

# Latent Diffusion for Protein-Targeted Small Molecule Generation

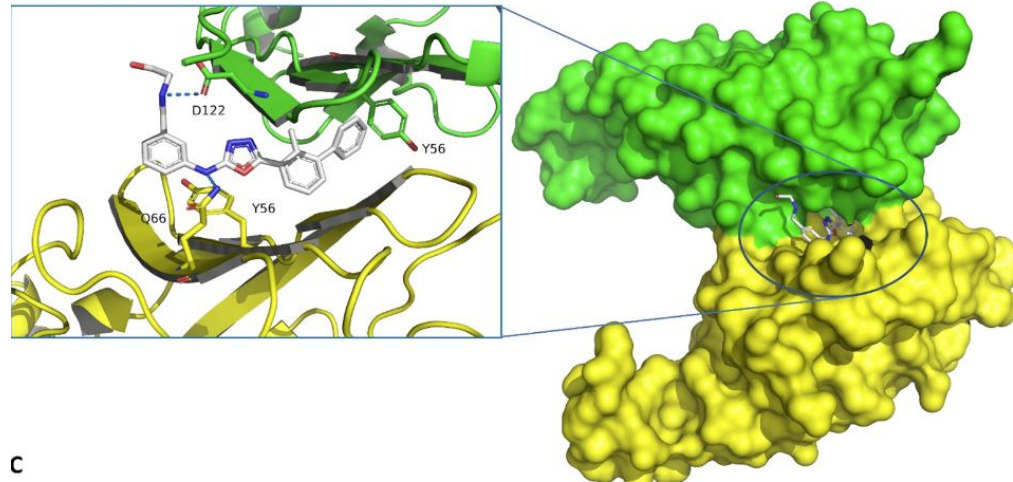
Srikar K., Rushil Y., Arnav S., Tyler K.

# Background



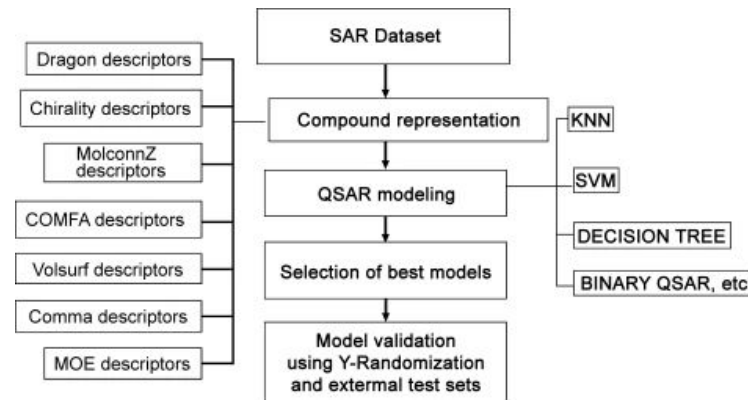
# Protein and Small Molecule Interactions

- modulate protein function via targeting active sites → **high specificity** and **selectivity** for therapeutics & **diversity of function**
- size (normally >900 Da) corresponds to **good diffusion properties** in vivo
- **high structural diversity** of small molecule candidates

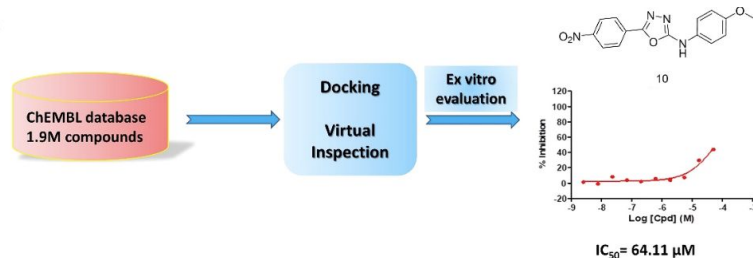


# Small Molecule Inhibitor Development

- **virtual screening** on existing databases (ChEMBL, DrugBank, etc.)
  - Ligand-based
  - Structure-based (docking & other simulations)
- **high-throughput screening** — high volume in vitro testing of candidates
- **QSAR** models — prediction of ligand characteristics from structural/computed features



A



# Current Limitations

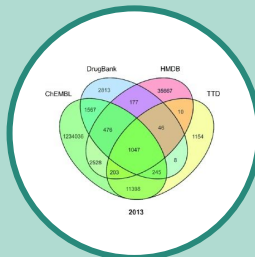


# Challenges in Small Molecule Discovery

- current methods are **limited in scope** to already-discovered & highly-adjacent compounds
  - ChEMBL: ~2.4 million
  - DrugBank: ~16.5K
- simulations can be **expensive** and **inaccurate**
- expensive and **impractical** to do in vitro analysis of large # of chemicals

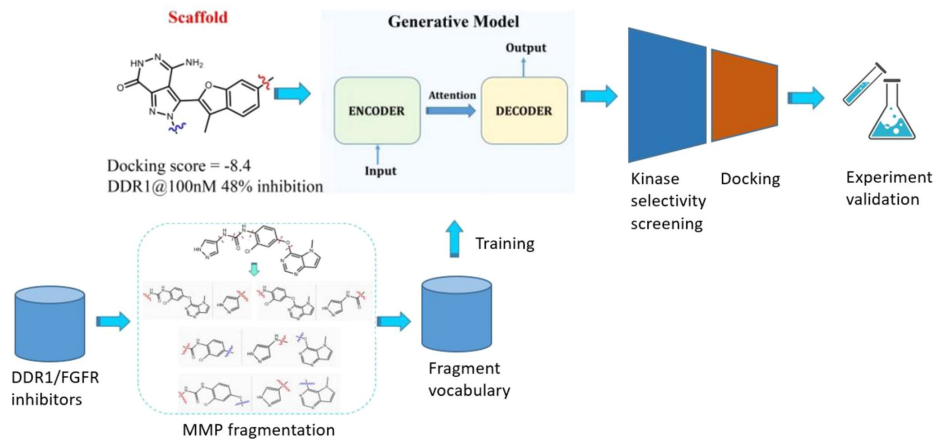
RO5-space ~  $10^{60}$  possible molecules

