

Generative models for protein design and evaluation

How to choose one generative model (stick) over the other options?



Scout Van den Bergh

Supervisor: Gaetan De Waele

Promotor: Prof. Dr. Willem Waegeman

Presentation Thesis first semester

Ghent University

Generative models for protein design and evaluation

Overview



01 Generative AI meets Proteins

02 Generative model

03 New evaluation metric

**Generative models for protein design
and evaluation**

Generative AI meets Proteins

Autoregressive Models

Probability of sampling the current residue depends on previously sampled residues

Generative models for protein design and evaluation

Generative AI meets Proteins

Autoregressive Models

Probability of sampling the current residue depends on previously sampled residues

Diffusion Models

Iterative denoising mechanism, learn to predict the noise of your data.

Generative models for protein design and evaluation

Generative AI meets Proteins

Autoregressive Models

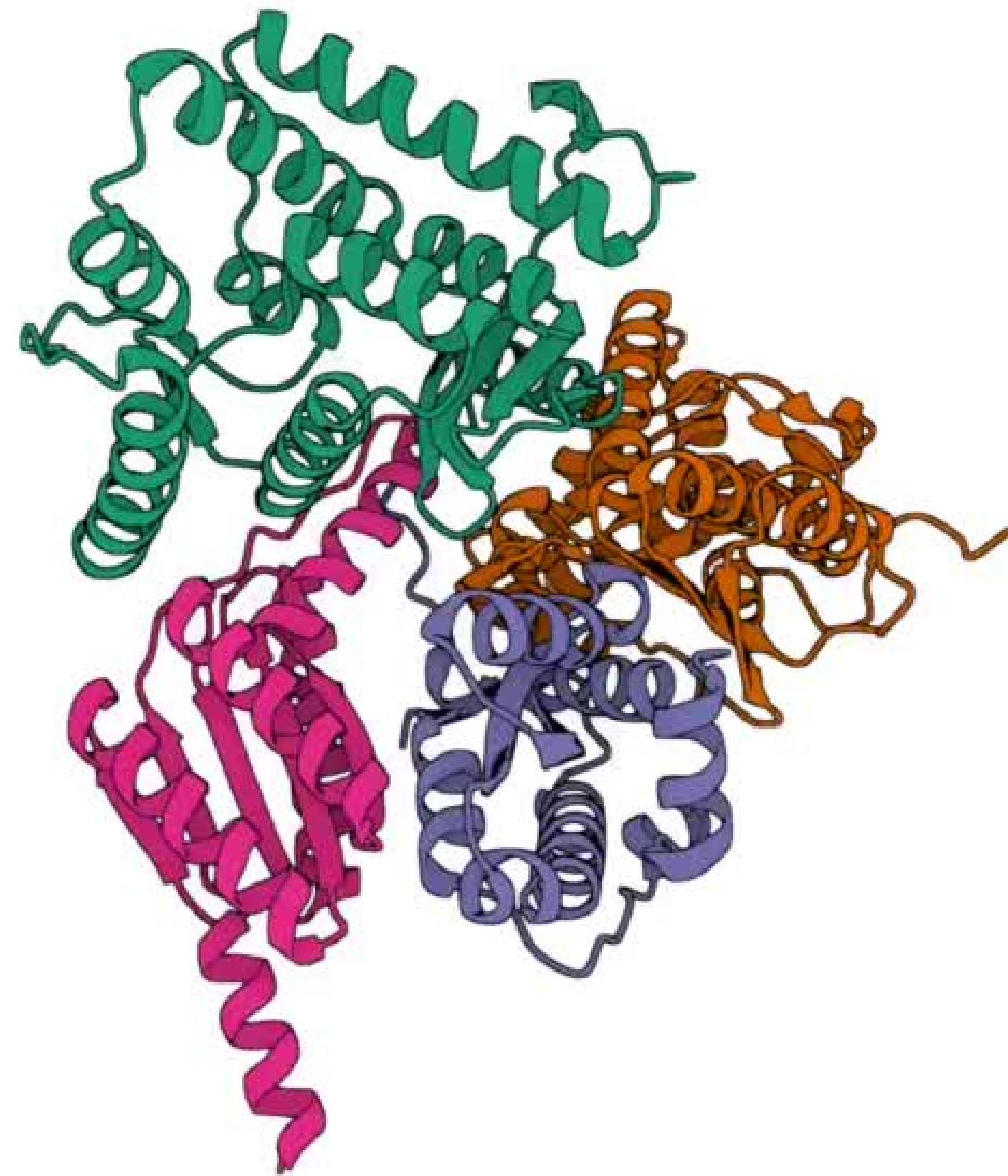
Probability of sampling the current residue depends on previously sampled residues

Diffusion Models

Iterative denoising mechanism, learn to predict the noise of your data.

Bi-directional masking protein language models

BERT-like approach.



**Generative models for protein design
and evaluation**

Generative model

Dataset

Subset of Protein Data Bank (PDB)
Contains > 400k protein sequences

Generative models for protein design and evaluation

Generative model

Dataset

Subset of Protein Data Bank (PDB)

