

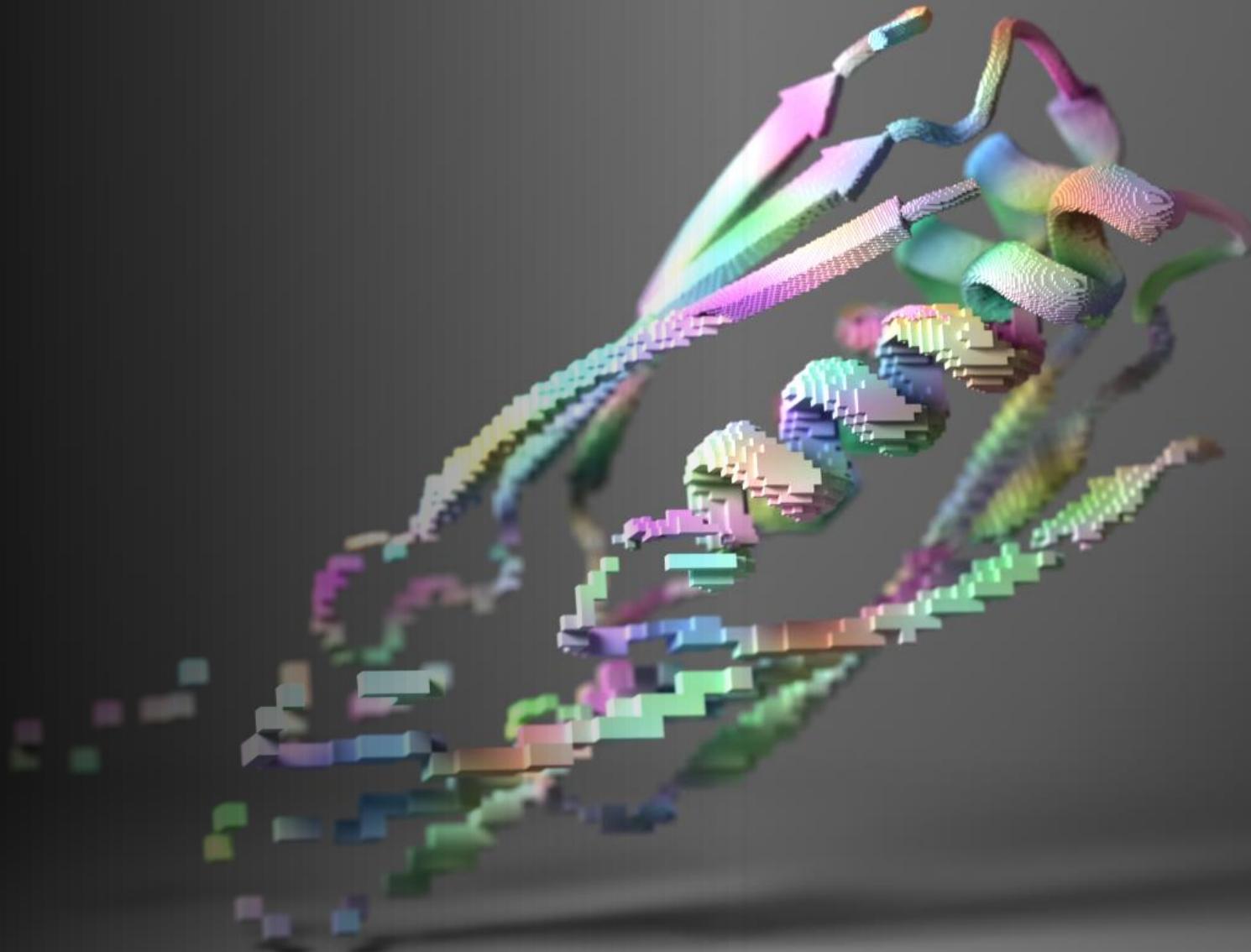
# Generating Protein Conformation Ensemble Through A Transformer- based Diffusion Model

YI YANG

November 13, 2024



SMU®



# OUTLINE

## Part I

Protein structure  
and function,  
protein dynamics



## Part II

How to build a  
transformer-  
based diffusion  
model

I O I O  
I O I O

## Part III

Computational  
result for applying  
our model to vivid  
protein system



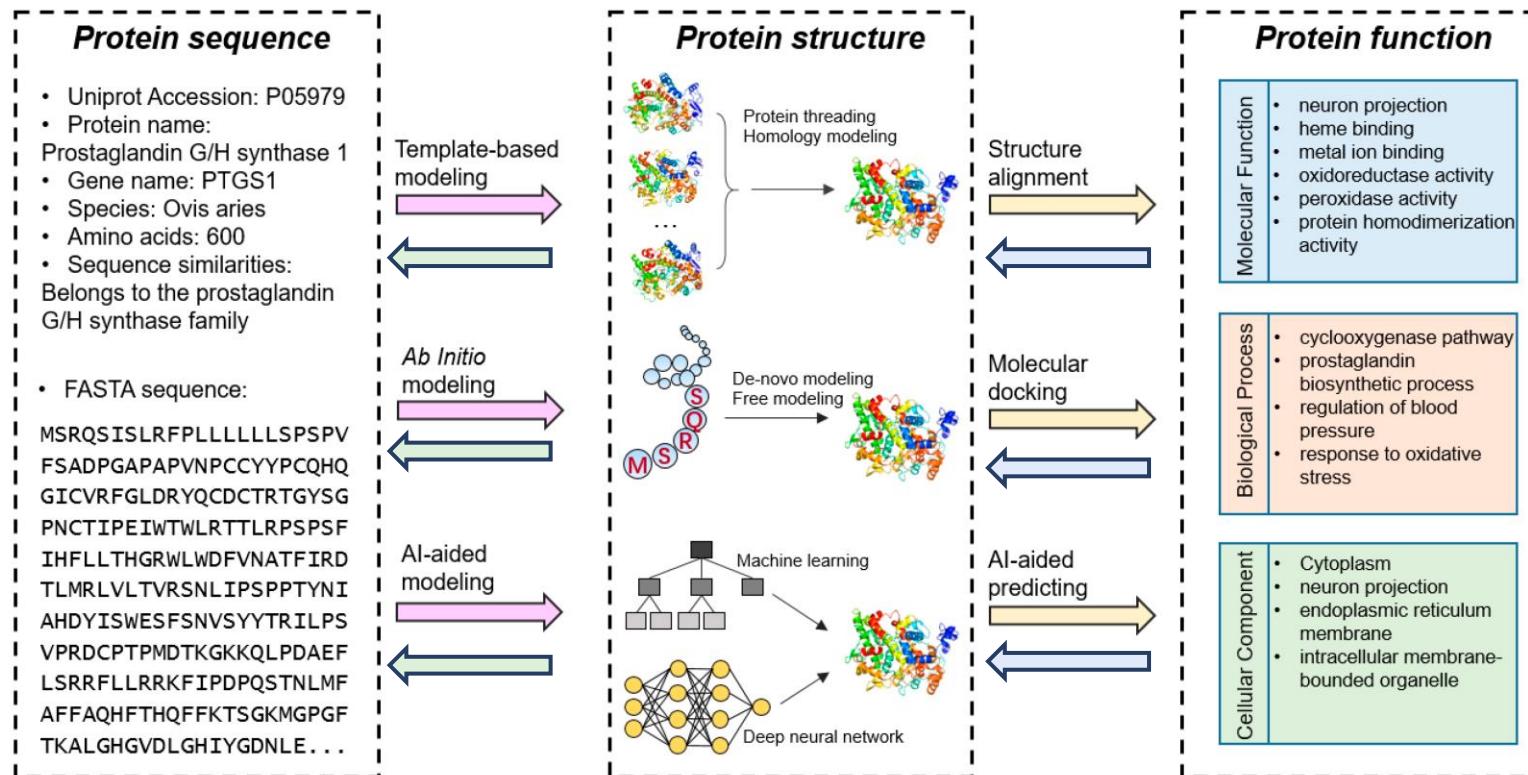
## Part IV

Conclusion



# The relationship between Protein structure and function

Generally, protein sequence determines structure and structure determines function.

