# Helix Synth: A Machine Learning Framework for Protein Secondary and Tertiary Structure Prediction

Allan

May 2025

## Abstract

Protein structure prediction is a cornerstone of computational biology, traditionally reliant on resource-intensive methods like X-ray crystallography. Helix Synth introduces a machine learning framework that leverages Convolutional Neural Networks (CNNs), Bidirectional Long Short-Term Memory networks (BiLSTM), Variational Autoencoders (VAEs), and Diffusion Models to predict both secondary (helix, beta-sheet, coil) and tertiary protein structures efficiently. Using real datasets from the Protein Data Bank (PDB) for secondary structure prediction, Helix Synth achieves an overall accuracy of 71.01%. For tertiary structure prediction, synthetic protein data is used to train a generative model, producing 5,003 novel structures with a mean squared error (MSE) of 0.0931 and mean absolute error (MAE) of 0.2523. This paper details the datasets, technical implementation, performance metrics, visualizations, ethical governance, and applications in drug discovery, mutation analysis, and synthetic biology.

## 1 Introduction

Protein structure prediction is pivotal for understanding biological processes and advancing biotechnology, yet traditional methods like X-ray crystallography and NMR spectroscopy are costly and time-consuming, often requiring months to determine a single structure. Helix Synth addresses these challenges by employing deep learning techniques, including CNNs, BiLSTM, VAEs, and Denoising Diffusion Probabilistic Models (DDPM), to predict secondary (helix, beta-sheet, coil) and tertiary protein structures with high accuracy and efficiency.

The framework first predicts secondary structures using a CNN-BiLSTM model, leveraging real protein datasets from the Protein Data Bank (PDB). It then extends to tertiary structure prediction through generative modeling with VAEs and diffusion models, using synthetic protein data to generate novel structures. This paper outlines the datasets, technical implementation, performance metrics, visualizations for both secondary and tertiary structure prediction, ethical governance, and potential applications in drug discovery, mutation analysis, and synthetic biology.

## 2 Core Objectives

Helix Synth aims to:

- Develop a deep learning model for accurate secondary structure prediction using real PDB datasets.

- Extend the framework with generative AI (VAEs and diffusion models) to synthesize novel tertiary structures using synthetic data.

- Establish an ethical governance model for responsible AI-driven biotechnology.

- Enable advancements in drug discovery, mutation analysis, synthetic biology, and distributed machine learning.

# 3 Datasets and Preprocessing

Helix Synth utilizes distinct datasets for secondary and tertiary structure prediction, reflecting the different requirements of each phase.

## 3.1 Secondary Structure Prediction: Real PDB Datasets

For secondary structure prediction, Helix Synth uses datasets from the Protein Data Bank (PDB), intersected with PISCES culling, spanning multiple years and filtering criteria:

- **2018 Dataset (ss_2018)**: 9,078 sequences, 25% identity cutoff, 2.0 Å resolution, sequence length $\geq$ 20.

- **2020 Datasets**: Updated through mid-2020 with a minimum sequence length of 40:

  - ss_2020_25_20: 7,320 sequences, 25% identity, 2.0 Å resolution.
  - ss_2020_25_25: 9,646 sequences, 25% identity, 2.5 Å resolution.
  - ss_2020_30_25: 13,406 sequences, 30% identity, 2.5 Å resolution.

- **2022 Datasets**: Updated through late 2022, also with a minimum length of 40:

  - ss_2022_25_20: 8,313 sequences, 25% identity, 2.0 Å resolution.
  - ss_2022_25_25: 10,931 sequences, 25% identity, 2.5 Å resolution.
  - ss_2022_30_25: 15,079 sequences, 30% identity, 2.5 Å resolution.

Relaxing the identity and resolution cutoffs (e.g., to 30% and 2.5 Å in ss_2022_30_25) increased the number of sequences by approximately 47% compared to the 2018 dataset, enhancing data diversity while maintaining quality. Sequence lengths were analyzed, showing a shift toward longer sequences (50–500 amino acids) in the 2022 datasets, particularly in ss_2022_30_25.