Contains > 400k protein sequences

Variational Autoencoder

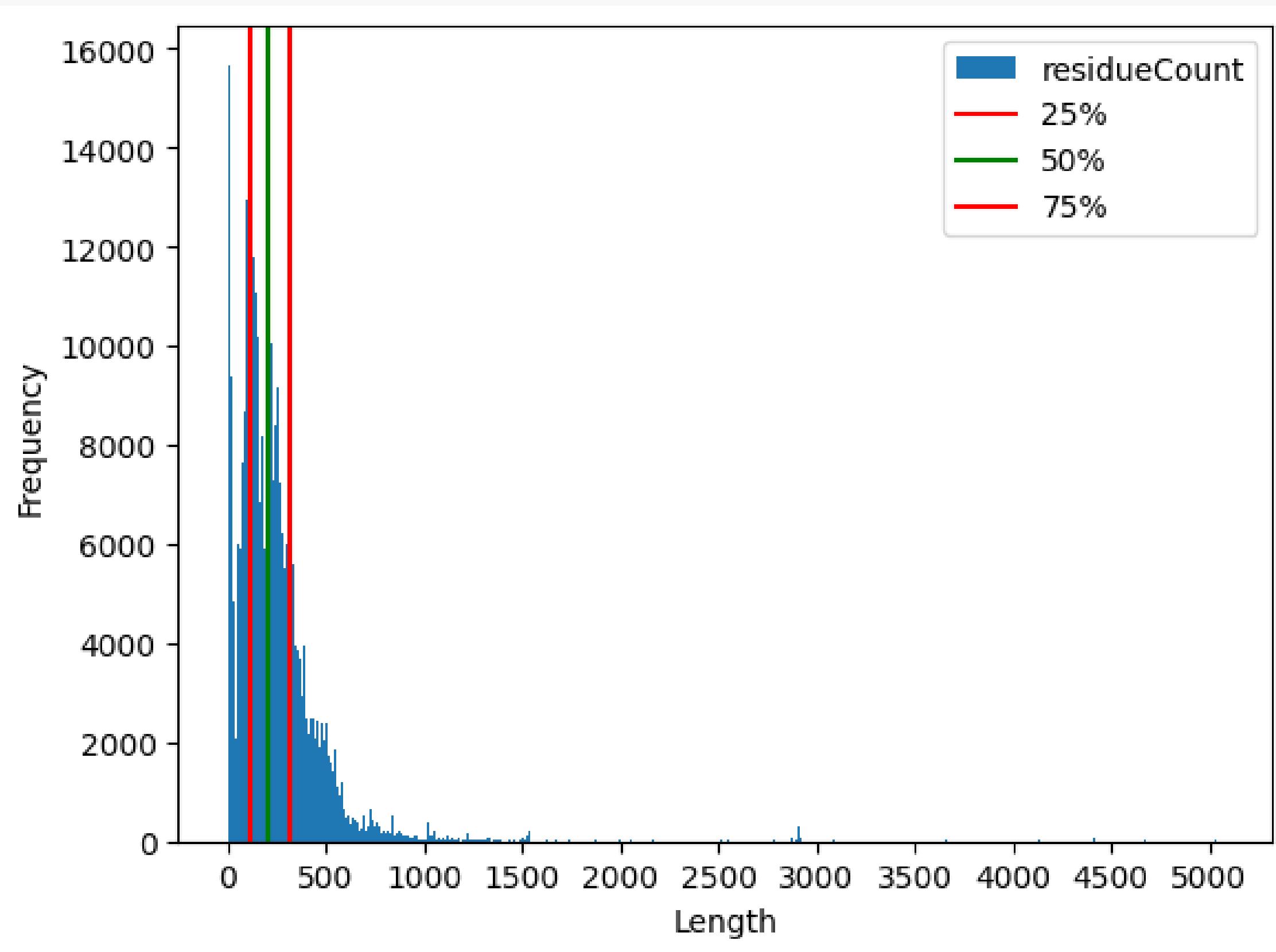
Learns μ of distribution

Learns σ of distribution

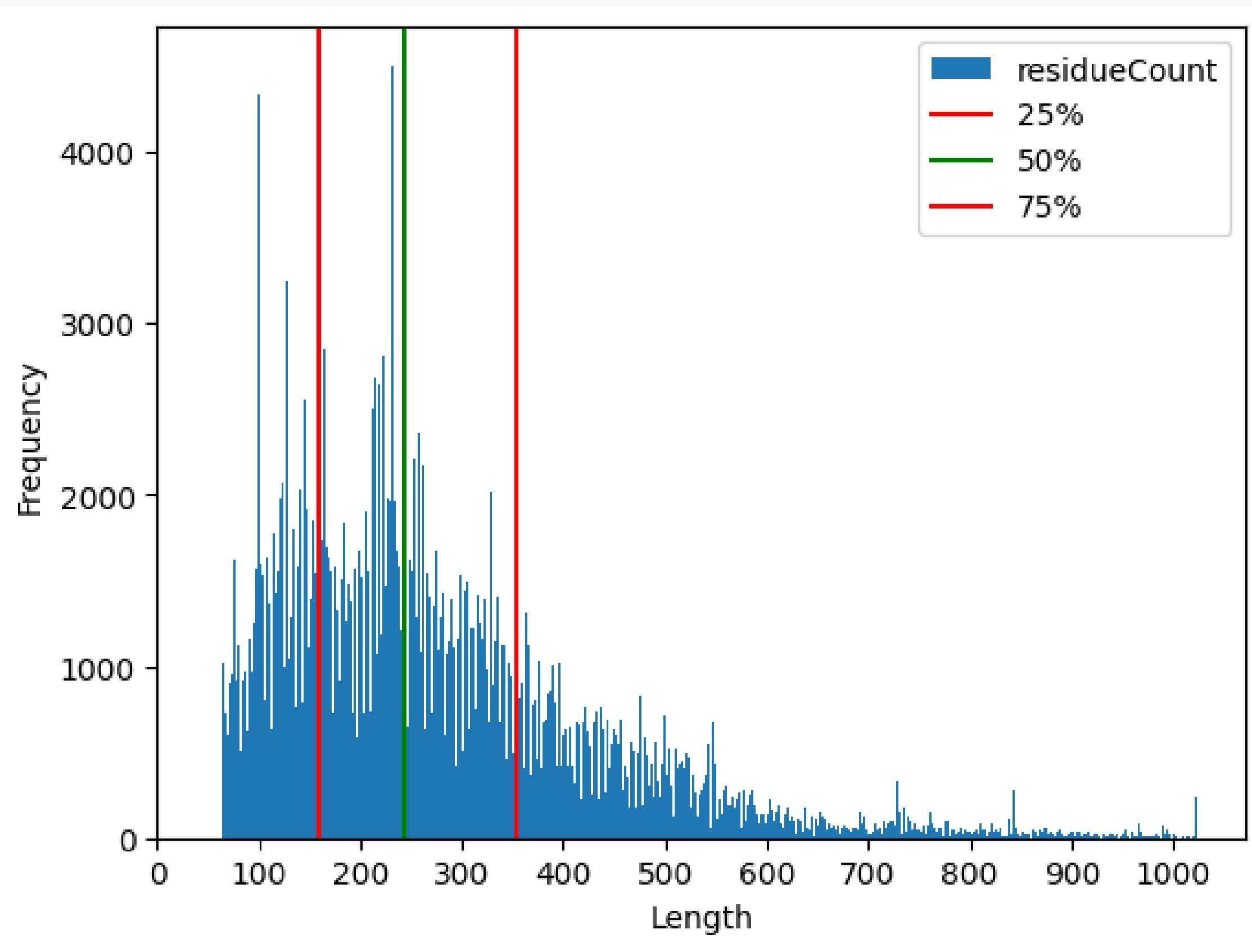
Reparameterization trick

Reconstruction loss + KL divergence

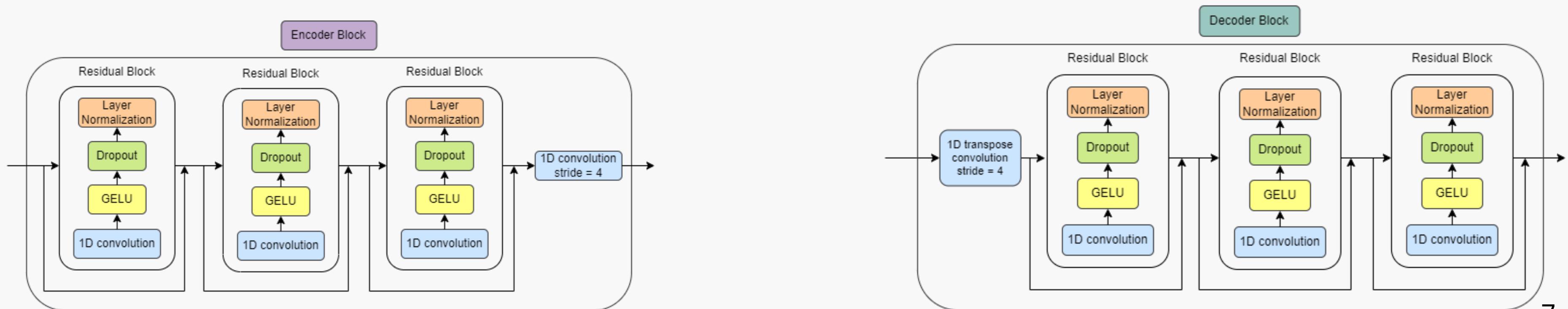
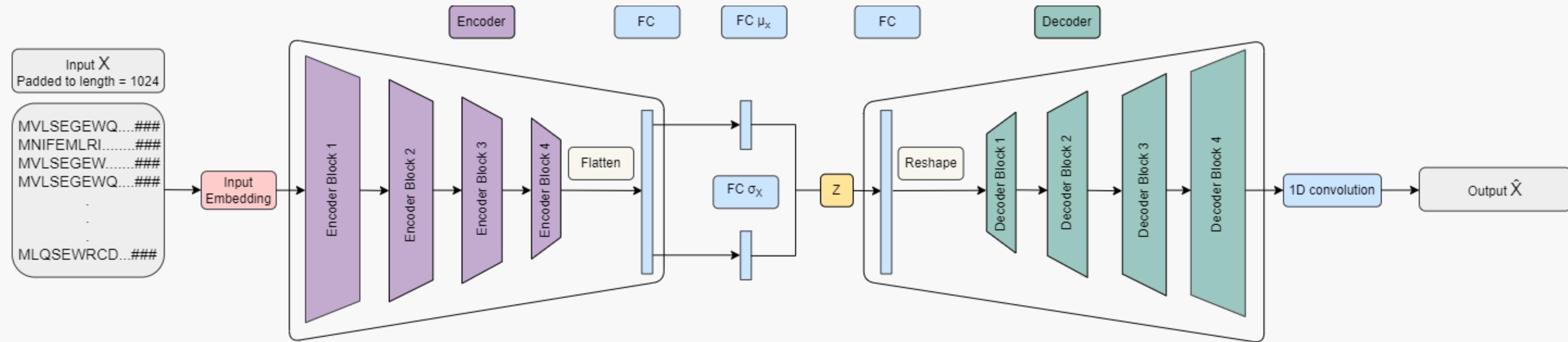
Distribution of dataset