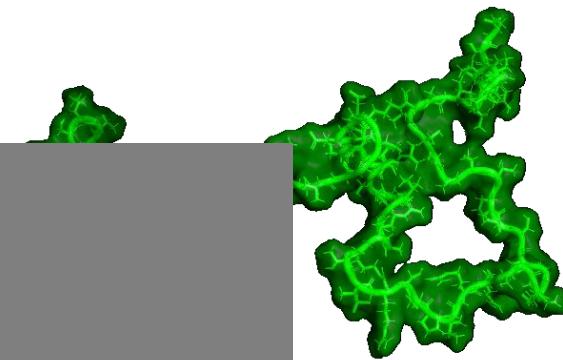
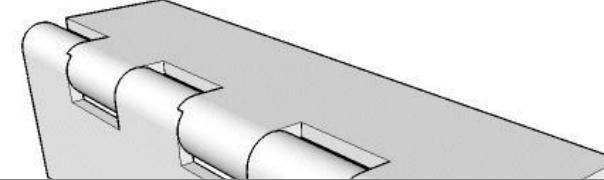
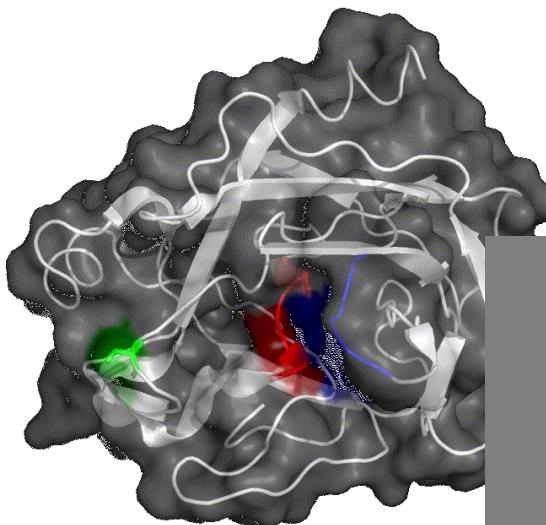


## Protein in Dynamics



## Protein Function

But some proteins exhibit functional flexibility through dynamic structural changes



## Obtaining protein dynamic ensemble is important !!!

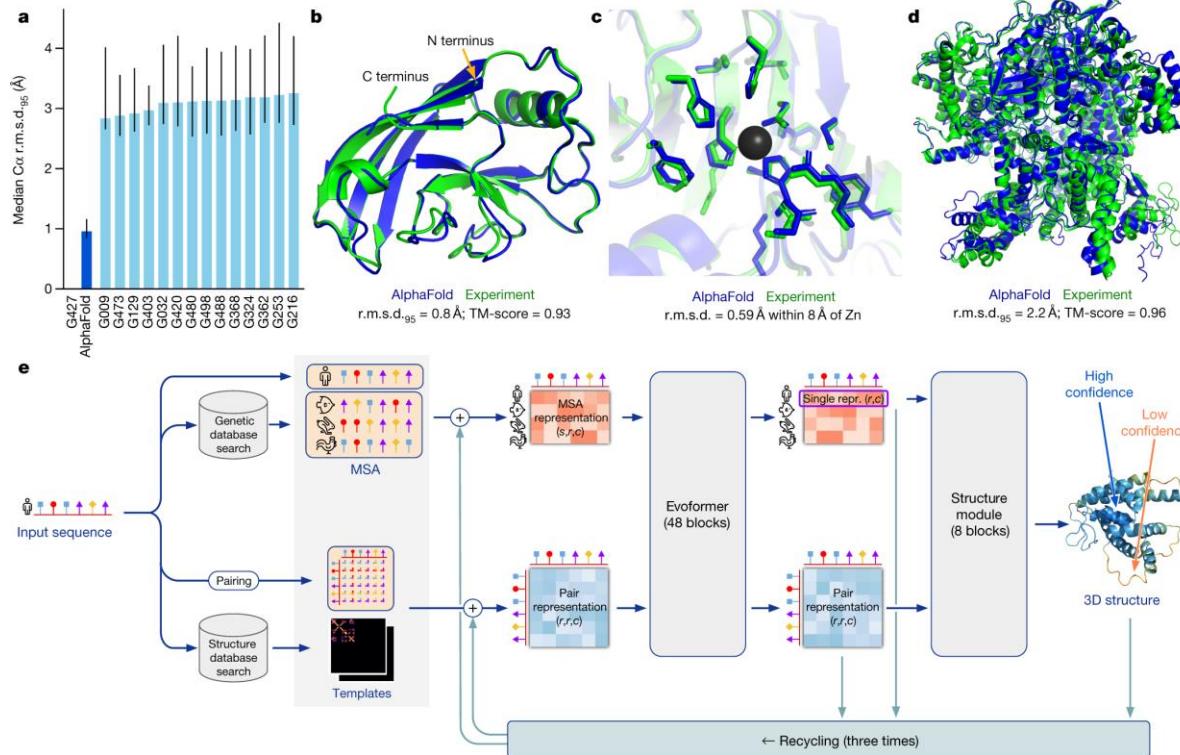
- **Hinge motion** in disordered activation domain in Trypsinogen (PDB ID: 2PTN)
- Biological process: digestion

- An intrinsically disordered protein (IDP)
- A neuron-specific microtubule-associated protein that plays an important role in maintaining neuronal structure and function(MAP2C)

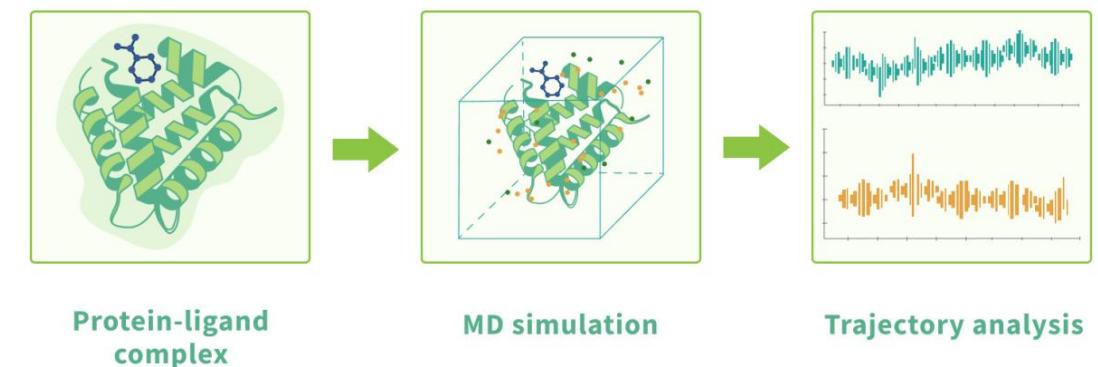
## AlphaFold 2

AlphaFold2 can only give us static structures

## AlphaFold 2



## MD simulation



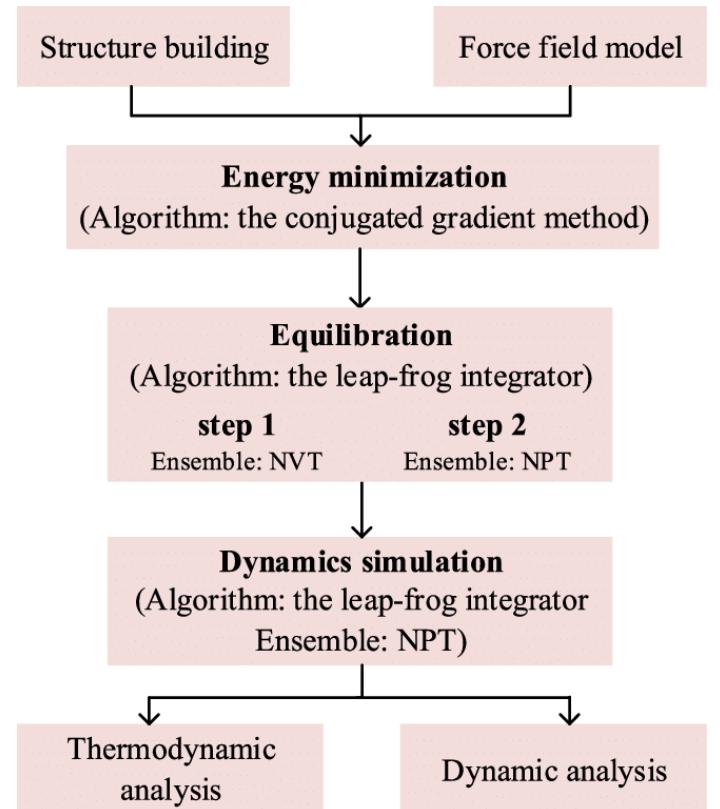
# MD simulation

## MD simulation:

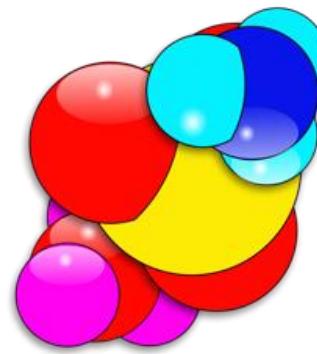
- Newton's laws of motion
- Repeatedly update the forces on each atom by all other atoms
- Capture position and motion of every atom at every point in time