Preprocessing involved labeling proteins into Q3 states (H: Helix, E: Beta Sheet, C: Coil) using the Dictionary of Secondary Structure of Proteins (DSSP). Features were extracted using one-hot encoding and pretrained embeddings (e.g., ProtBERT, TAPE, ESM2), followed by tensor preparation with NumPy and Pandas, batching, and shuffling for GPU-optimized training.

## 3.2 Tertiary Structure Prediction: Synthetic Dataset

For tertiary structure prediction, Helix Synth generates a synthetic dataset of 1,000 protein structure samples, each with 400 dimensions, to simulate protein sequences. The generation process, performed on the CPU, involves:

- Creating a base sequence of length 100 (input_dim/4) using random values.

- Tiling the base sequence four times and adding noise (0.05 scale) and a sinusoidal perturbation (0.1 scale).

- Normalizing the resulting sequence to the range [0, 1].

- Saving each sample as a NumPy array (e.g., protein_0.npy).

The dataset is split into 70% training (700 samples), 15% validation (150 samples), and 15% testing (150 samples). A custom PyTorch Dataset class (`ProteinStructureDataset`) loads the data, and DataLoader is used with a batch size of 128 for training, validation, and testing.

# 4 Technical Breakdown

Helix Synth is structured into three phases: secondary structure prediction, tertiary structure generation via VAEs, and structure refinement using diffusion models.

## 4.1 Phase 1: Model Development (Secondary Structure Prediction)

### 4.1.1 Model Architecture

The secondary structure prediction model combines CNNs and BiLSTM to classify protein sequences into Q3 states, as shown in Table 1.

Table 1: Model architecture for secondary structure prediction.

| Component | Purpose | Reason |
|---|---|---|
| CNN | Feature Extraction | Captures local sequence patterns |
| BiLSTM | Sequence Learning | Captures long-range dependencies |
| Fully Connected | Classification | Maps features to Q3 states |
| Softmax | Probabilities | Assigns confidence scores |
| Adam Optimizer | Optimization | Fast, adaptive learning |
| Cross-Entropy | Loss Function | Suited for multi-class prediction |

### 4.1.2 Training Pipeline

Training occurs on Kaggle T4 GPUs with CUDA acceleration, using:

- **Extreme Garbage Collection**: Memory is cleared using `torch.cuda.empty_cache()`.

- **Batch Processing and Data Caching**: Optimizes latency and resource usage.

- **Epochs and Early Stopping**: 30 epochs with early stopping to prevent overfitting.

## 4.2 Phase 2: Generative Model - Variational Autoencoder (Tertiary Structure Prediction)

The VAE generates tertiary structures from synthetic sequences, with the architecture defined as:

- **Encoder**: Maps input sequences (400 dimensions) to a 32-dimensional latent space through two linear layers (400 to 256, 256 to 256) with BatchNorm, LeakyReLU (0.2), and Dropout (0.1), followed by separate linear layers for mean (`fc_mu`) and variance (`fc_var`).

- **Reparameterization**: Samples latent variables using the reparameterization trick.

- **Decoder**: Reconstructs sequences from the latent space through three linear layers (32 to 256, 256 to 256, 256 to 400) with BatchNorm, LeakyReLU (0.2), and Dropout (0.1).

- **Time Embedding**: Incorporates diffusion time steps via a two-layer MLP (1 to 32, 32 to 32) with SiLU activation.

Training uses a combined loss function:

$$\text{Loss} = \text{Reconstruction Loss (MSE)} + \beta \cdot \text{KL Divergence} + \text{Diffusion Loss (MSE)},$$

where $\beta = 1.0$. The model is trained for 75 epochs using the AdamW optimizer (learning rate 1e-4, weight decay 1e-5), a cosine annealing scheduler, and mixed precision training with GradScaler. Data preparation is CPU-based, with GPU acceleration for encoder-decoder processes.

## 4.3 Phase 3: Diffusion Model (Tertiary Structure Refinement)

Inspired by Denoising Diffusion Probabilistic Models (DDPM), the diffusion model refines VAE-generated sequences by iteratively improving 3D protein folds, enhancing tertiary structure accuracy. Noise is added to the input during training, and the model learns to predict and remove this noise, guided by the diffusion loss term.

