecture of ESMS VAE. An input sequence is processed by a 4-layer Transformer Encoder, which outputs the parameters (mean and log-variance) of the latent distribution. A latent vector is sampled and then reconstructed into the output sequence by a 4-layer Transformer Decoder.

ESMS VAE is a comparably lightweight transformer with 5.5M parameters, composed of 4-layer transformer encoders and decoders (Figure 1). The hyperparameters are detailed in Table 1. The Adam [12] optimizer was used.

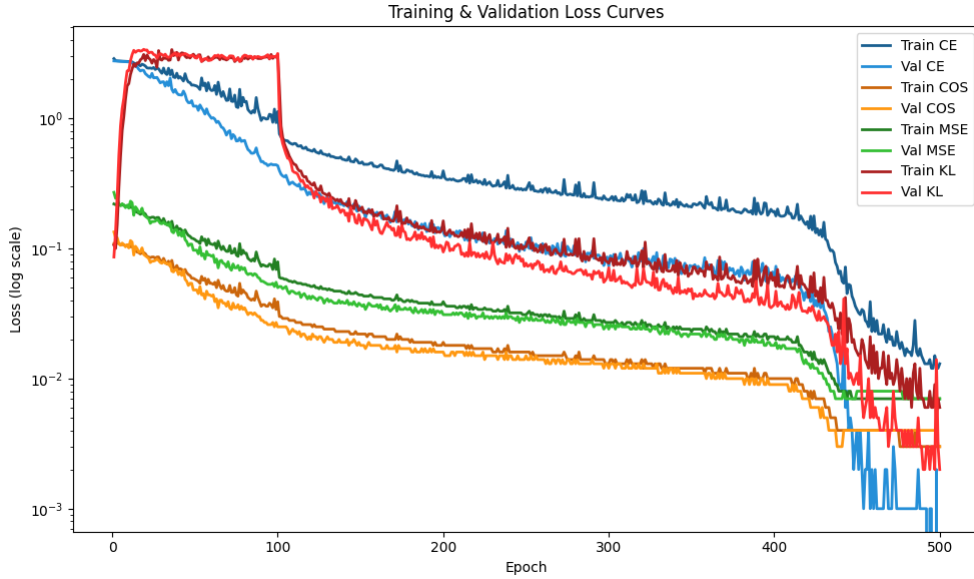
## 2.2 Training and Model Selection

The model was trained using two T4 GPU sessions provided by Kaggle, using a random subsample of the Uni-ref 50 dataset. Monitored learning curves did not show any signs of overfitting (Figure 2). The model saved at epoch 500 (VAE\_500) had

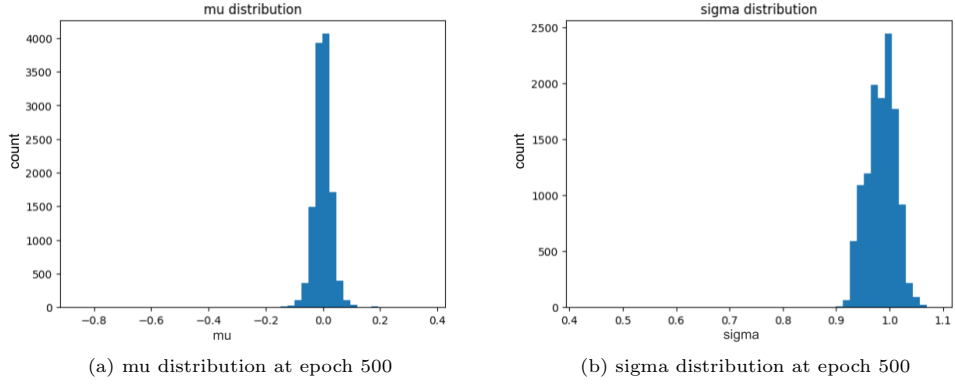
Parameter	Value
Vocab Size	33
d_model	256
Latent Dim	256
n_heads	4
Feed Forward	512
Dropout	0.3