$$F = ma = -\nabla U(x)$$

$$\frac{dx}{dt} = v \quad \frac{dv}{dt} = a$$



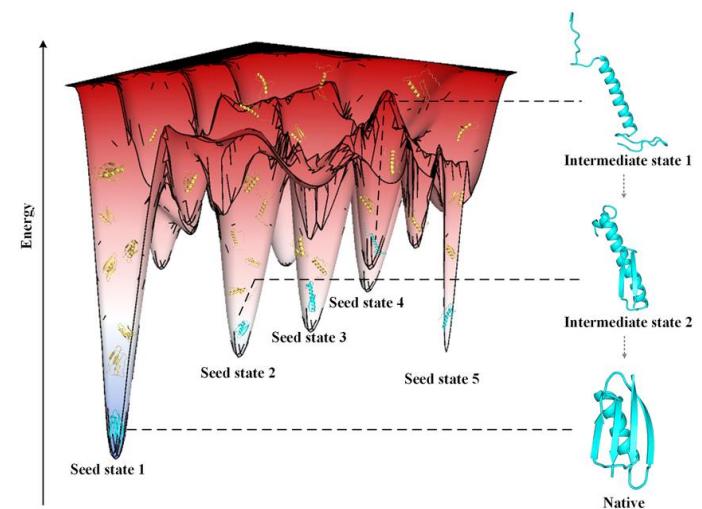
## Computational expensive



## Folding@home

Folding@home is a distributed computing project for simulating protein dynamics.

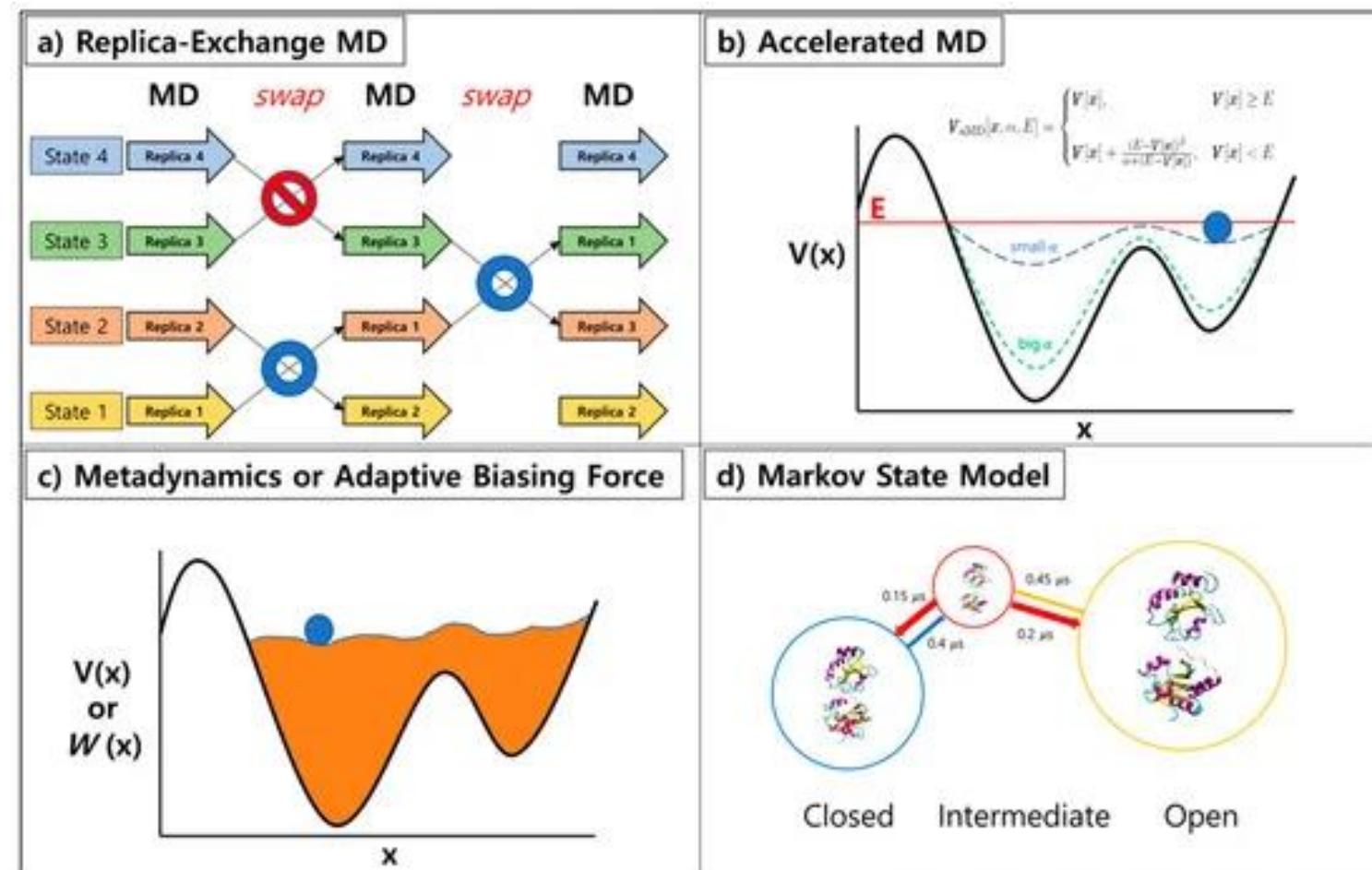
Low efficiency for sampling all possible conformation