on-demand synthesis ~  $10^{15}$  possible molecules

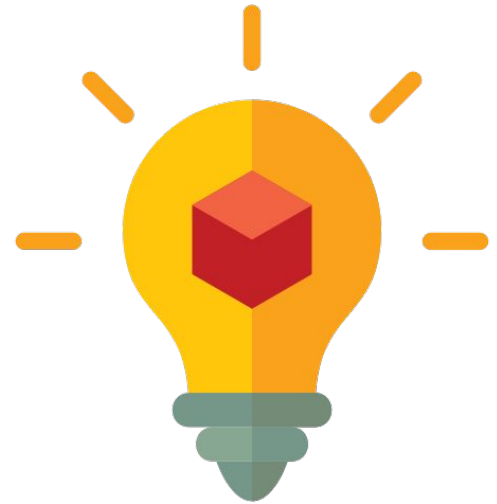


# Generative Models for Discovery

- generative models **not limited** to already-explored drug space → new chemistry & increased diversity
- Example generative model types include **RNNs**, **GANs**, and **VAEs**
- turn toward **NLP approaches** (SMILES, sequence info) → transformer-based architecture & other approaches



# Model Goals & Novelty





# Model Goals

## Simple Terms

Input Protein Sequence, output SMILES of protein-binding small molecule

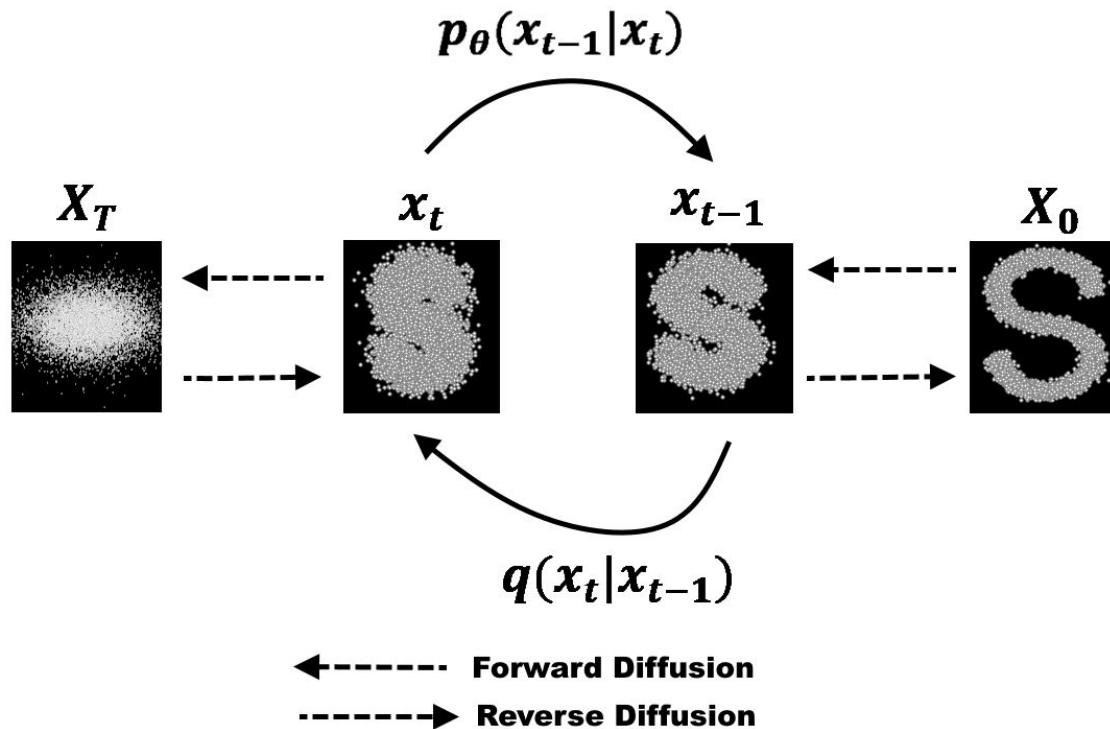
## Academic Speak

Generation molecular embeddings directly from protein embeddings via **latent sequence diffusion**, thereby creating a novel link between protein sequences and potential interacting molecules.

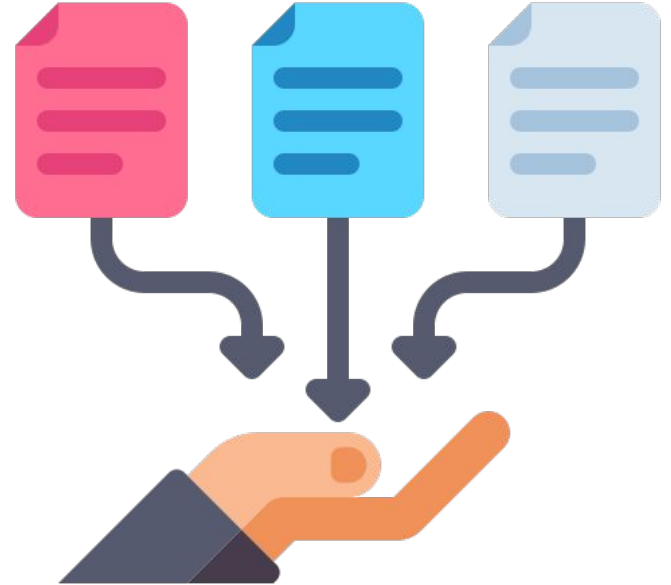
# Our Novelty

- **Innovative Approach to Molecule Generation:**
  - Unlike traditional methods that focus on molecular structures, our model explores latent sequence diffusion, offering a new perspective in molecule generation.
  - To our knowledge, the first of its kind in protein-targeted small molecule research.