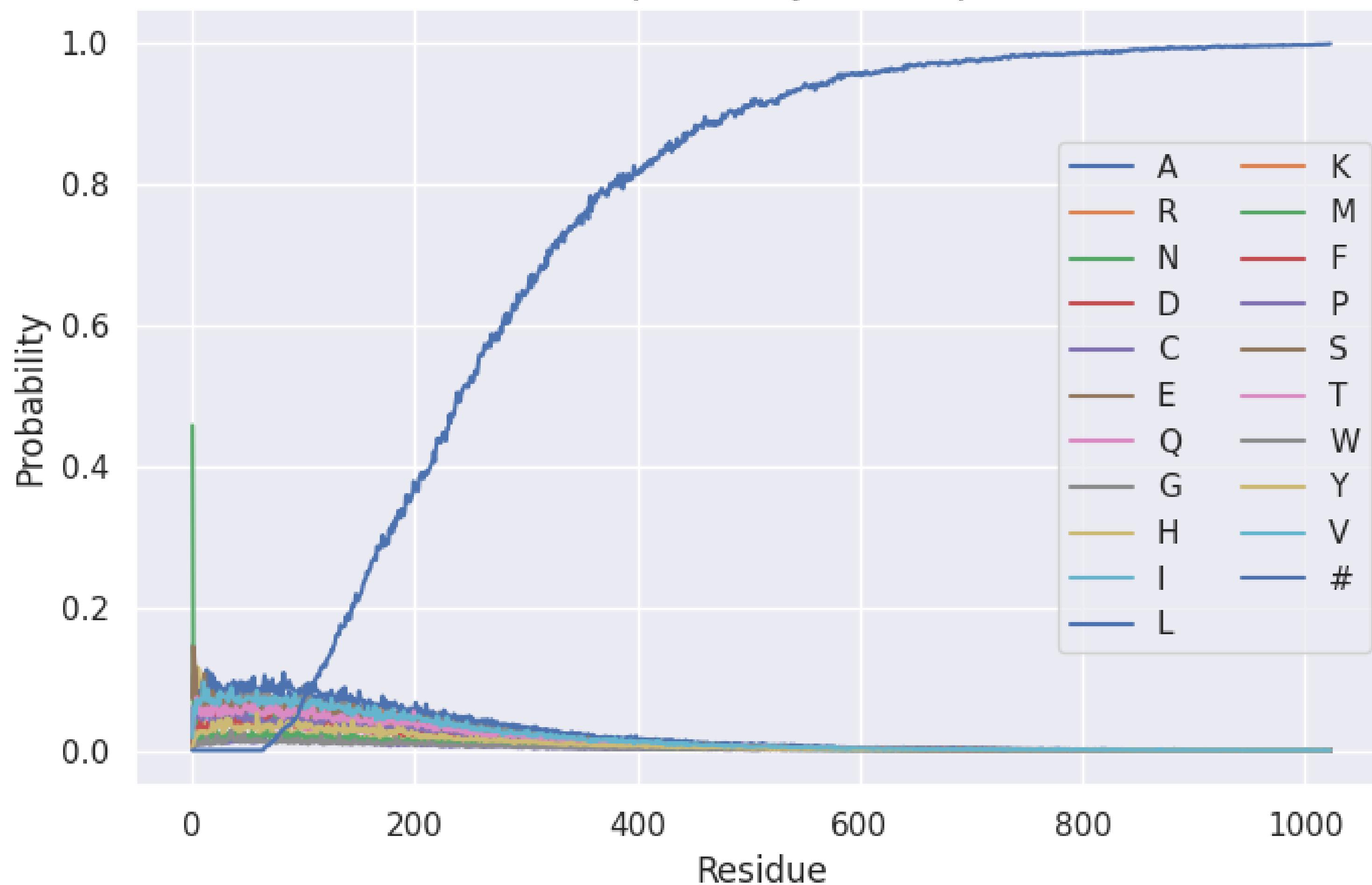
Distribution of dataset



VAE model



Amino acid probability at each position

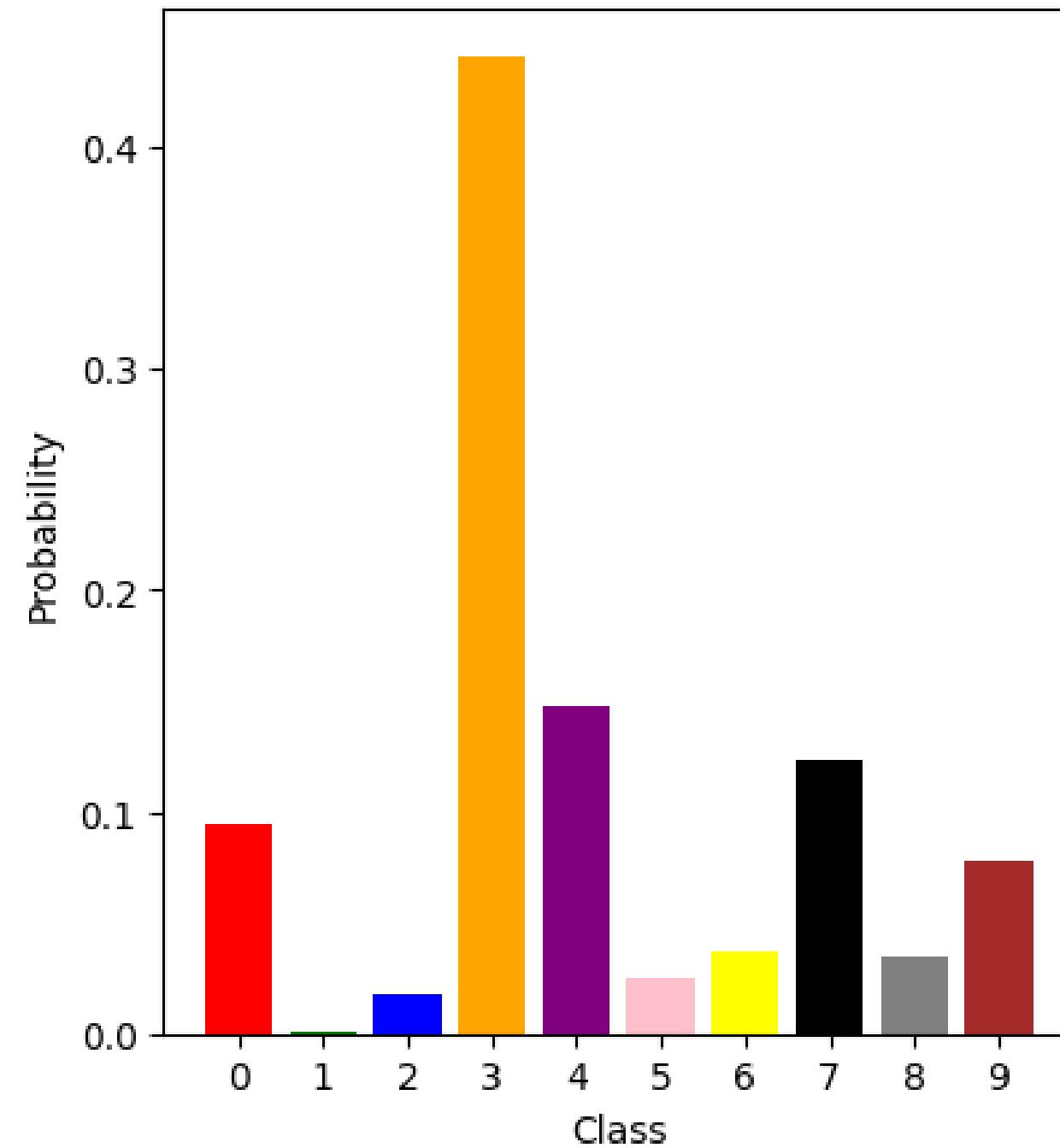