**Table 1:** Hyperparameters for ESMS VAE.

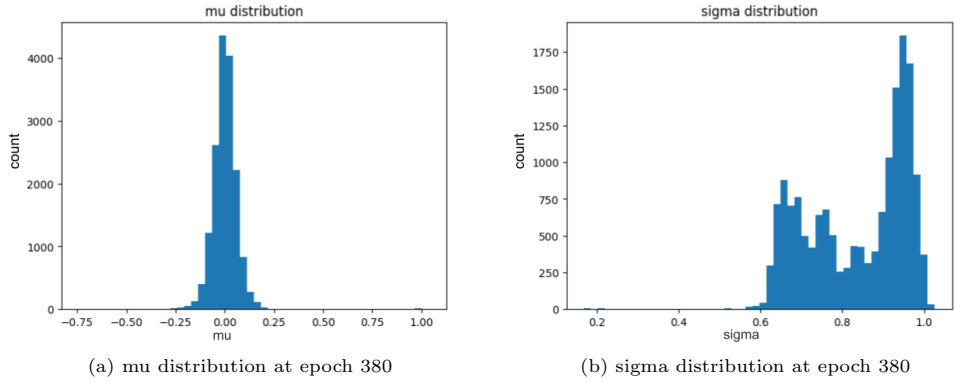
a final reconstruction rate of 99.976%. Validation loss values were: Val CE=0.000, COS=0.003, MSE=0.007, and KL=0.002. It exhibited very small latent space values, and the very low KL value suggested potential KL Vanishing (Figure 3). When noise was added to the latent space, it did not show significant changes in the reconstruction rate. Instead, the model from epoch 380 (VAE\_380) was chosen because it had a KL value closer to the recommended active value of 0.05 and had the lowest validation Cross-Entropy (CE) loss. At epoch 380, the validation losses were: Val CE=0.072, COS=0.010, MSE=0.020, and KL=0.048.



**Fig. 2:** Training and validation loss curves for Cross-Entropy (CE), Cosine Similarity (COS), Mean Squared Error (MSE), and KL Divergence (KL) over 500 epochs. The y-axis is on a log scale.



**Fig. 3:** The mean ( $\mu$ ) and standard deviation ( $\sigma$ ) distribution of the latent space at epoch 500, showing signs of KL vanishing.



**Fig. 4:** The mean ( $\mu$ ) and standard deviation ( $\sigma$ ) distribution of the latent space at epoch 380, showing healthier distributions.

### 3 Results

#### 3.1 Model Reconstruction and Robustness

VAE\_380 achieved a final reconstruction accuracy of 97.17% on a test set, demonstrating that ESMS VAE generalizes well to different kinds of protein (Figure 4). To check if the latent space had learned meaningful representations, the following tests were conducted in a teacher-forced manner on completely unseen sequences.

VAE\_380's reconstruction rate was tested after adding noise to the latent space. For a test set of about 1000 sequences, the reconstruction rate increased with noise (Table 2), suggesting the VAE is capable of distinguishing and potentially correcting

large noise. This denoising capability was further tested by substituting amino acids with 'X' in a random sample of 100 sequences. The ESMS cosine similarity between the original and reconstructed sequences was 0.9592, which is very close to 1, meaning the VAE is capable of identifying and eliminating noise while not altering structure.

Noise Level ( $\sigma$ )	Mean Reconstruction Accuracy
0.1	97.77%
1.0	98.34%
2.0	98.58%

**Table 2:** Reconstruction accuracy under noisy conditions.

The KL divergence per dimension was checked to ensure the model had not suffered from posterior collapse. The mean KL was 0.04998 with an RMSE of 0.07027. A KL value around 0.05 is considered active, indicating that the latent space is being utilized by the decoder.