# What is Diffusion?



# Data Curation



# BindingDB

- **Components**

- High-Affinity Protein Interactions: Proteins with ligand interactions stronger than 100 nM IC50/Kd/EC50.
- Ligand Criteria: Compounds with PubChem CIDs, representable by SMILES, <1000 Da, linked to Uniprot IDs.
- Protein Characteristics: Sequences ranging from 80 to 1000 amino acids, functioning as monomers

- **Selection**

- A dataset of 5000 protein, molecule pairs were chosen at random
- Train: Test: Validation Split → 70:15:15



# Improvements to Data Curation

- **Refine via clustering**

- Refine the dataset by sequence clustering at the 30% identity level using the tool Mmseqs2, aiming to cover a diverse range of sequences for training, validation, and testing.

- **Validate via PDB**

- The validation datasets should be carefully assembled by pairing PDB protein-ligand co-crystal structures with BindingDB entries.



# ESM-2 + Uni-Mol Embeddings

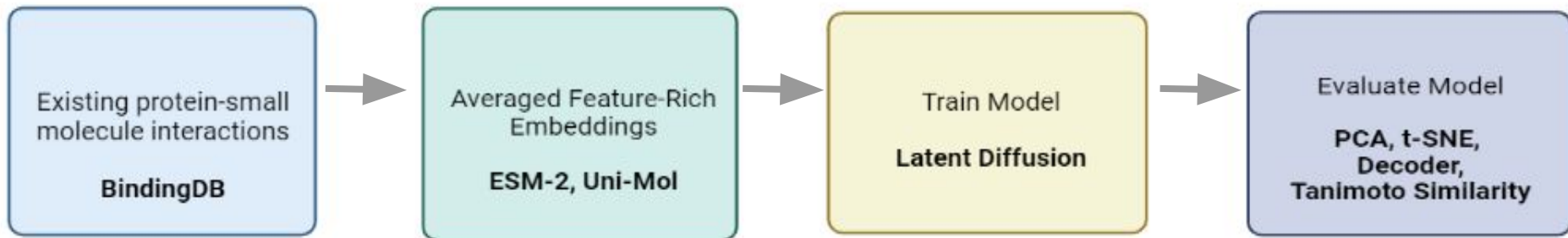
- **ESM-2:**

- Condenses High-Dimensional Data: It effectively reduces the complexity of protein sequences into a single embedding dimension shape vector, preserving essential information in a more manageable form.
- Captures Sequential Information: ESM-2's embeddings are adept at encapsulating the sequential nature of protein sequences, which is crucial for understanding the structural and functional aspects of proteins. However, we used *average* ESM-2 embeddings.

- **Uni-Mol:**

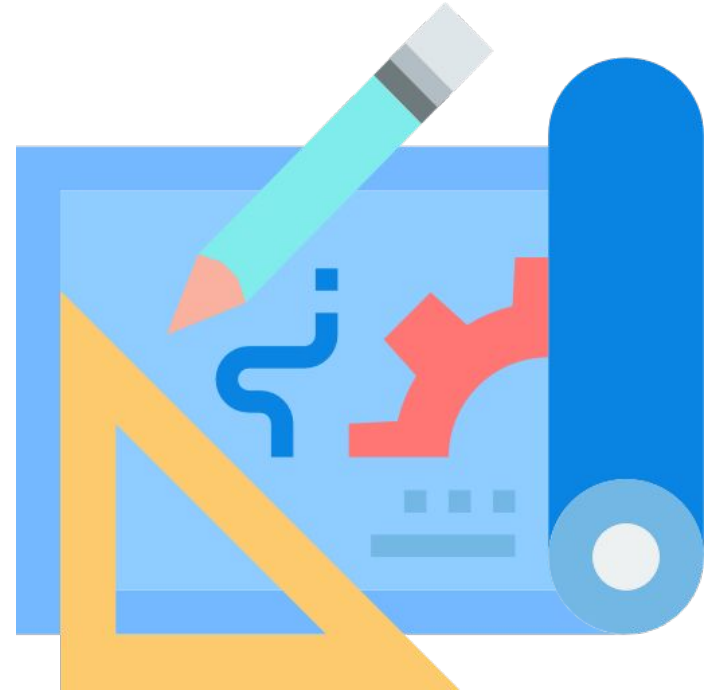
- Averaging of Tokens: Uni-Mol uses an approach that averages the embeddings of SMILES tokens, which are representations of molecular structures, effectively capturing the chemical characteristics of molecules.
- Reduction to Single Vector: The model condenses the entire sequence and SMILES tokenized sequence length into a single embedding dimension shape vector, efficiently representing complex molecular information.