Softmax sampling using temperature

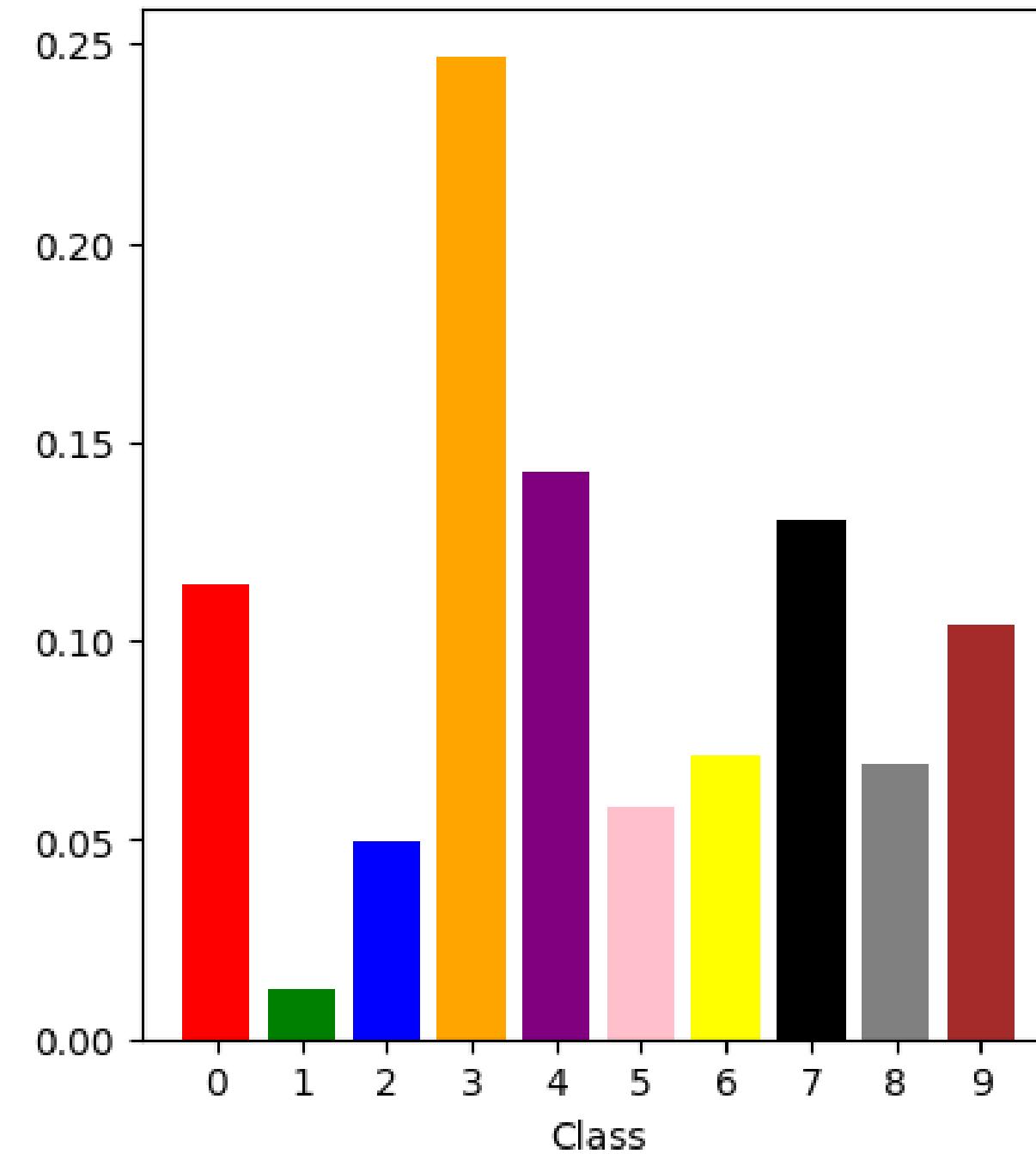
$$\text{Softmax}(y)_i = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