# 5 Results Summary

## 5.1 Performance Metrics Overview

The performance of Helix Synth is summarized in Table 2.

Table 2: Summary of Helix Synth performance metrics.

| Metric | Value |
|---|---|
| *Secondary Structure Prediction (CNN-BiLSTM)* | |
| Overall Accuracy | 71.01% |
| H-Structure Accuracy | 76.21% |
| E-Structure Accuracy | 63.26% |
| C-Structure Accuracy | 70.92% |
| *Tertiary Structure Prediction (VAE + Diffusion)* | |
| Generated Proteins | 5,003 |
| Mean Squared Error (MSE) | 0.0931 |
| Mean Absolute Error (MAE) | 0.2523 |

## 5.2 Secondary Structure Prediction Results

The CNN-BiLSTM model achieves an overall accuracy of 71.01%, with specific accuracies of 76.21% (H), 63.26% (E), and 70.92% (C). Prediction confidence metrics are detailed in Table 3.

Table 3: Prediction confidence for secondary structures.

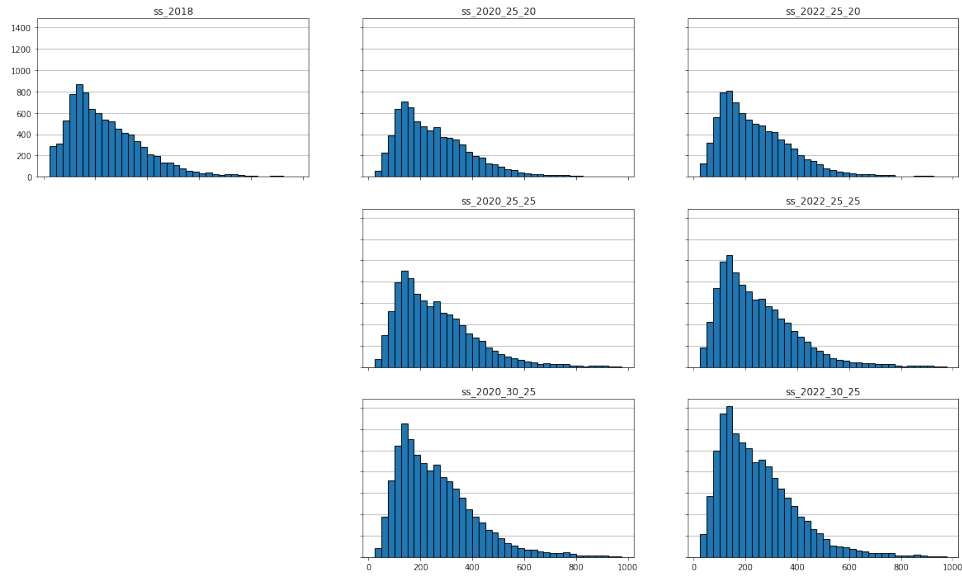| Structure Type | Accuracy | Mean | Std | Min | Max |
|---|---|---|---|---|---|
| H | 0.7621 | 0.8013 | 0.1763 | 0.3354 | 0.9998 |
| E | 0.6326 | 0.7272 | 0.1723 | 0.3355 | 0.9987 |
| C | 0.7092 | 0.6969 | 0.1511 | 0.3349 | 0.9970 |
| Overall | 0.7101 | – | – | – | – |

## 5.3 Secondary Structure Visualizations



Figure 1: Distribution of sequence lengths across datasets. This set of histograms compares sequence length distributions across datasets (ss_2018, ss_2020_*, ss_2022_*), with bins from 25 to 1,000. The 2022 datasets, especially ss_2022_30_25, show a shift toward longer sequences (50–500 amino acids), reflecting the relaxed cutoffs.
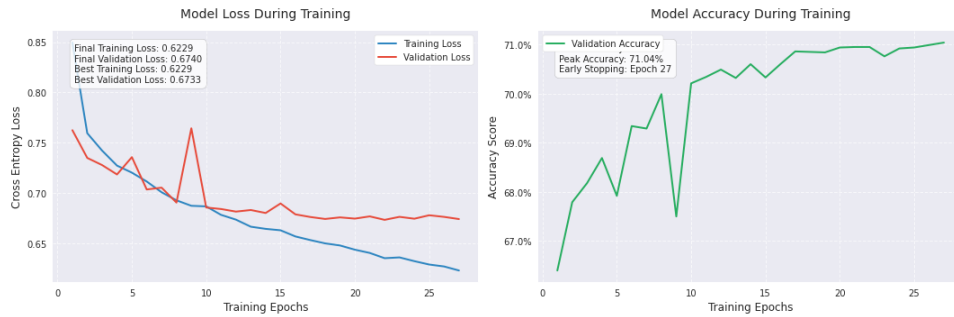