# Workflow Diagram



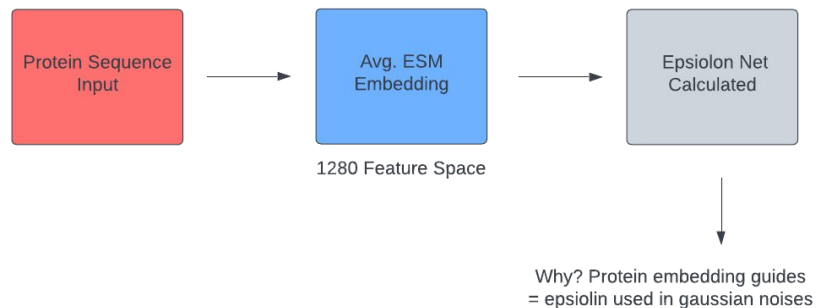


# Model Architecture

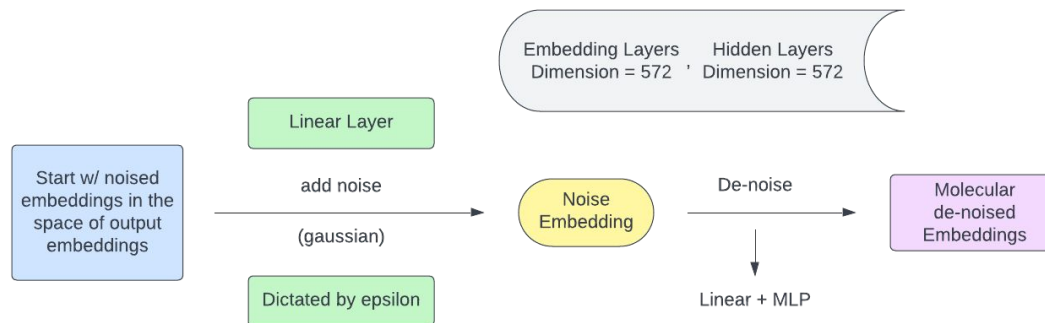


# Model Diagram

1



2

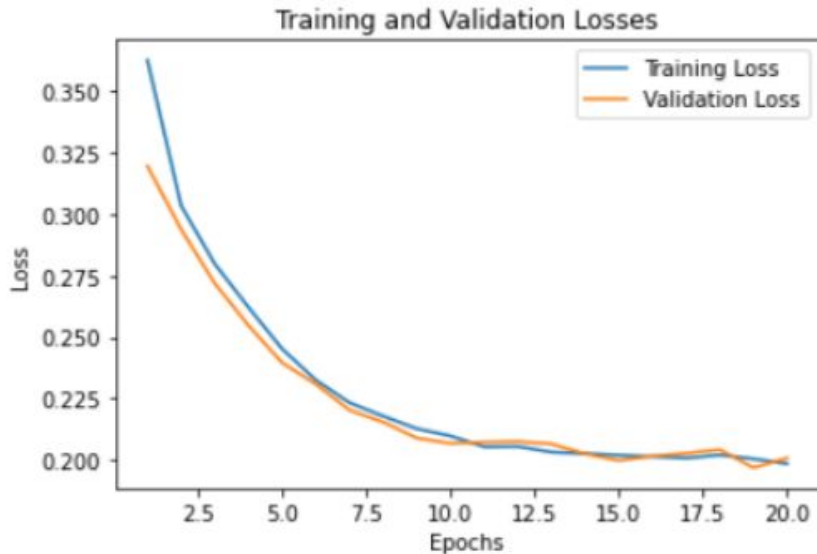


# Results



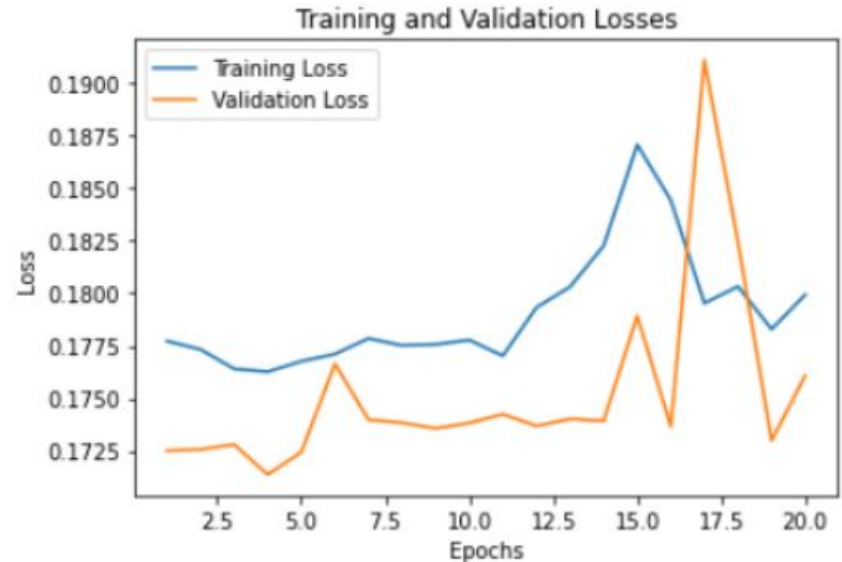
# Model Loss (L1 Loss)

Finalized Model (Linear Epsilon Net)