$$\text{Temp-Softmax}(y)_i = \frac{e^{y_i/T}}{\sum_j e^{y_j/T}}$$

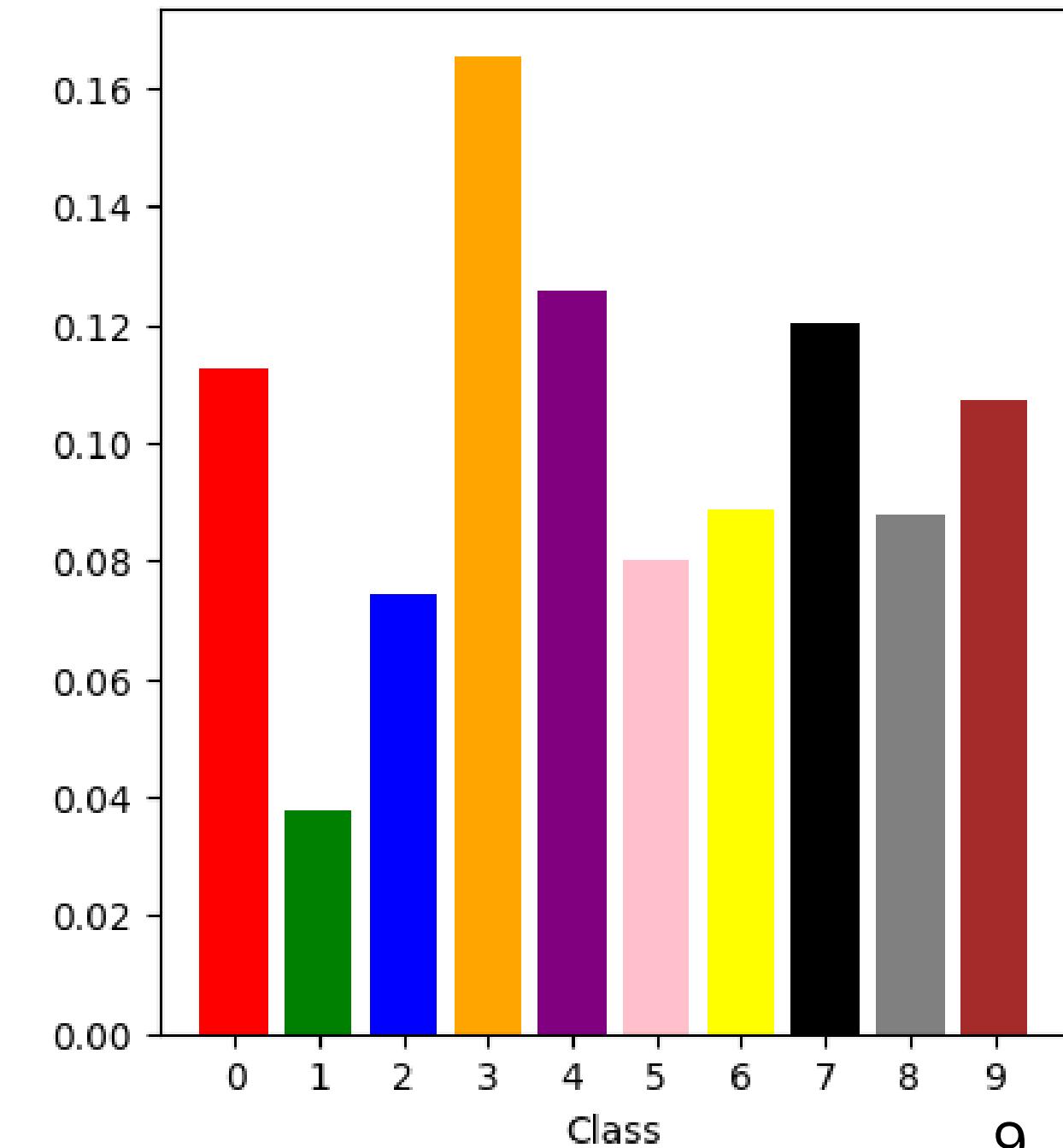
Softmax using temperature = 0.5



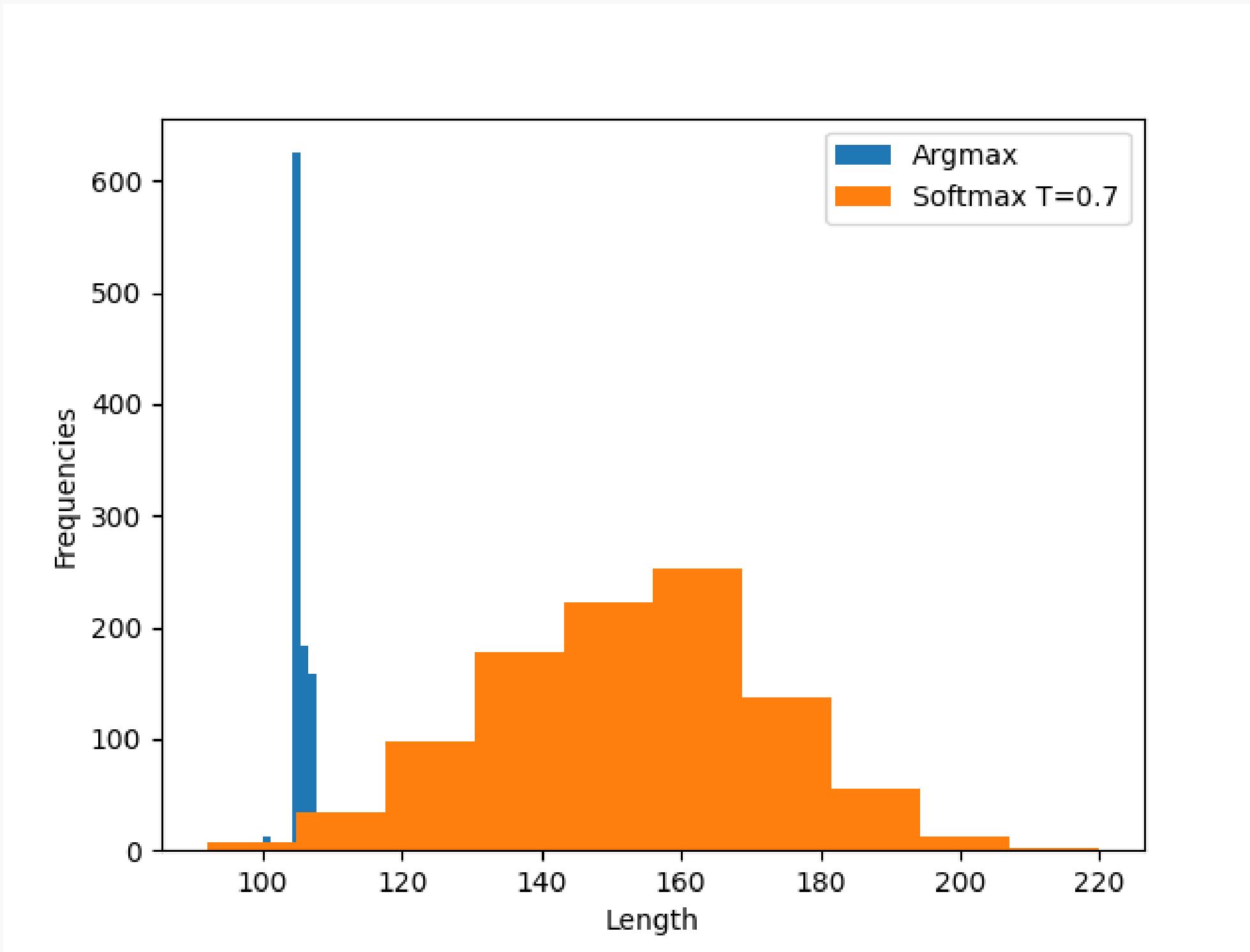
Softmax using temperature = 1