Figure 2: Model loss during training for secondary structure prediction. This plot displays the training and validation loss over 30 epochs. The loss curves decrease steadily, with early stopping preventing overfitting, aligning with the 71.01% accuracy.
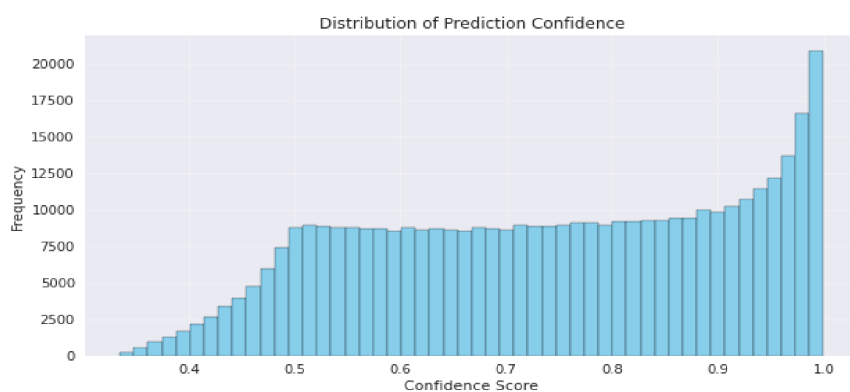
Figure 3: Distribution of prediction confidence for secondary structure predictions. This histogram illustrates prediction confidence scores for H, E, and C states, with means of 0.8013, 0.7272, and 0.6969, respectively, and tight distributions (standard deviations of 0.1763, 0.1723, 0.1511).

Figure 4: True and predicted secondary structure distribution. This bar plot compares true (DSSP) and predicted structure distributions, highlighting higher accuracy for helices (76.21%) compared to beta sheets (63.26%).

## 5.4 Tertiary Structure Prediction Results

The VAE generates 5,003 synthetic proteins, evaluated on a test set with an MSE of 0.0931 and MAE of 0.2523, indicating high reconstruction accuracy. Training history shows convergence of train and validation losses over 75 epochs, with final values of 5,083.7703 (train) and 2,978.4311 (validation).
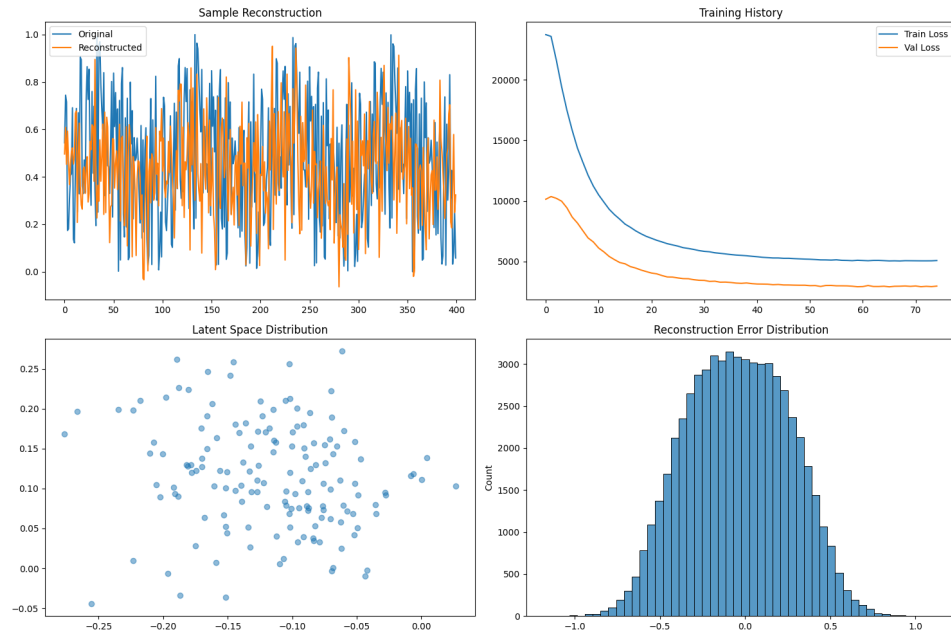
## 5.5 Tertiary Structure Visualizations