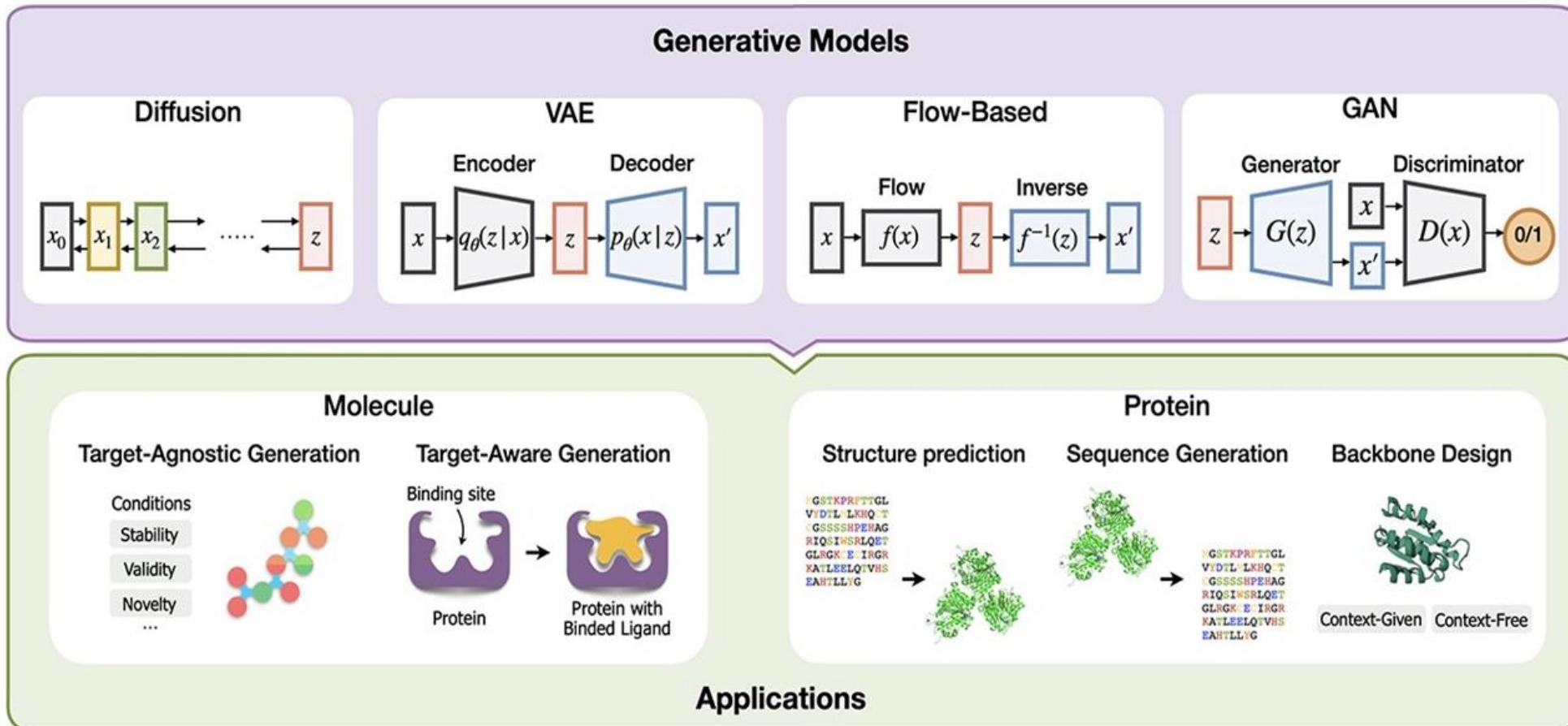


## Enhance MD simulation

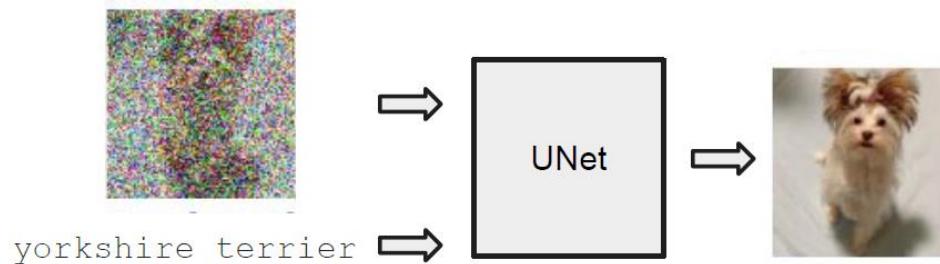
## Overview of Enhanced Sampling Methods

- (a) replica-exchange MD simulation
- (b) accelerated MD
- (c) metadynamics or adaptive biasing force
- (d) Markov state model.

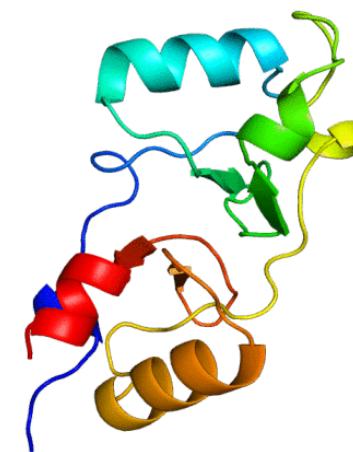




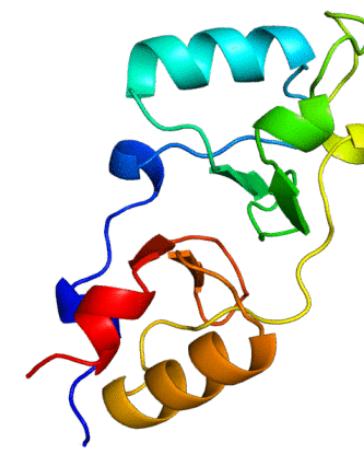
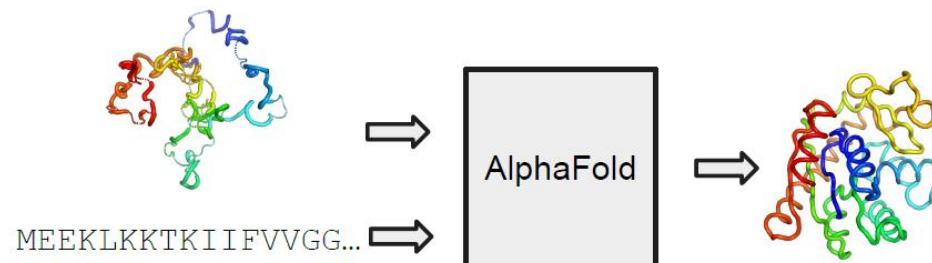
## AlphaFlow

**Text-to-image generative model**

300 ns MD

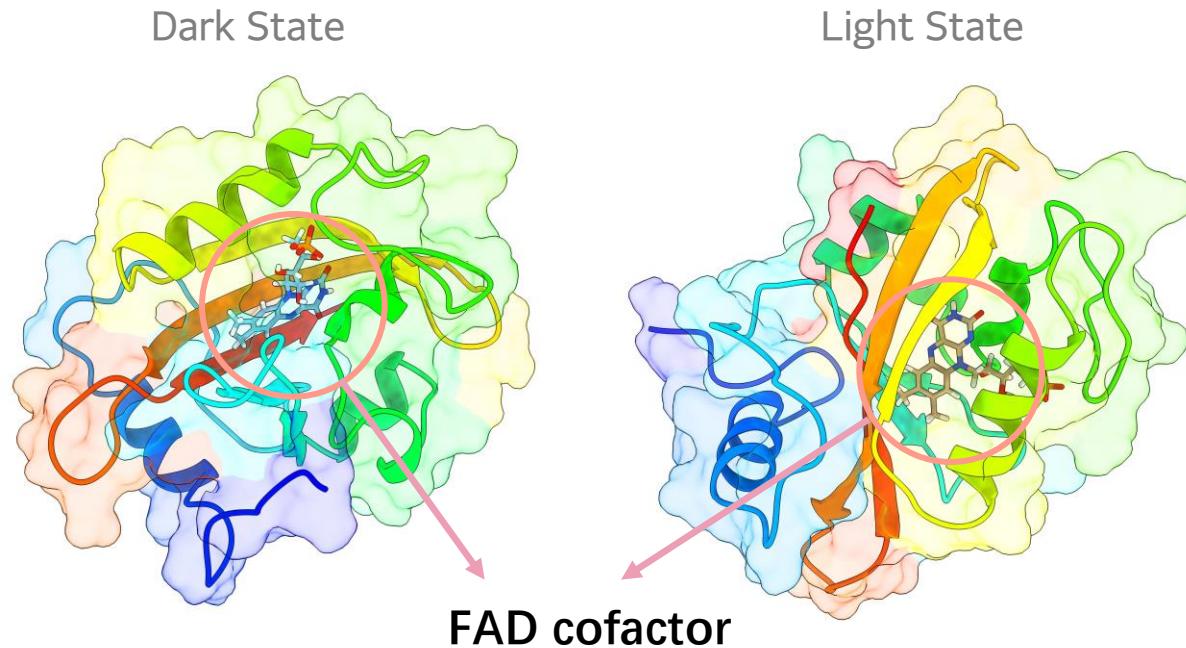


AlphaFlow

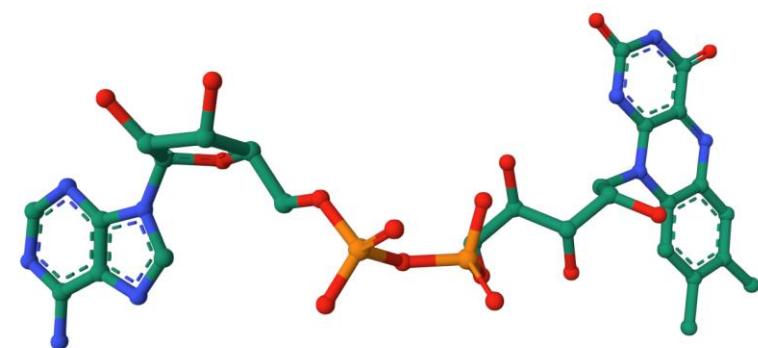
**Sequence-to-structure generative model**

The presence of small molecules or cofactors that significantly influence the protein's conformational ensemble