Softmax using temperature = 2

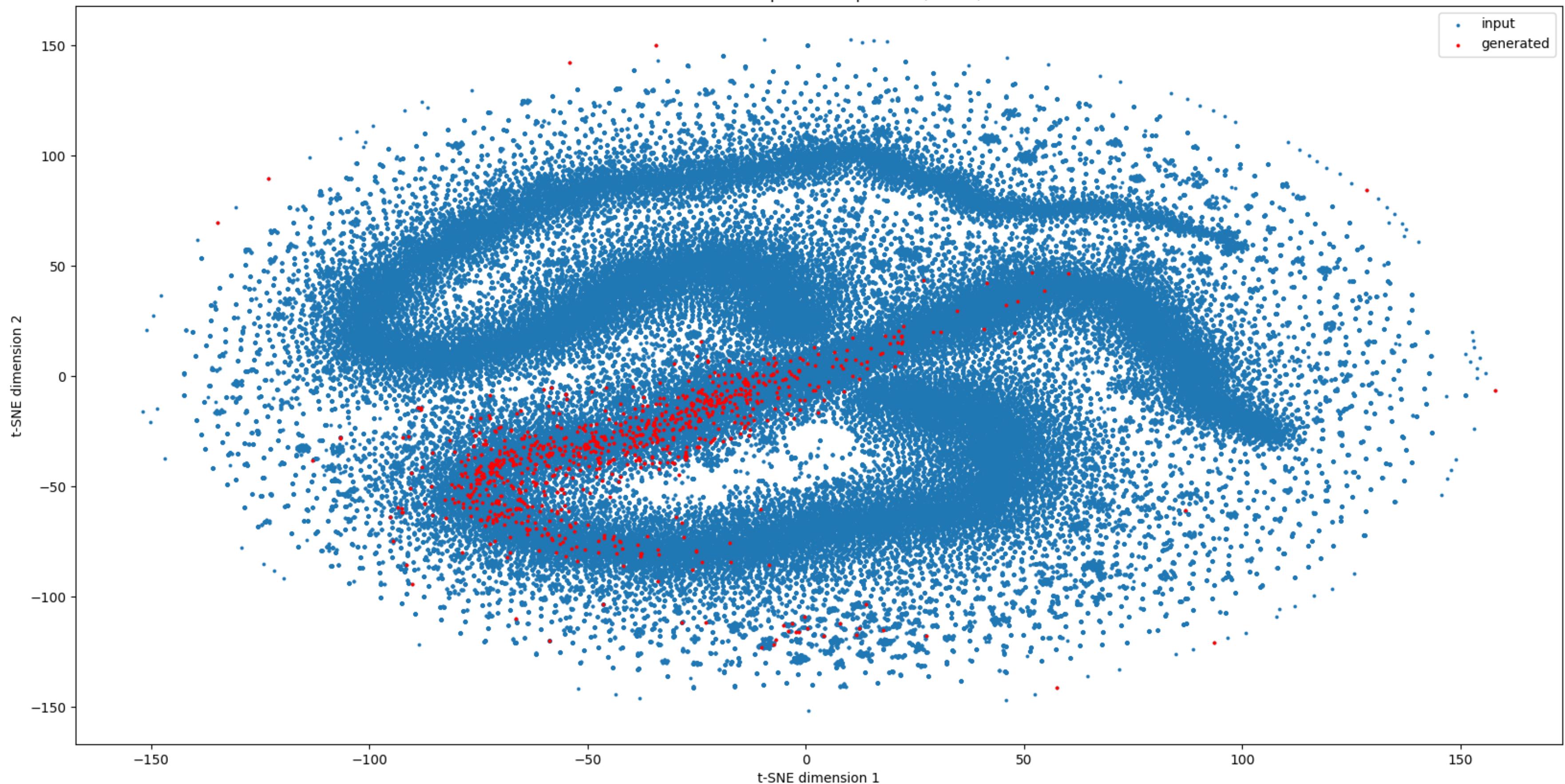


Histogram of length of 1000 generated proteins



t-SNE visualization

t-SNE of 311k protein sequences (T=0.7)