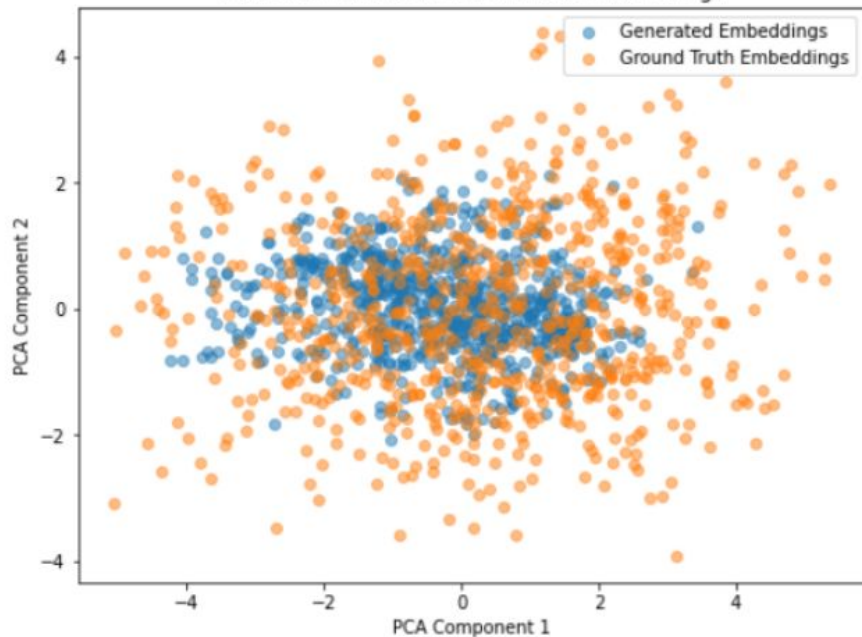
**Important Note:** Denoising architecture did not change model loss performance significantly.

One of tested variations (RNN Epsilon Net)

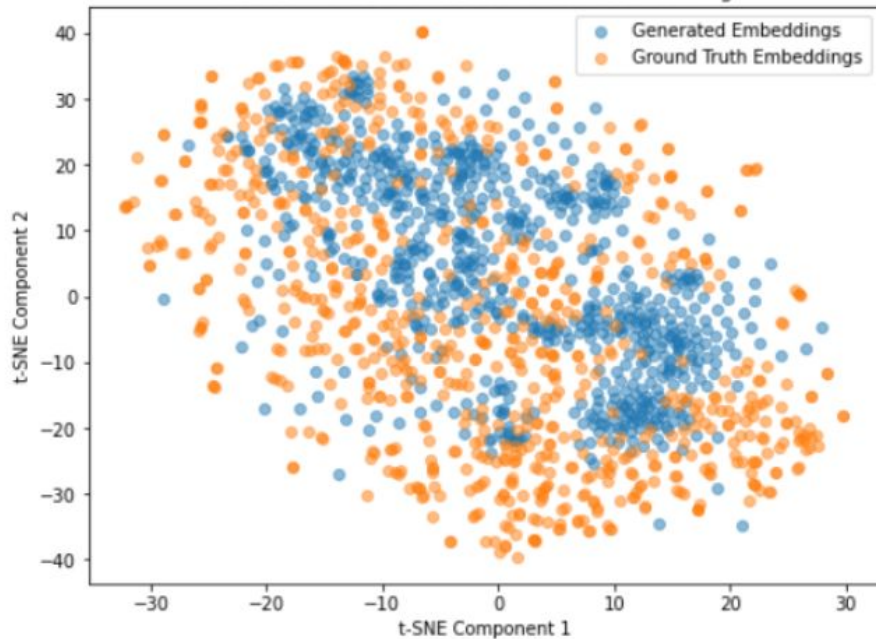


# Generated vs. Ground Truth Latent Visualization

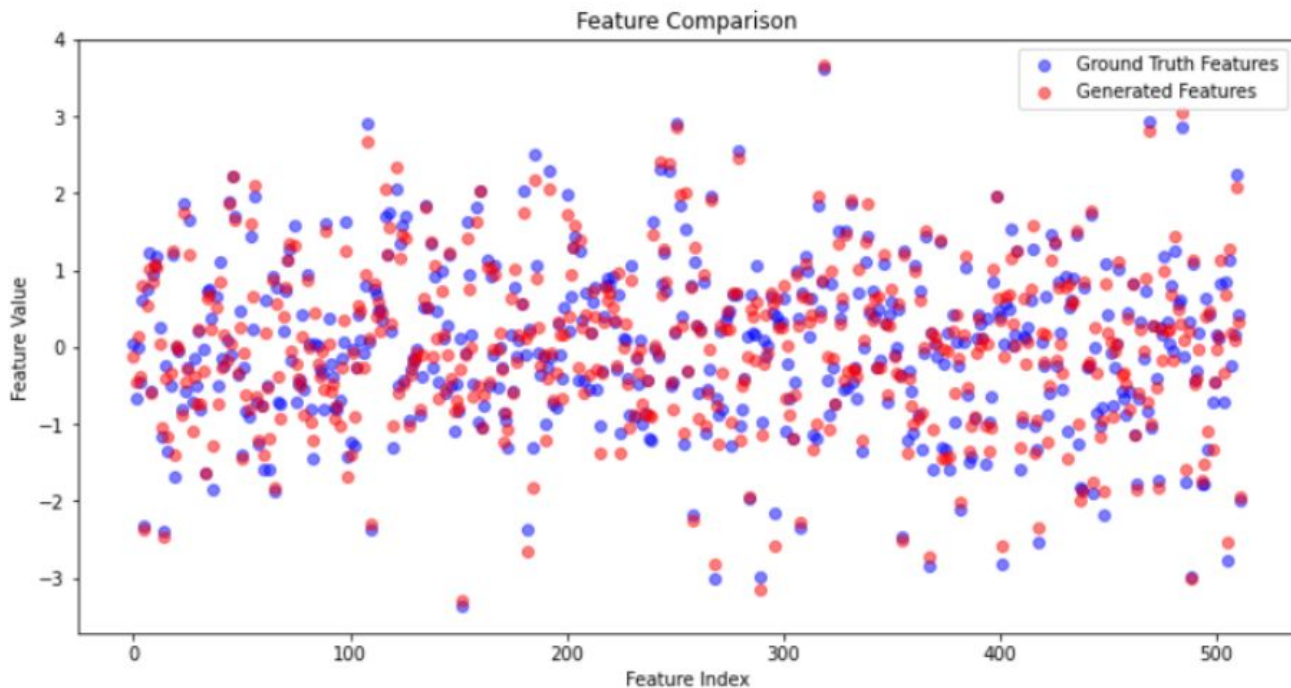
PCA of Generated vs Ground Truth Embeddings