### The vivid Photoreceptor protein system

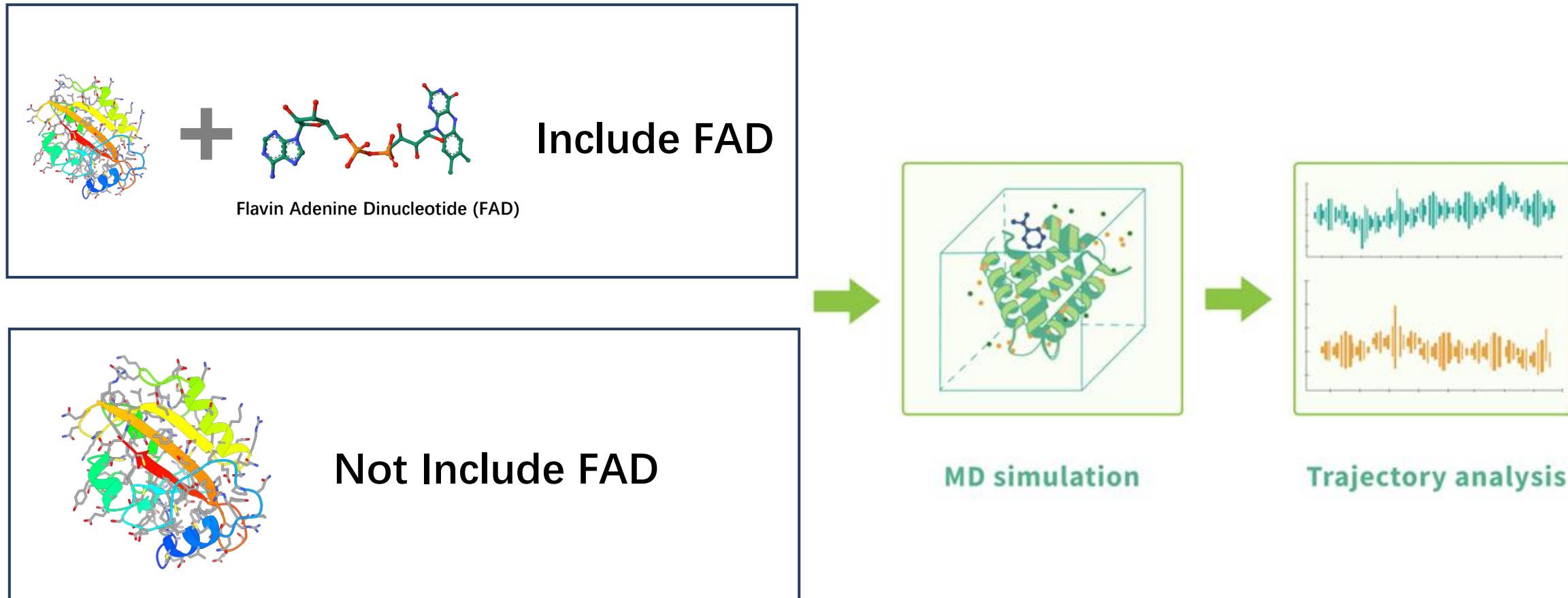


### Flavin Adenine Dinucleotide (FAD)



## Complex Protein Systems

The presence of small molecules or cofactors that significantly influence the protein's conformational ensemble



# Complex Protein Systems

**So only use the sequence as input, It is challenging to produce accurate ensembles.**

The presence of small molecules or cofactors that significantly influence the protein's conformational ensemble



## Use AlphaFlow

Official Neurosnap webserver for accessing AlphaFlow online.

[Protein Conformations](#) [Molecular Dynamics](#)

### Overview

Use AlphaFlow's large variety of models to generate protein structures that closely reflect experimental and physiological conditions using ESMFlow and AlphaFlow models trained on experimental structures and molecular dynamics trajectories.

### Neurosnap Overview

The AlphaFlow online webserver allows anybody with a Neurosnap account to run and access AlphaFlow, no downloads required. Information submitted through this webserver is kept confidential and never sold to third parties as detailed by our strong terms of service and privacy policy.

[View Paper](#)

### Features

- Generate many protein conformations resembling experimental and physiological ensembles.
- Supports sequence input for generating conformations.
- Provides four model choices to fit experimental design parameters.

### Configuration & Options

#### Model Inputs

Input Sequence [Select Proteins](#) 0 / 1 proteins added.

The primary sequence of the protein you wish to predict conformations for.

Model Weights AlphaFlow - Molecular Dynamics

The following options are the available model weights available for AlphaFlow.

Number of Conformations 100

The number of conformations to generate for the input protein. The more conformations the longer the job runtime.

# SUMMARY

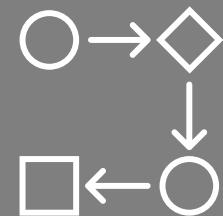
- Protein function is tied to its structure, with both static and dynamic conformations playing critical roles in defining biological activity.

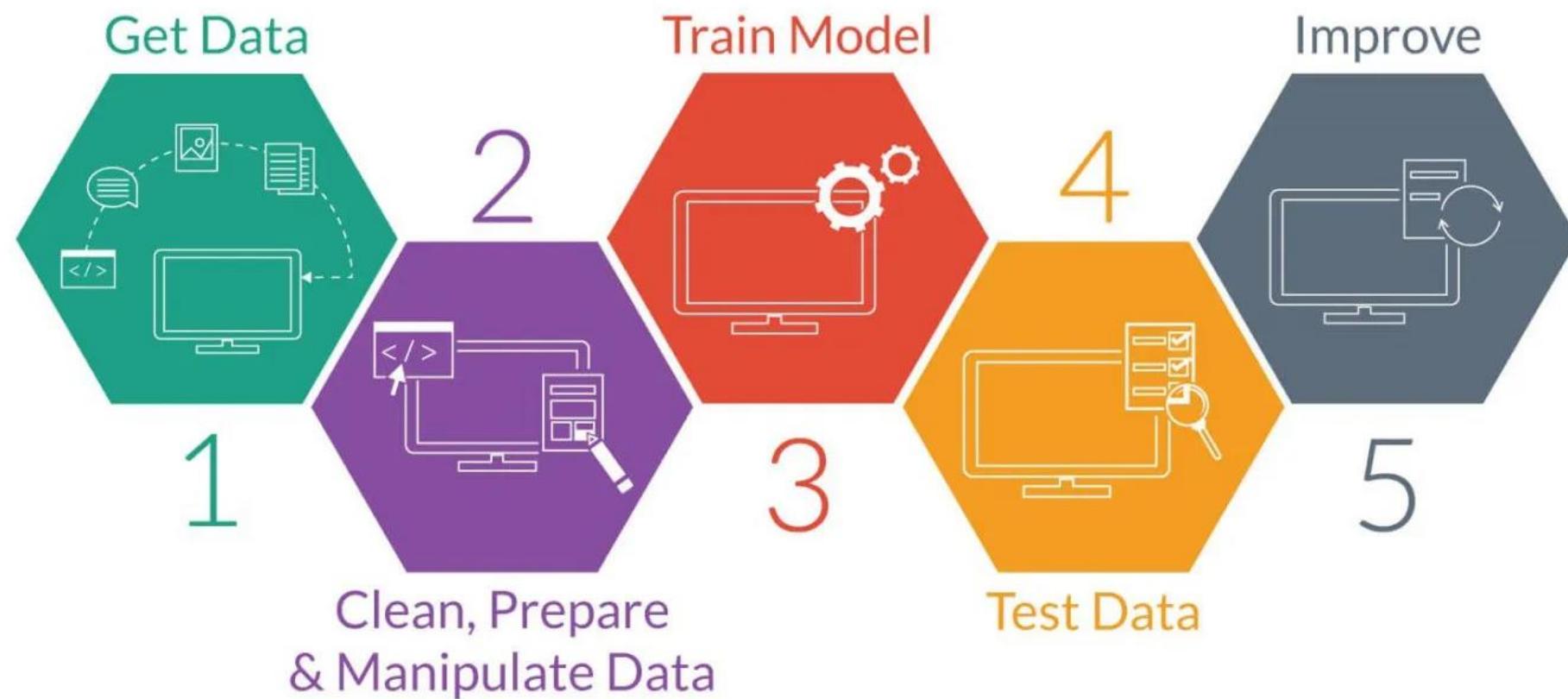
- MD simulation is the primary tool for exploring the dynamic properties of protein structures, the major limitation of MD simulations is their requirement for extensive computational resources.

- Generative AI has opened new avenues for protein conformation generation, but the current models struggle to accurately generate conformations for complex protein systems.