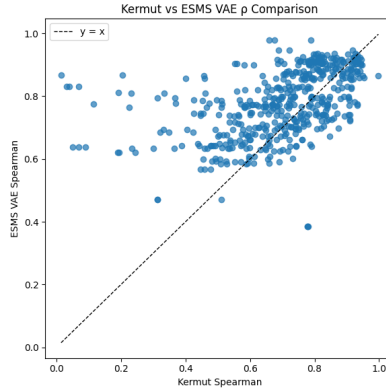
Random vectors were sampled from the latent space and decoded. These randomly generated sequences were compared with the original sequences and calculated maximum identity; identity mentioned below will mean maximum identity. As shown in Table 3, ESMS VAE is not generating copies of the trained sequence but is capable of generating completely different protein sequences.

Statistic	Value
Median identity	0.111
90th percentile identity	0.147
Max identity	0.2121

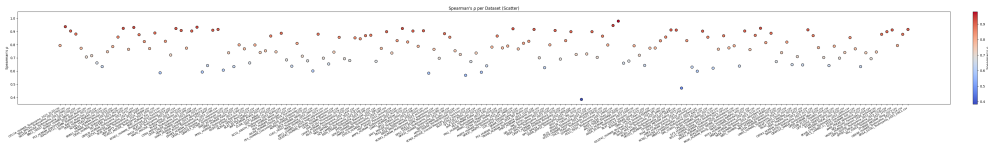
**Table 3:** Novel generation test results.

### 3.2 Mutation Effect Prediction on DMS dataset

Model as tested on DMS substitution dataset from protein gym [13]. Out of all 217 datasets, 162 that had length shorter than 512 were selected. For this, we used a shallow MLP head of seven layers, while using the latent space as an embedding and ran for 300 epochs for each dataset. The result is shown in Figure 6. Mean of 162  $\rho$ (Spearman) value is 0.7779. This result was compared with supervised DMS substitution model, Kermut. Kermut is a model that has first place in the supervised leaderboard. Kermut and ESMS VAE as compared on 162 datasets, that was compared on had a mean Spearman constant of 0.6982.



**Fig. 5:** Comparison of ESMS VAE and Kernut on Spearman correlation.



**Fig. 6:** Spearman  $\rho$  correlation for DMS datasets.

### 3.3 Case Study: Fluorescent Protein (FP) Analysis and Generation

The utility of the learned latent space was evaluated on downstream tasks involving Fluorescent Proteins (FPs). FP sequences were collected by hand from FPbase [14] and uploaded on kaggle. This was done by hand because there was no suitable API for collection purpose on FPbase.

The latent space of the VAE was applied to the classification of FPs and non-FPs, and to estimating maximum absorption and emission wavelengths. A classic Gaussian Process (GP) model was used as the model head [15]. The GP Classifier (GPC) had a 0.987 5-fold CV accuracy. The GP Regressor (GPR) had an RMSE of 2.70 nm for absorption wavelength ( $A_{abs}$ ) and 3.80 nm for emission wavelength ( $x_{em}$ ). The classification reports on the train and test sets show the model performs well and is not overfit (Tables 4 and 5).

A t-SNE map [16] shows that the latent space successfully captures structural information, separating the two classes very well. Non-fluorescent proteins are on the inner side of two curves and fluorescent proteins are on the outside. Visualizations for emission and absorption wavelengths show continuous color gradients, meaning the latent space has separated proteins with different wavelengths and gathered proteins with similar wavelengths together (Figure 7).

Collected FP sequences were projected onto the latent space and clustered using K-means. Three clusters were found, with statistics detailed in Table 6. Consensus vectors