t-SNE of Generated vs Ground Truth Embeddings



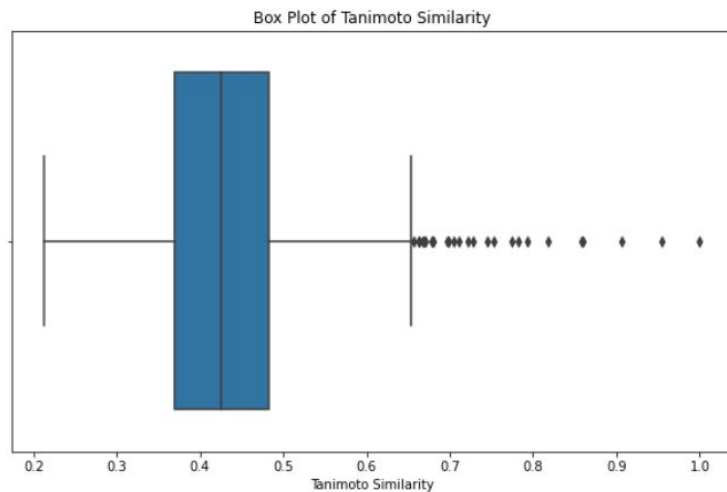
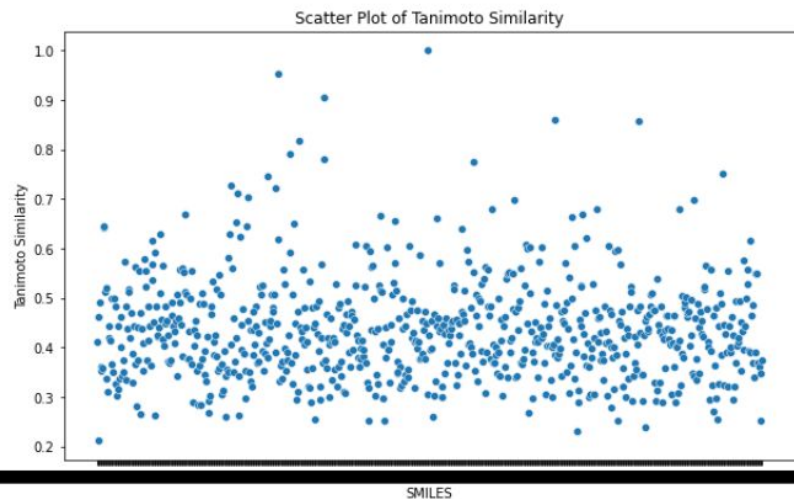
# • One Sample Generation (Expanded Chemical Space) •