Figure 5: Master plot for tertiary structure prediction. Top left: Sample reconstruction comparing original and reconstructed sequences. Top right: Training history showing train and validation loss over 75 epochs. Bottom left: Latent space distribution visualized using a scatter plot. Bottom right: Reconstruction error distribution.

The master plot (Figure 5) integrates four components of the VAE's generative process:

- **Top Left: Sample Reconstruction**: Compares an original synthetic sequence (blue) with its VAE-reconstructed counterpart (orange) over 400 dimensions. High agreement indicates effective reconstruction, aligning with the low MSE (0.0931) and MAE (0.2523).

- **Top Right: Training History**: Plots training (blue) and validation (orange) losses over 75 epochs. Both losses decrease, converging to 5,083.7703 (train) and 2,978.4311 (validation), indicating stable training.

- **Bottom Left: Latent Space Distribution**: A scatter plot of the 32-dimensional latent space projected onto two dimensions, showing the distribution of synthetic proteins. Clustering reflects structural similarity.

- **Bottom Right: Reconstruction Error Distribution**: A histogram of reconstruction errors across the test set, with a mean MSE of 0.0931, showing most errors are near zero, confirming high reconstruction accuracy.
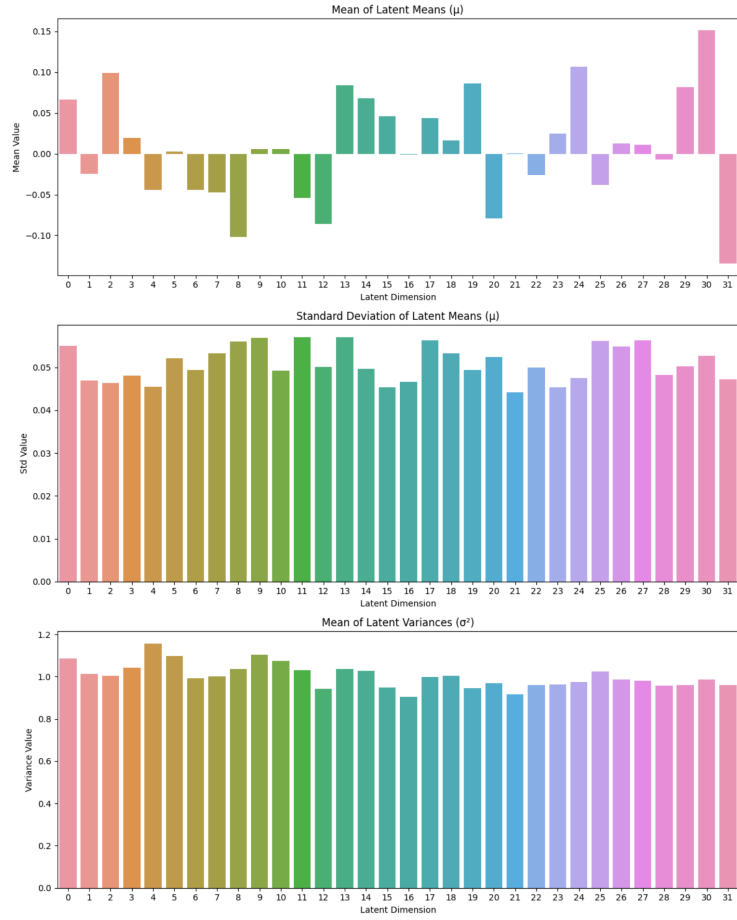
Figure 6: Mean of latent variances for VAE latent space. This plot displays the mean variances across the 32 latent dimensions, indicating well-defined features that support the generation of 5,003 proteins.
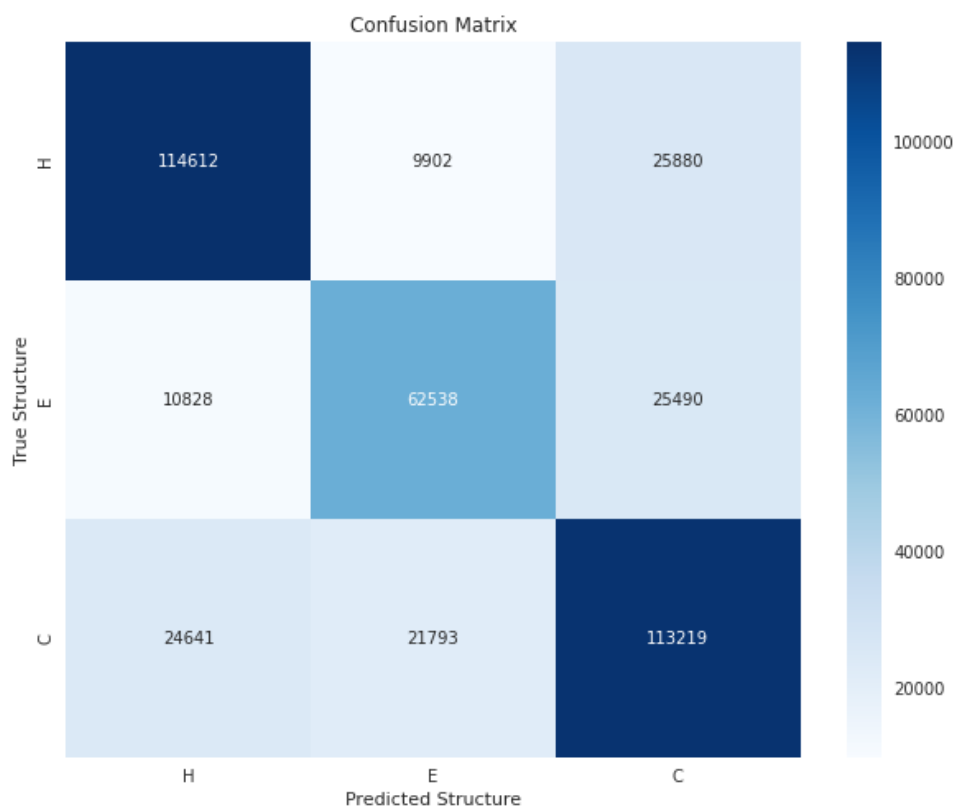
Figure 7: Correlation matrix of latent dimensions. This heatmap shows correlations among the 32 latent dimensions, providing insight into the VAE's ability to capture diverse structural properties.