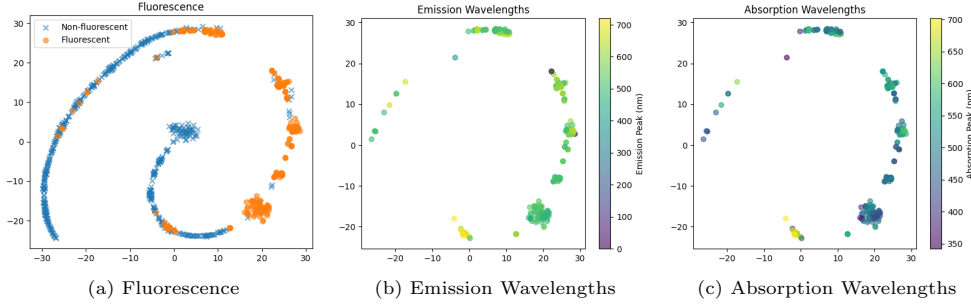


	Precision	Recall	F1-score	Support
Non-FP (0)	0.9920	0.9940	0.9930	501
FP (1)	0.9880	0.9840	0.9860	250
Accuracy			0.9907	751
Macro Avg	0.9900	0.9890	0.9895	751
Weighted Avg	0.9907	0.9907	0.9907	751

**Table 4:** Classification report on the training set.

	Precision	Recall	F1-score	Support
Non-FP (0)	0.9840	0.9840	0.9840	125
FP (1)	0.9683	0.9683	0.9683	63
Accuracy			0.9787	188
Macro Avg	0.9761	0.9761	0.9761	188
Weighted Avg	0.9787	0.9787	0.9787	188

**Table 5:** Classification report on the test set.



**Fig. 7:** t-SNE visualization of the FP latent space. (a) Clear separation of fluorescent (orange) and non-fluorescent (blue) proteins. (b, c) Continuous gradients for emission and absorption wavelengths, respectively.

were calculated based on these clusters and decoded to generate novel sequences. Generated sequences were truncated to the mean length of their cluster since VAEs usually generate meaningless sequences after the original sequence length. These sequences were classified using the trained GPC. As shown in Table 7 and Figure 8, Cluster 1 yielded a 100% success rate. This suggests that the quality and consistency of samples within a cluster are more important than the sheer number of samples for generating new functional proteins. Generated proteins had 96%–75% identity with FP sequences that was collected from FPbase. This is likely because the vectors were very alike.

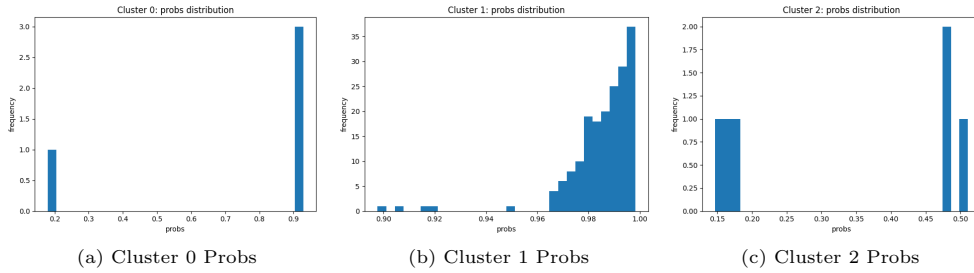
The 3D structures of generated vectors were predicted using the AlphaFold server. For Cluster 1 proteins, iconic beta barrels were found, maintaining a high pLDDT

Cluster	n (samples)	Mean Length	Std (Length)
0	22	316	5.28
1	318	234	10.00
2	14	118	29.71

**Table 6:** Statistics for FP clusters identified by K-means.

Cluster	Total Generated	Count as GFP	Ratio
0	4	3	0.7500
1	181	181	1.0000
2	6	1	0.1667

**Table 7:** Generation of GFP-like sequences from cluster consensus vectors.



**Fig. 8:** Probability distribution of generated sequences being classified as GFP for each cluster. Cluster 1 shows high-confidence predictions.

above 90. Proteins from Cluster 0 have two sets containing three beta strands on the inside covered by alpha helices. Proteins from Cluster 2, although not showing a beta-barrel structure, also show structural similarity among themselves, with three beta strands surrounded by alpha helices. This confirms that close vectors in the latent space encode for proteins with similar structures (Figure 9).