## Part II

A transformer-based diffusion model to generate protein conformations

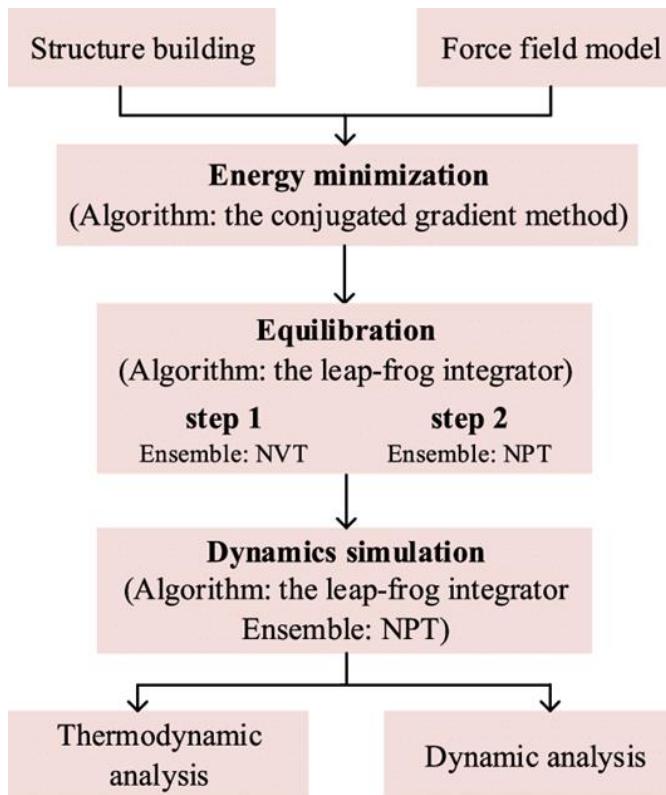




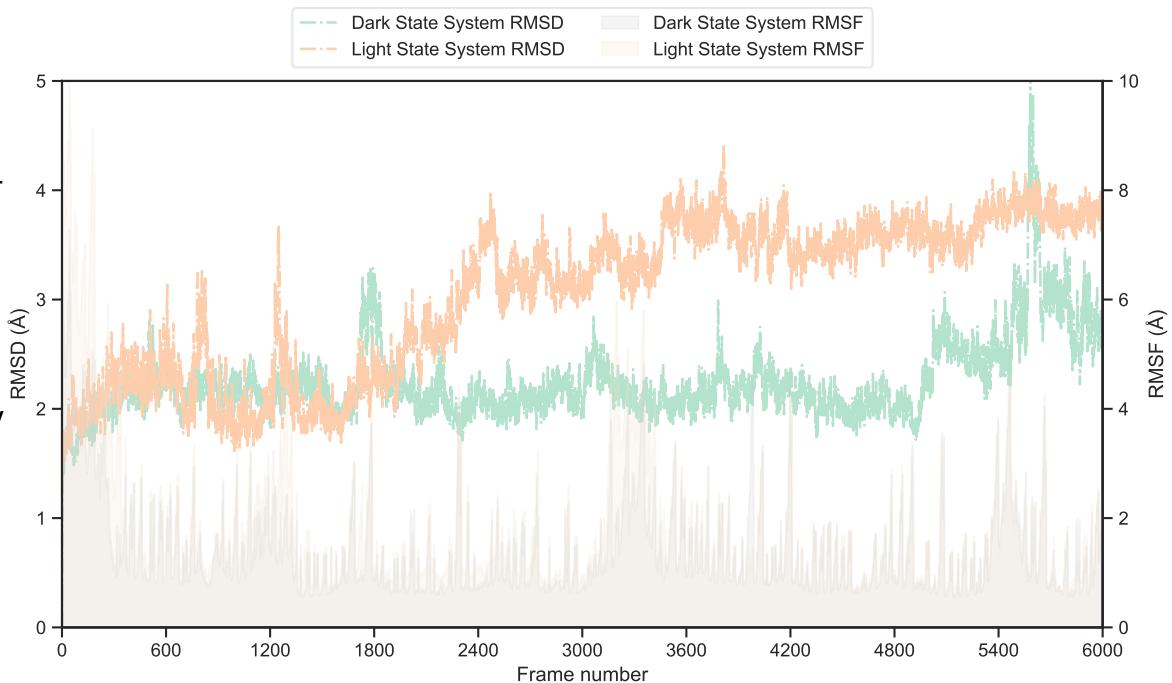
# 01 GET DATA

We collected a total of **6,001** distinct conformations of the VVD protein of each system.

A total of two simulation systems were constructed

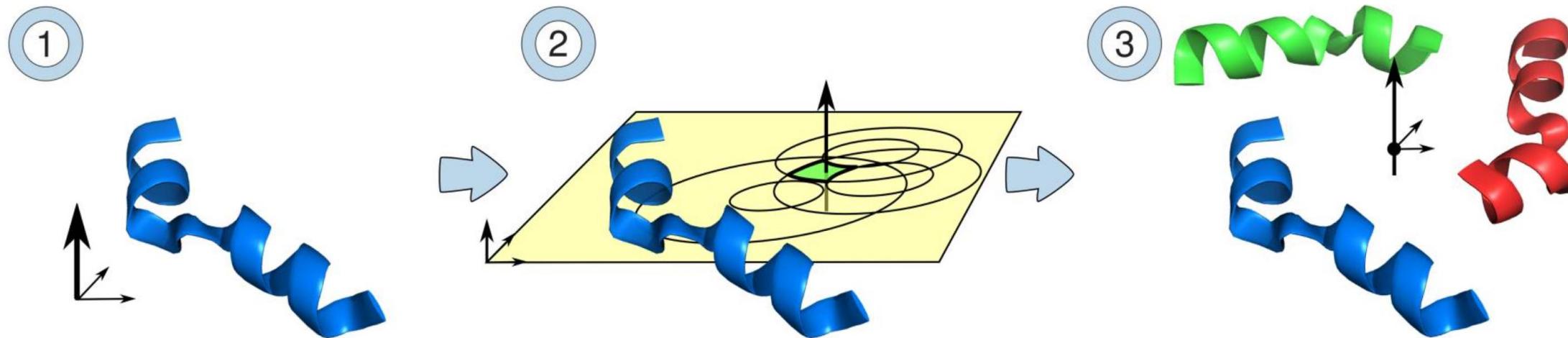


- | 10 nanoseconds (ns) of NPT MD simulations for Equilibrium
- | 1.1 microseconds ( $\mu$ s) of NVT MD simulations for production



The input feature and output feature in our model

## SE(3) symmetry of protein structures

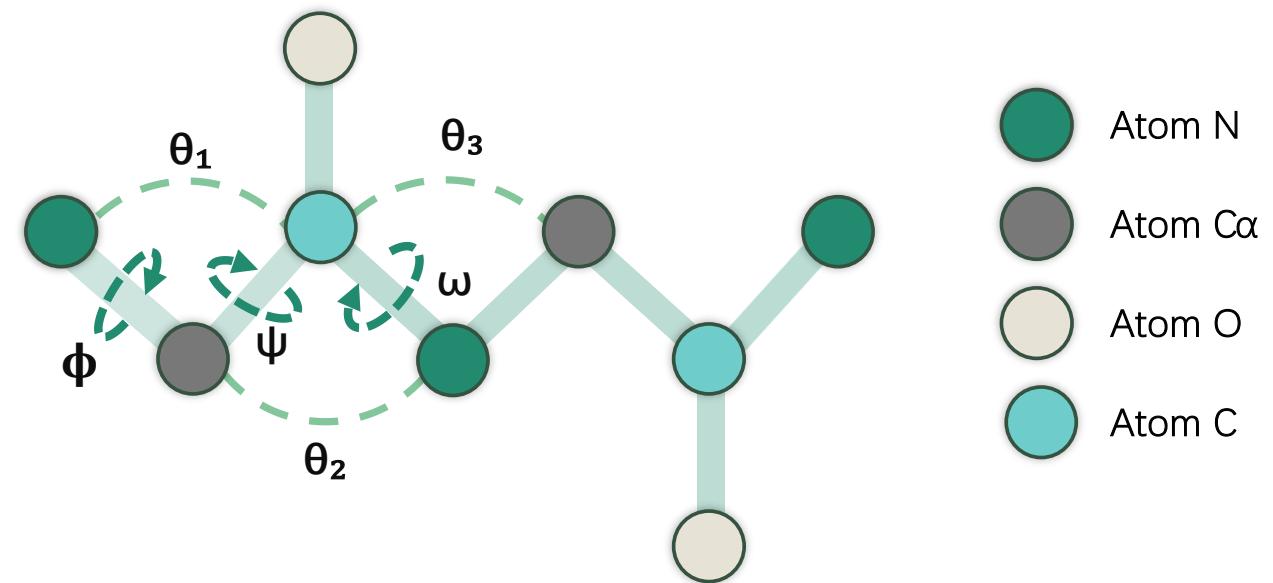


## 02 Clean, Prepare, and Manipulate data

The input feature and output feature in our model

### Six types of angles from the protein backbone structures during the MD simulation

- dihedral torsions (' $\phi$ ', ' $\psi$ ', ' $\omega$ ')
- bond angles (' $\theta_1$ ', ' $\theta_2$ ', ' $\theta_3$ ')
- Data shape (N,147,6)