# 6 Governance Model

Helix Synth adheres to an ethical governance model:

1. **Open-Access Development**: Initial models and datasets are public under the Apache 2.0 license, promoting collaboration.

2. **Independent Review**: External biologists validate synthetic proteins via lab testing.

3. **Controlled Release**: Core methods are open-sourced, with premium features access-controlled for industry labs.

4. **Regulatory Compliance**: Continuous monitoring ensures adherence to bioethical and biosecurity standards.

# 7 Future Applications

Helix Synth will impact:

- **Mutation Analysis**: Predicting structural effects of mutations for disease research.

- **Drug Discovery**: Modeling protein-ligand interactions for drug development.

- **Synthetic Biology**: Engineering novel proteins for industrial and medical uses.

- **Distributed Machine Learning**: Leveraging decentralized frameworks for scalable training.

# 8 Conclusion

Helix Synth advances protein structure prediction by integrating deep learning and generative modeling. It achieves 71.01% accuracy in secondary structure prediction using real PDB datasets and generates 5,003 synthetic tertiary structures with an MSE of 0.0931 and MAE of 0.2523. The reorganized visualizations in the results section provide clear insights into both the secondary structure prediction and tertiary structure generation capabilities. Ethical governance ensures responsible application, positioning Helix Synth to transform synthetic biology, drug discovery, and molecular design.

# 9 Next Steps

Next steps include:

- Deploying inference APIs for biotech labs.

- Validating synthetic structures experimentally.

- Expanding governance with regulatory bodies.