Three sequences were chosen for future in vitro validation based on: 1) GPC prediction, 2) structural likelihood to an original sequence, and 3) whether it needs other fluorescent molecules. Two sequences were chosen from Cluster 1 and one from Cluster 2. The selected sequences are:

```
>seq_cluster1_1
MASTPFKFQLKGTINGKSFTVEGEGEGNSHEGSHKGKYVCTSGKLPMSWAALGTS
FGYGMKYYTKYPSGLKNWFHEVMPEGFT
>seq_cluster1_2
MVSTGEELFTGVVPKFQLKGTINGKSFTVEGEGEGNSHEGSHKGKYVCTSGKLP
MSWAALGTSFGYGMKYYTKYPSGLKNWFF
>seq_cluster2_1
```

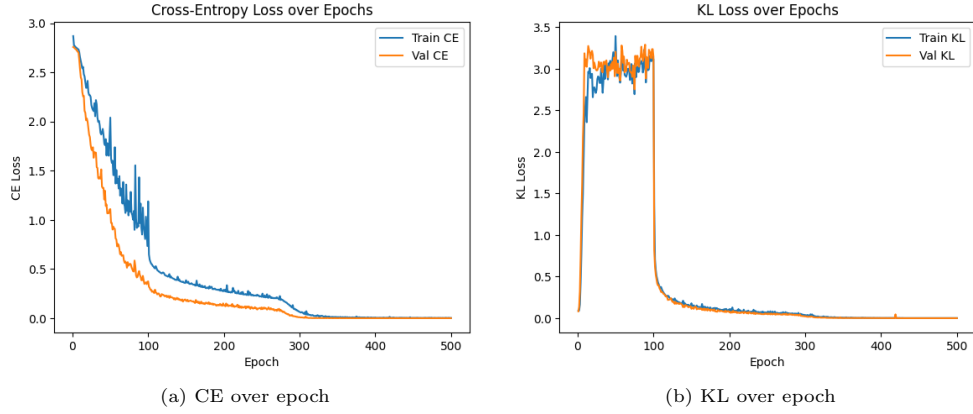


**Fig. 9:** Predicted 3D structures of generated proteins from different clusters show high intra-cluster structural similarity. Cluster 1 notably forms the GFP-like beta barrel.

MPRISDKLMKTRWRGFHSIPSIPDLGGIYGIGEKTSRRKTTEHLYTGRAKDIKS  
RLMKHKYGHQAIDRKIRSNIKQKKLSDLRFKFVE

### 3.4 Ablation Study

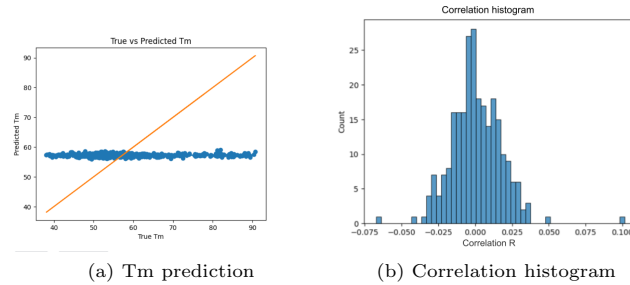
An ablation study as done in this case loss function is  $L = \alpha CE + \beta KL$ . The weights are set as  $\alpha = 30, \beta = 0$  for epochs  $\leq 30$  else  $\alpha = 0.1, \beta = 0.1$ . In the ablation study KL vanishing occurred right after epoch 100 proving role of structural loss in training.



**Fig. 10:** CE and KL value over epochs for the ablation study model without structural loss.

### 3.5 Limitation: Thermo stability (Tm) Prediction

ESMS VAE was used as an embedding, and an MLP head was used to predict Tm value of proteins from multiple classes. As shown in Figure 11a, the MLP did not perform better than just returning the mean value of Tm. This is because none of the 256 latent vectors had a strong correlation with Tm, as seen in Figure 11b. This is likely because ESMS VAE was trained to learn structural likeness rather than the thermo stability of proteins, which explains why the model excelled on the DMS dataset.



**Fig. 11:** (a) Tm prediction results using an MLP head, showing performance similar to a mean-prediction baseline. (b) A histogram of the correlation coefficients between latent space dimensions and Tm values, indicating no strong correlators.