Next steps

1. More data --> UniRef (311k --> 43M)
2. Implement Transformers
 - a. Use rotary positional embeddings
 - b. Use SwiGLU activation function
3. Implement Gradient Descent: The Ultimate Optimizer
4. TBD

ESMFold

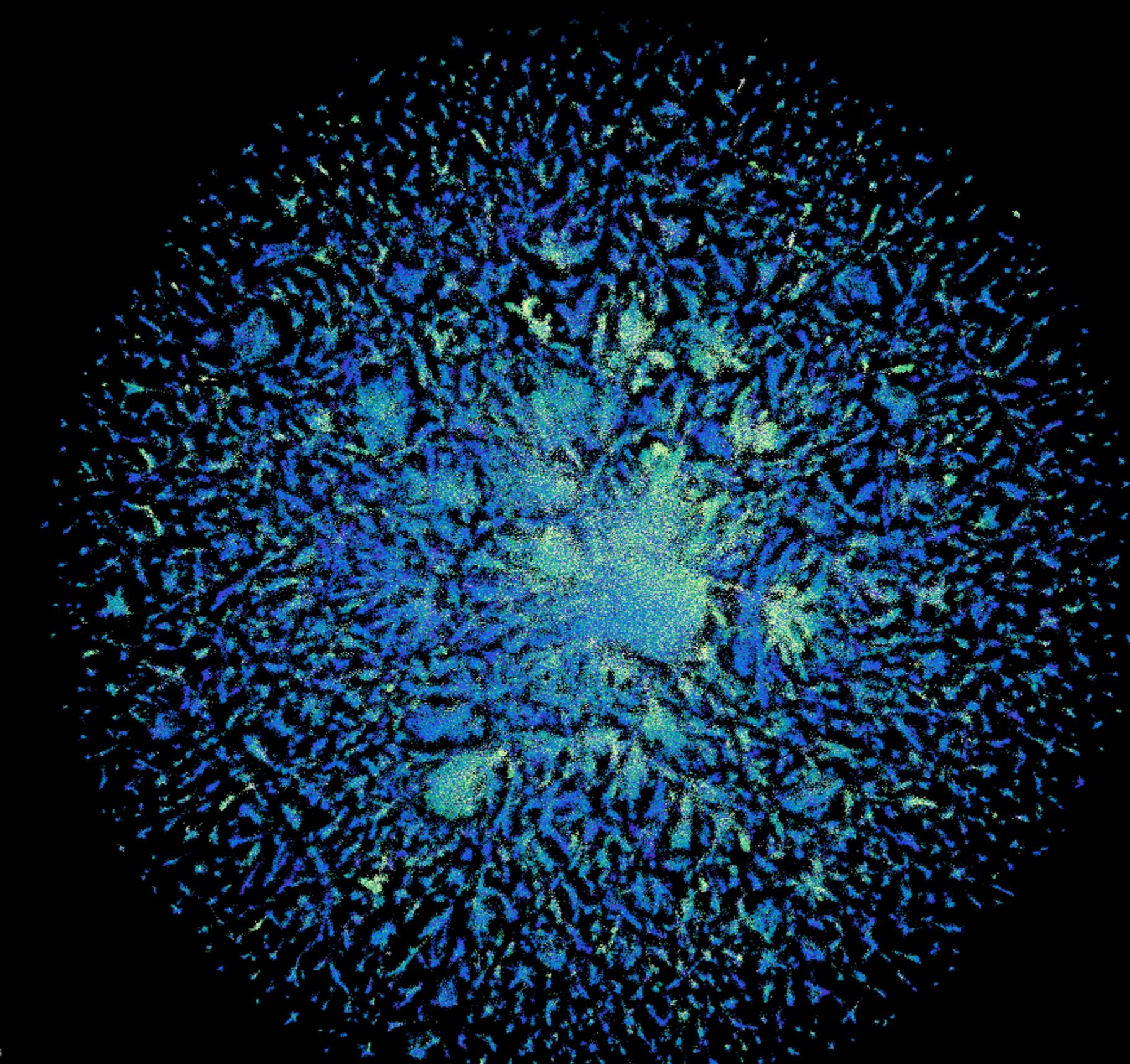
ESM Metagenomic Atlas

Explore

Resources ▾

About

Blog ▾



+ | - | Drag to pan
Scroll to zoom in/out
Click to select a protein

Exploring 1 million out of 617M proteins
Unknown Known

Meta AI

AlphaFold 2 vs ESMFold

214 683 829 structures

Requires MSA

~617 000 000 structures

Single sequence
60x faster



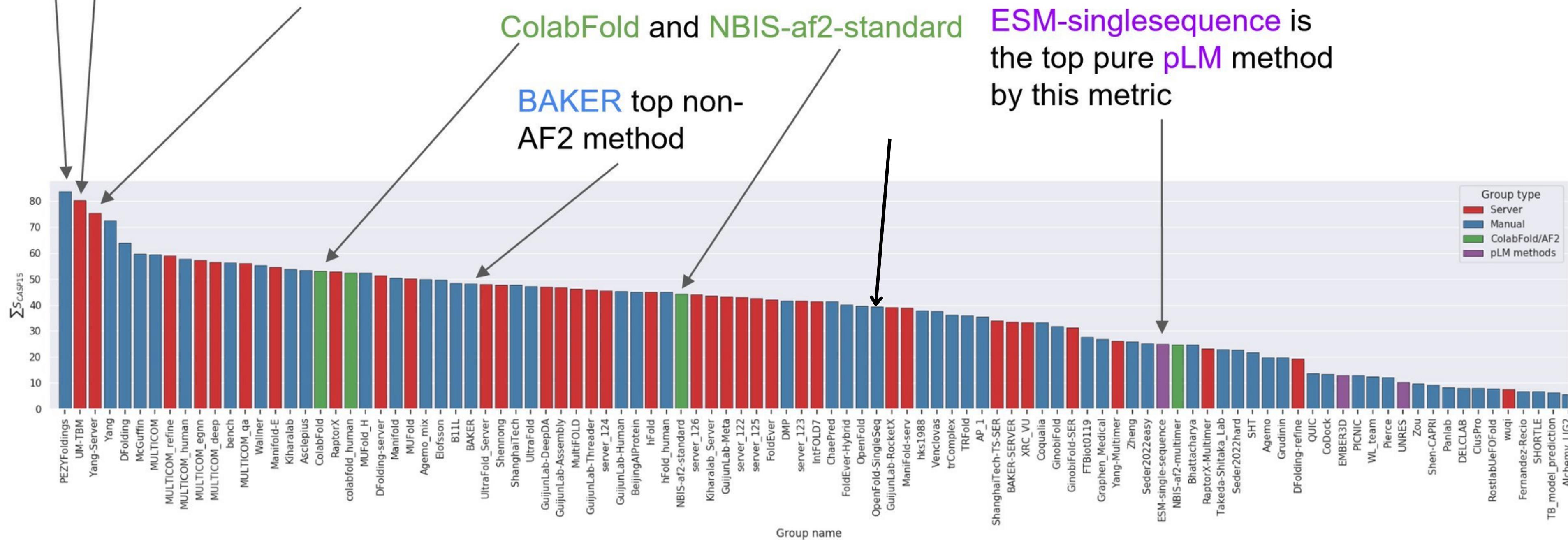
pIDDT: ■ Very low (<50) ■ Low (60) ■ OK (70) ■ Confident (80) ■ Very high (>90)