# Decoded SMILES Output Comparisons (Look-up Table)

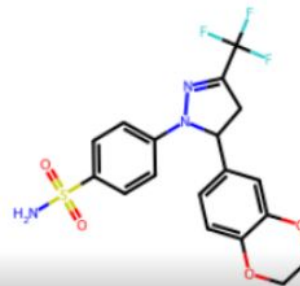
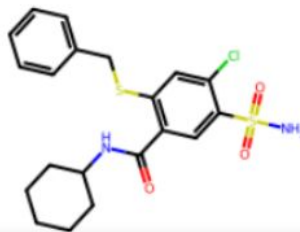
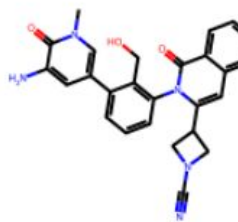
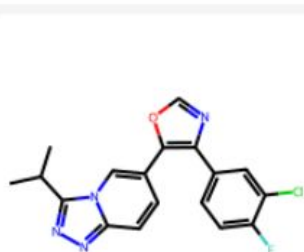
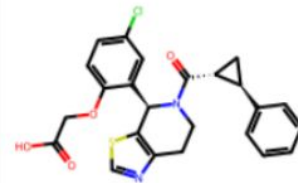
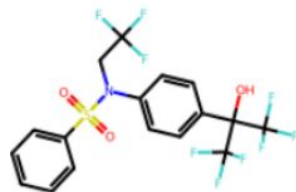
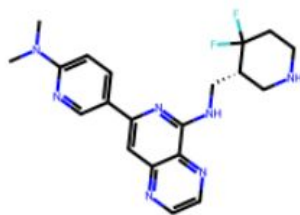
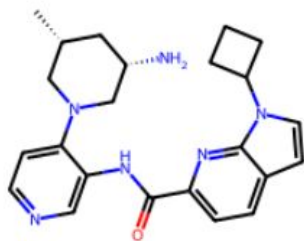
smiles	Decoded_SMILES	Min_Distance
<chem>C[C@@H]1C[C@H](N)CN(C1)c1ccncc1NC(=O)c1ccc2ccn...</chem>	<chem>CN(C)c1ccc(cn1)-c1cc2nccnc2c(NC[C@H]2CNCCC2(F)...</chem>	4.233210
<chem>CC(C)c1nnc2ccc(cn12)-c1ocnc1-c1ccc(F)c(Cl)c1</chem>	<chem>Cn1cc(cc(N)c1=O)-c1cccc(c1CO)-n1c(cc2cc(ccc2c1...</chem>	4.535190
<chem>CSc1nc(c([nH]1))-c1ccc(F)cc1)-c1ccnc(NC2CCCCC2C)c1</chem>	<chem>CC(C)c1cc(C(=O)N2Cc3ccc(CN4CCN(C)CC4)cc3C2)c(O...</chem>	4.616539
<chem>CCN(CC)CCN1C(=O)[C@](O)(c2c1cc(cc2C(F)(F)F)C(N...</chem>	<chem>COc1ccc(Cn2c(nnc2[C@H])(Cc2ccccc2)NC(C)=O)[C@@H...</chem>	4.581069
<chem>NS(=O)(=O)c1cc(C(=O)NC2CCCCC2)c(SCc2ccccc2)cc1Cl</chem>	<chem>NS(=O)(=O)c1ccc(cc1)N1N=C(CC1c1ccc2OCCOc2c1)C(...</chem>	3.765682
...	...	...
<chem>NC(=O)c1ccc2cc(ccc2c1)C1(O)CCn2cncc12</chem>	<chem>Cc1c(-c2ccnc3c(F)cccc23)c2cc(C)ccc2n1CC(O)=O</chem>	3.492570
<chem>OC(c1ccc(cc1)N(CC(F)(F)F)S(=O)(=O)c1ccccc1)(C(...</chem>	<chem>OC(=O)COc1ccc(Cl)cc1[C@H]1N(CCc2ncsc12)C(=O)[C...</chem>	3.526263
<chem>CCc1cc2c(cc1C(=C)c1ccc(cc1)C(O)=O)C(C)(C)CCC2(C)C</chem>	<chem>Cc1c(-c2ccnc3c(F)cccc23)c2cc(C)ccc2n1CC(O)=O</chem>	3.844205
<chem>O=C1NCc2cc(ccc12)S(=O)(=O)NCCNC\C=C\c1ccc(cc1)...</chem>	<chem>Cn1ccc(Nc2nc(N)cc(n2)-c2cccc(c2CO)-n2ccc3cc(cc...</chem>	4.524452
<chem>CN1CC[C@H](CC1=O)c1ccncc1Oc1ccc(Nc2nc3ccccc3s2...</chem>	<chem>COc1ccccc1-c1nccnc1C1CN(C1)C(=O)c1nc2ccccc2[nH]1</chem>	4.654112

# Tanimoto Similarity (Generated vs. Ground Truth)

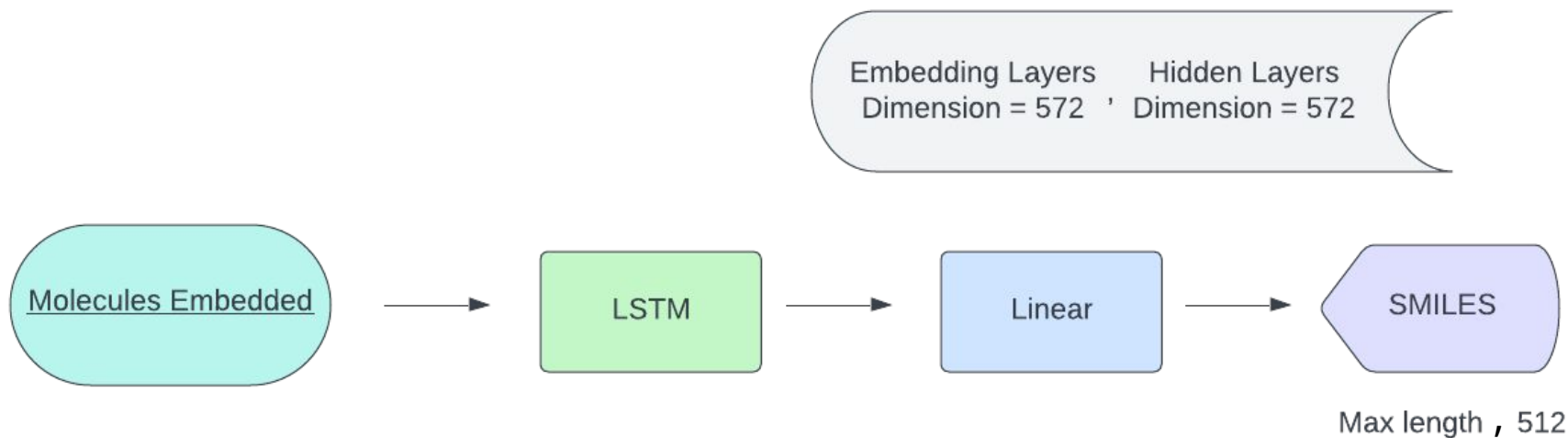




# Decoded Structure Output Comparisons (Look-up Table)



# LSTM Decoder Architecture



# Limitation of Architecture

- **Limitations:**

- Static Context: Lacks temporal dynamics, providing the same context at each timestep.
- Over Reliance on Initial Embeddings: Limits sequence diversity and adaptability.
- Risk of Redundancy: Repetitive, unvaried outputs in longer sequences.

- **Alternatives:**

- Positional Encodings: To add temporal information.
- Direct Sequence Generation: Fastest using seq2seq internal model.
- Autoregressive Generation: One sequence at a time for dynamic context.
- Learned Transformations: Simplest method from vector to sequence.
- Recommended for Our Application: Combine learned transformation with positional encoding for optimal results.