## 4 Discussion

ESMS embedding seems to concentrate on structural integrity and likelihood. Thus, the latent space contained structural information and had high Spearman scores in the DMS dataset; however, it seems to have missed Tm-related information.

It can be thought of using the AlphaFold structure and contact map as a structural loss; however, AlphaFold commonly takes about 10 to 20 minutes per sequence, which is too long to be applied to mass data. On the other hand contact map is too costly, using sequences that have a known structure limits the data pool, and the contact map itself requires large memory. It may be suggested that using other PLMs, like ProTrans 5 and ProtBert. Like ESMS and ESM2, it is highly recommended to use smaller versions of PLM.

The latent space contains structural information; thus, like protein space itself, most of the latent space would hold functionless protein. The latent space of ESMS VAE is still too vast, and functional proteins are too scarce. It is highly recommended to use a regressor and a classifier together. The regressor is used to predict your objective, and the classifier will identify a functional protein. So, the regressor will guide your search while the classifier limits your search. This combination will be critical for efficiently navigating the latent space to engineer novel proteins.

## 5 Conclusion

ESMS VAE is the first protein VAE to be trained with a loss function that considers structural information. This concept of structural loss was valid and was capable of creating a latent space that contained structural information.

## Acknowledgements

This research was supported by Daegu Science High School.

This research was supported by Anseong Topbob

## References

- [1] Kingma, D. P. & Welling, M. Auto-encoding variational bayes (2022). URL <https://arxiv.org/abs/1312.6114>. arXiv:1312.6114.
- [2] Greener, J. G., Sternberg, M. J. E. & Jones, D. T. Protein sequence-based inference of structure, function and stability using a deep vae. *Nature Communications* **9**, 3095 (2018).
- [3] Lucas, J., Tucker, G., Grosse, R. & Norouzi, M. Understanding posterior collapse in generative latent variable models. *Journal of Machine Learning Research* **21**, 1–39 (2019).

- [4] Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deepsequence: Deep generative models of genetic variation capture the effects of mutations. *Nature Methods* **15**, 816–822 (2018).
- [5] Sinai, S., Hochreiter, S. & Obermayer, K. Lc-vae: Latent-consensus variational autoencoder for protein sequences. *arXiv preprint arXiv:1712.03346* (2017).
- [6] Elnaggar, A., Heinzinger, M., Dallago, C. *et al.* Prottrans: Towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [7] Rives, A., Meier, J., Sercu, T. *et al.* Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **118**, e2016239118 (2021).
- [8] Lin, Z., Akin, H., Rao, R. *et al.* Esmfold: End-to-end single-sequence protein structure prediction. *bioRxiv* (2022).
- [9] Jumper, J., Evans, R., Pritzel, A. *et al.* Highly accurate protein structure prediction with alphafold. *Nature* **596**, 583–589 (2021).
- [10] Hinton, G., Vinyals, O. & Dean, J. Distilling the knowledge in a neural network (2015). URL <https://arxiv.org/abs/1503.02531>. arXiv:1503.02531.
- [11] Johnson, J., Alahi, A. & Fei-Fei, L. *Perceptual losses for real-time style transfer and super-resolution*, Vol. 9906 of *Lecture Notes in Computer Science*, 694–711 (Springer, 2016).
- [12] Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization (2017). URL <https://arxiv.org/abs/1412.6980>. arXiv:1412.6980.
- [13] Notin, P. *et al.* Oh, A. *et al.* (eds) *Proteingym: Large-scale benchmarks for protein fitness prediction and design*. (eds Oh, A. *et al.*) *Advances in Neural Information Processing Systems*, Vol. 36, 64331–64379 (Curran Associates, Inc., 2023). URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets_and_Benchmarks.pdf).
- [14] Lambert, T. J. *et al.* Fpbase: a community-editable fluorescent protein database. *Nature Methods* **19**, 169–170 (2022).
- [15] Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning* (MIT Press, 2006).
- [16] van der Maaten, L. & Hinton, G. Visualizing data using t-sne. *Journal of Machine Learning Research* **9**, 2579–2605 (2008).