#1 PEZYFoldings AF2-based. Diverse MSAs.
Custom, fine-tuned AF2 refinement

#2 UM-TBM Diverse MSAs. Threading then AF2 predictions guide I-TASSER REMC

#3 Yang-Server Diverse MSAs. AF2 predictions fed to trRosettaX2



The CASP15 rankings

Generative models for protein design and evaluation

Evaluation

Current evaluation

- HMMR score
- DCA statistical energy score
- z-score
- pLDDT
- GDT
- PM
- TM

Generative models for protein design and evaluation

Evaluation

Current evaluation

- HMMR score
- DCA statistical energy score
- z-score
- pLDDT
- GDT
- PM
- TM

Fréchet Inception Distance

- Metric to evaluate performance of generative model.
- Measure of similarity between input and generated data.
- Fréchet Distance = distance between two probability distributions

Generative models for protein design and evaluation

Evaluation

Current evaluation

- HMMR score
- DCA statistical energy score
- z-score
- pLDDT
- GDT
- PM
- TM

Fréchet Inception Distance

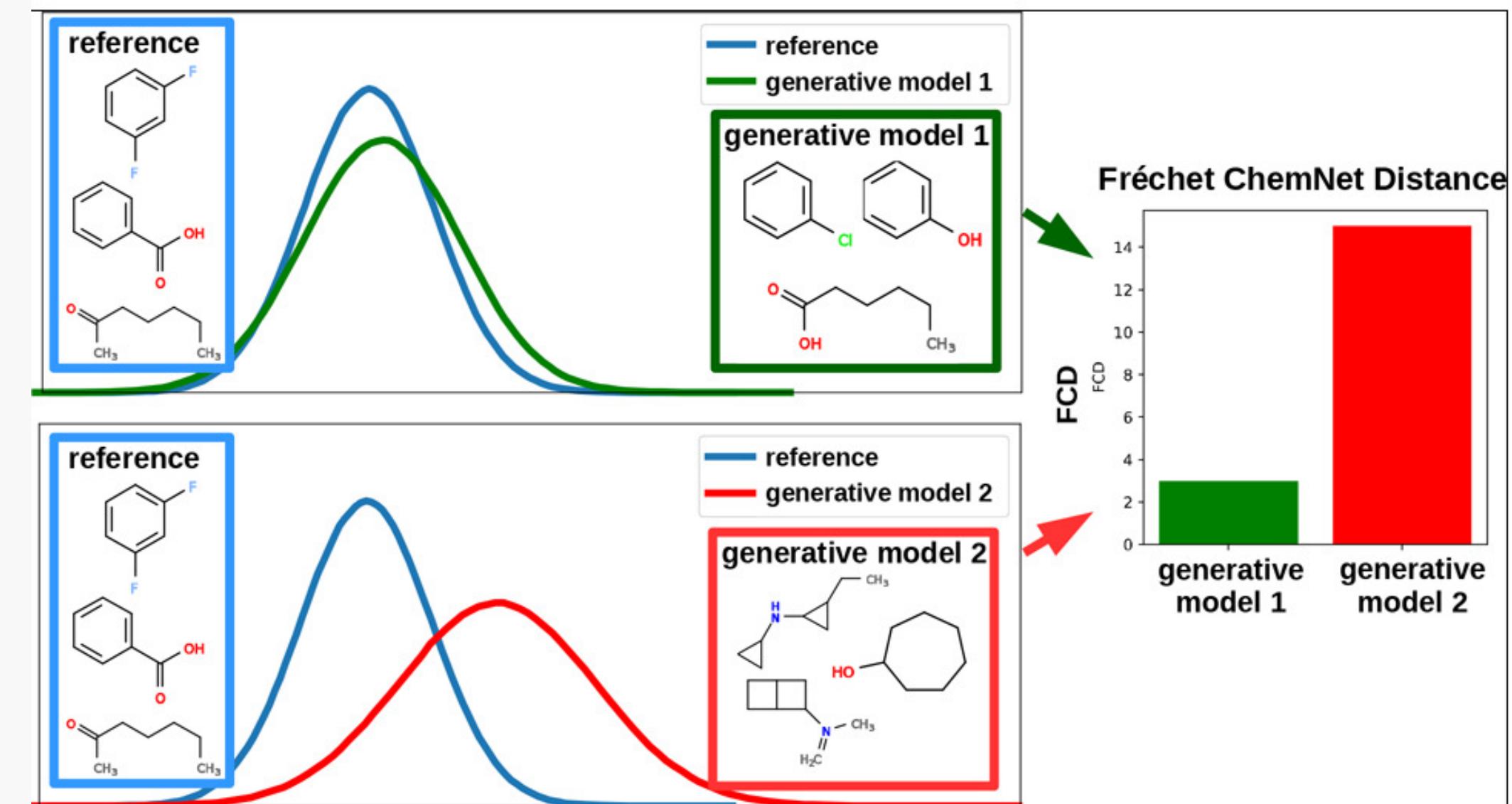
- Metric to evaluate performance of generative model.
- Measure of similarity between input and generated data.
- Fréchet Distance = distance between two probability distributions

Fréchet ESMFold Distance

Based on SOTA structure models

Fréchet Inception Distance

- Fréchet Inception Distance
- Fréchet Audio Distance
- Fréchet Video Distance
- Fréchet ChemNet Distance



Questions

Post-processing of generated sequence

Generated protein

Post-processed protein

MASEHAHRHDTVESEKFMVGVAEIEEYLILLPECLAHGLKMKTCKVSVADTLFIQPMTDHVESN
VGPKGGPQAEGTPKAGDTSLKNLTKYDDASNSETIHSKYIKFLLRILFPARNSSSAMIIALSGKGR
OLPPDKKSPLKVKGGLPGAD