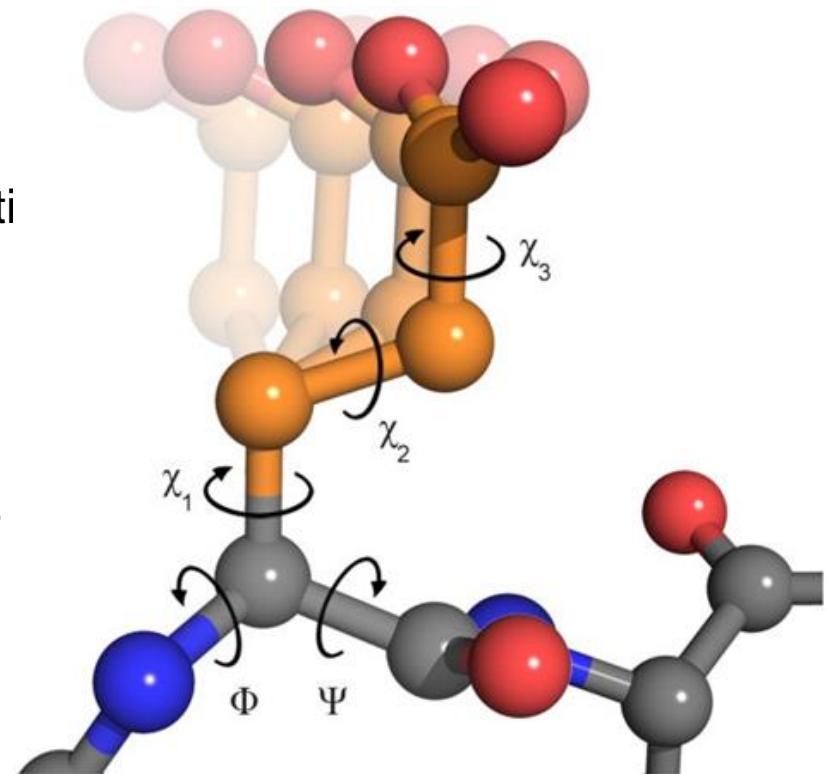
One of the Input and output features in our first model

## 02 Clean, Prepare, and Manipulate data

The input feature and output feature in our model

- The dynamic nature of the structure during MD simulation.
- Angular Deviations(changes) relative to the structure from the first time step of the MD conformation set.

**The other of the Input and output features in our second model**



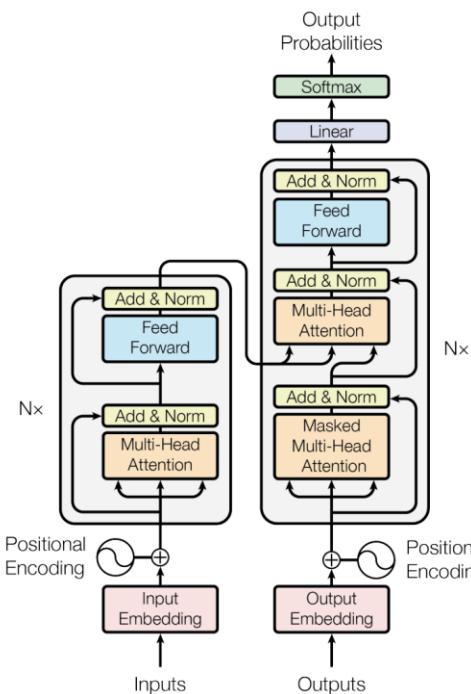
Angular Deviations

## 03 Select Model

A transformer-based diffusion model

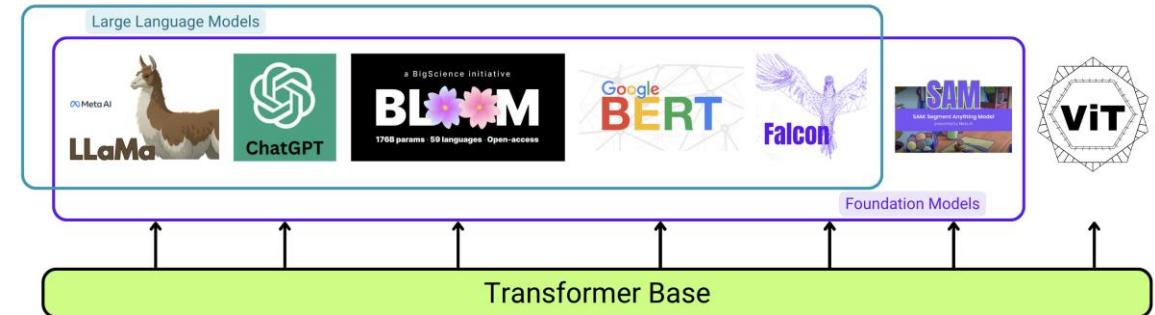
**BERT**

Encoder



**GPT**

Decoder

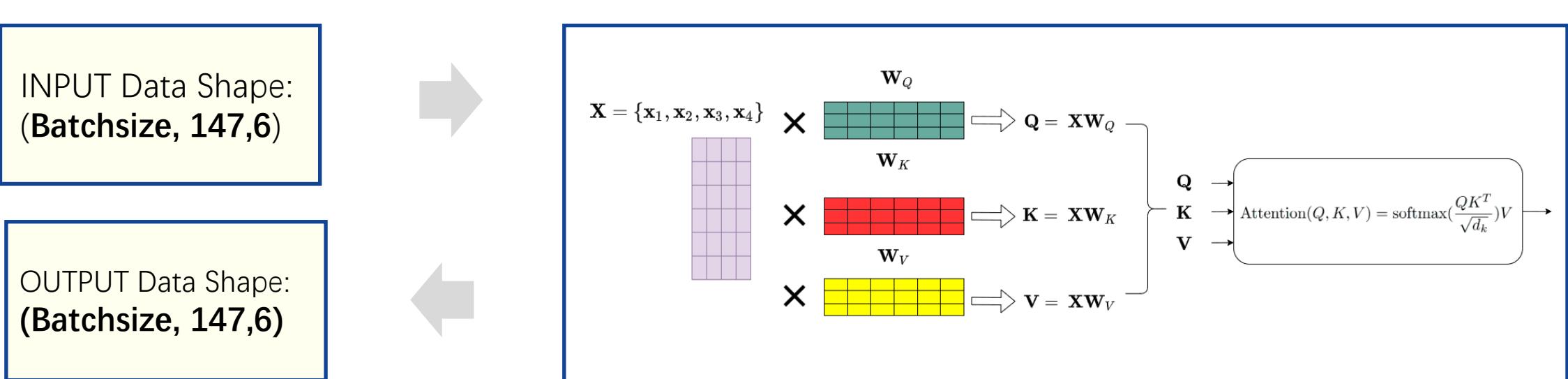


Hyperparameter	Value
Hidden Size	384
Number of Layers	12
Number of Attention Heads	12
Intermediate Size	768
Dropout Rate	0.1
Max Position Embeddings	147
Position Embedding Type	Relative Key
Temporal Embedding Type	Random Fourier Features