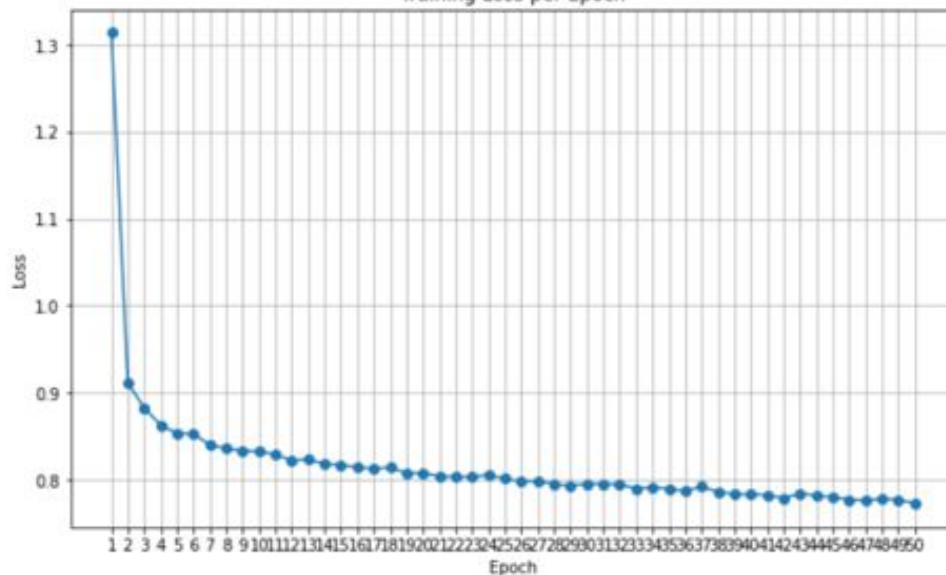
# Decoded SMILES (LSTM-Decoder)

Decoded\_SMILES

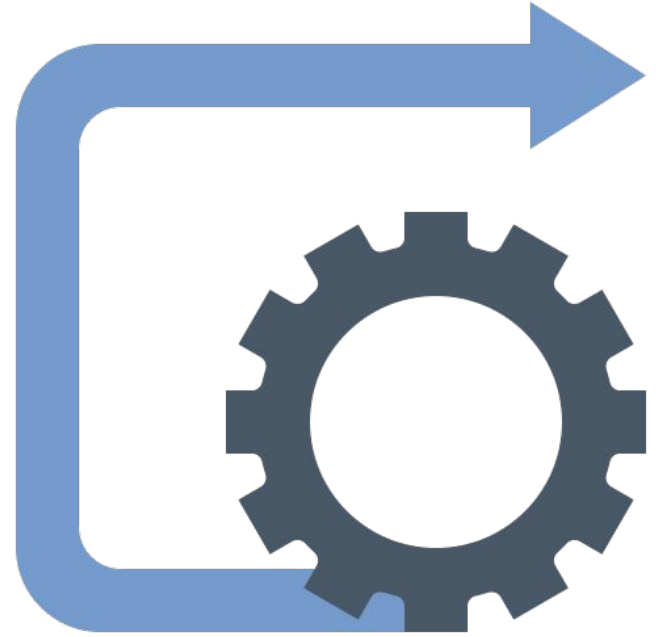
0	CCCCCCCCcccccccccccccccccccccccccccccccc...
1	Cc1cccccccccccccccccccccccccccccccccccc<P...
2	CC(Ccccccccccccccccccccccccccccccccccccc...
3	CC(Ccc(ccccccccccccccccccccccccccccccccc...
4	CC(==(ccccccccccccccccccccccccccccccccccc<...
...	...
745	CC(cccccccccccccccccccccccccccccccccc<PAD><PAD><...
746	FC(Fcccccccccccccccccc((((((((((((((((0))))))...
747	CC(=cccccccccccccccccccccccccccccccccc))))))<P...
748	Oc1cccccccccccccccccccccccccccccccccccc...
749	Cc1ccc(Ccccccccccccccccccccccccccccccccc...

750 rows × 1 columns

Training Loss per Epoch



# Limitations & Future Work



# Limitations

- **Featurized Embeddings are Averaged**
  - The averaging of the embeddings takes away the sequential nature of the proteins and molecules.
  - Suggestion: Experiment with tokenized embeddings, not average embeddings.
- **Transition Model Complexity:**
  - Current linear architecture might be too simplistic for modeling noise addition and removal in molecule generation.
  - Suggestion: Incorporate convolutional layers or attention mechanisms for enhanced spatial and sequential understanding.

# Limitations

- **Variance Schedule Rigidity:**

- Predetermined variance schedule may not fit all protein-molecule interactions, potentially limiting diffusion effectiveness.
- Suggestion: Implement an adaptive variance schedule for more effective diffusion processes

- **Output Interpretability:**

- Generated molecular embeddings are abstract and challenging to interpret in chemical structure terms.
- Suggestion: Add a decoding mechanism to convert embeddings into interpretable structures like SMILES.

- **Epsilon Network Design:**

- Linear layers may not capture complex, non-linear relationships between protein and molecule embeddings.
- Suggestion: Use transformer layers or gated recurrent units for better pattern recognition.

# Future Avenues

- **Enhancing Diversity and Generalization:**
  - Explore adversarial training, domain adaptation, and few-shot learning to improve generalization and molecular diversity.
- **Efficiency Improvements:**
  - Optimize model architecture and use efficient algorithms to reduce computational demands.
- **Integration with Other Biological Data:**
  - Integrate additional biological data like gene expression profiles to generate more targeted molecules and specific diseases.



The background features a light gray gradient with teal-colored wavy shapes at the top and bottom. Several small teal dots are scattered across the page, including three near the top and four near the bottom.

# Questions?