Only use the Encoder part like BERT

A transformer-based diffusion model

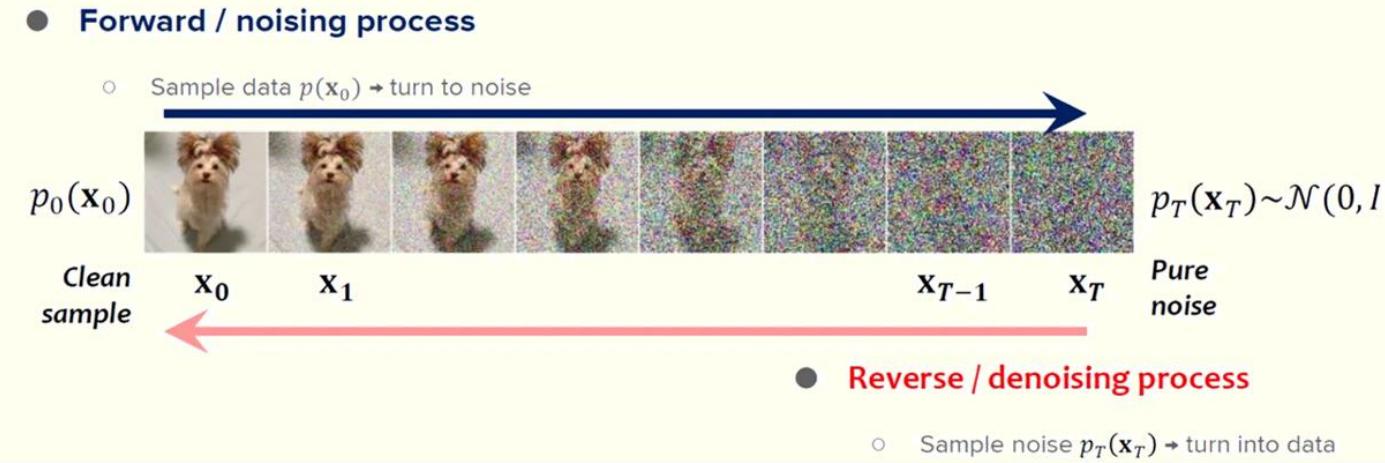
## Attention is all you need



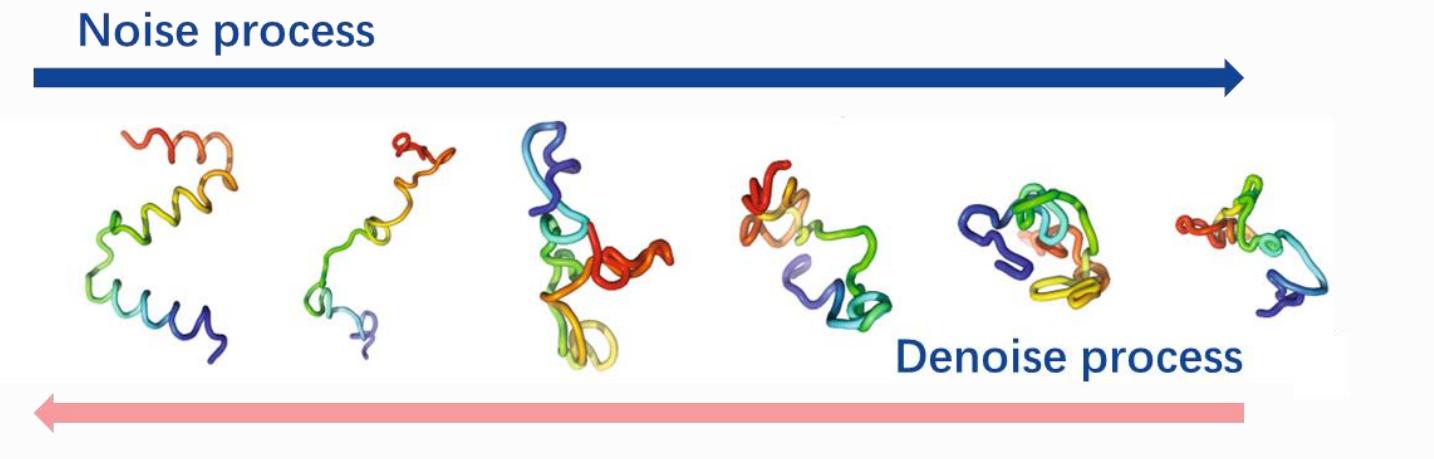
## 03 Diffusion Model

A transformer-based diffusion model

## Computer Vision



## Protein Structure



# 03 Model Algorithm

A transformer-based diffusion model

---

## Training Algorithm

```

1:  $x = x_{angle} - x_{ref}$  //calculate the angle deviations by subtracting the reference angle
2: do
3:  $x_0 \sim q(x)$  // sample the data from training data Set q(x)
4:  $t \sim U(\{1, \dots, T\})$  //sample the time t from a uniform distribution from 1 to T
5:  $\varepsilon \sim \mathcal{N}_{wrapped}(0, \mathbf{I})$  // sample the noise value from the wrapped gaussian distribution
6: gradient decent  $\nabla_{\theta}(L_w(\varepsilon, t, x_0))$  //gradient decent on loss function
7: until reaching the converge criterion

```

---



---

## Sampling Algorithm

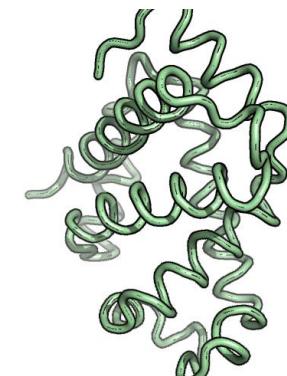
```

1:  $x_t \sim \mathcal{N}_{wrapped}(0, \mathbf{I})$  //sampling from the wrapped gaussian distribution
2: for  $t=T, \dots, 1$  do
3: if not last step:  $z \sim \mathcal{N}(0, \mathbf{I})$  else  $z = 0$ 
4: calculated  $x_{t-1}$  by equation  $x_{t-1} = \omega \left( \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \varepsilon_{\theta}(x_t, t)) \right) + \zeta z \sigma_t$ 
5: end for
6: generated angle =  $\omega(x_0 + x_{ref})$ 
7: reconstruct the Cartesian coordinate

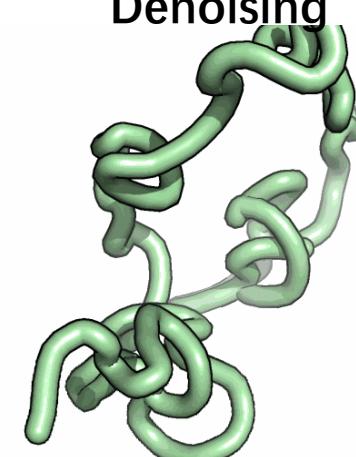
```

---

## Noising

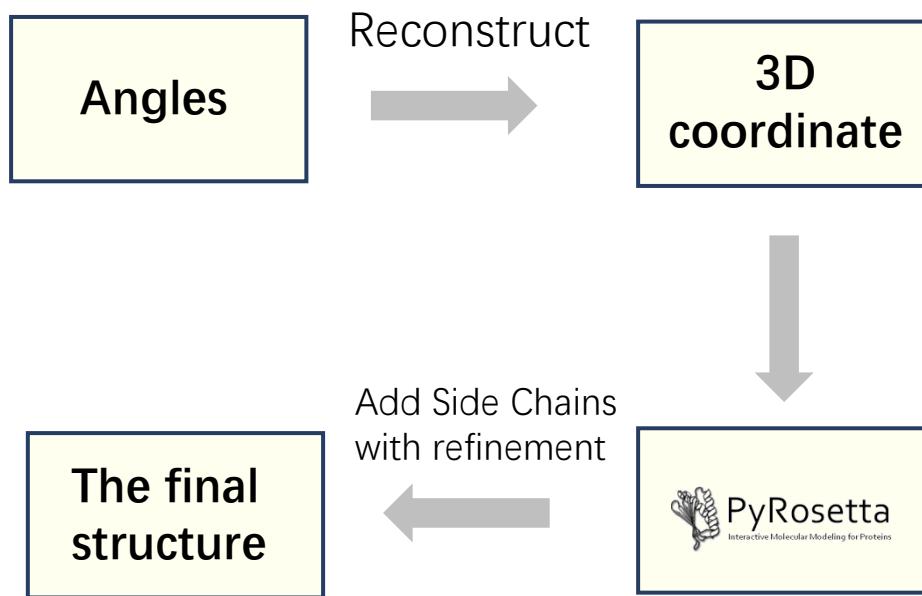


## Denoising

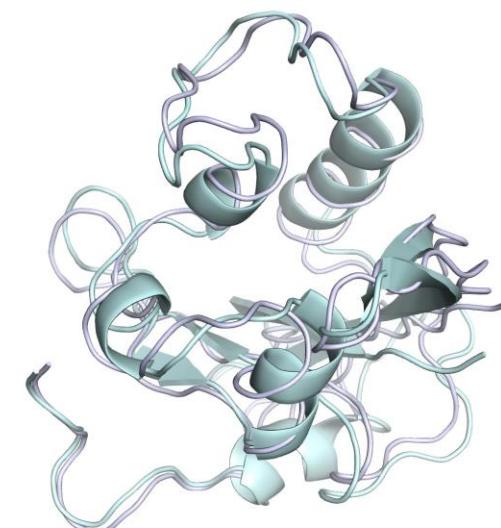


## 03 Add Side Chain

A transformer-based diffusion model

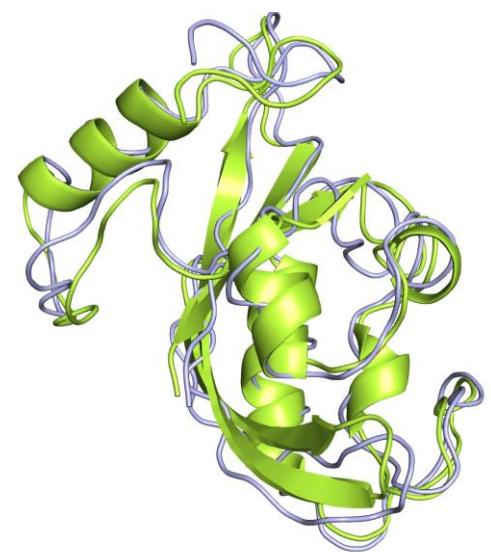


Dark State



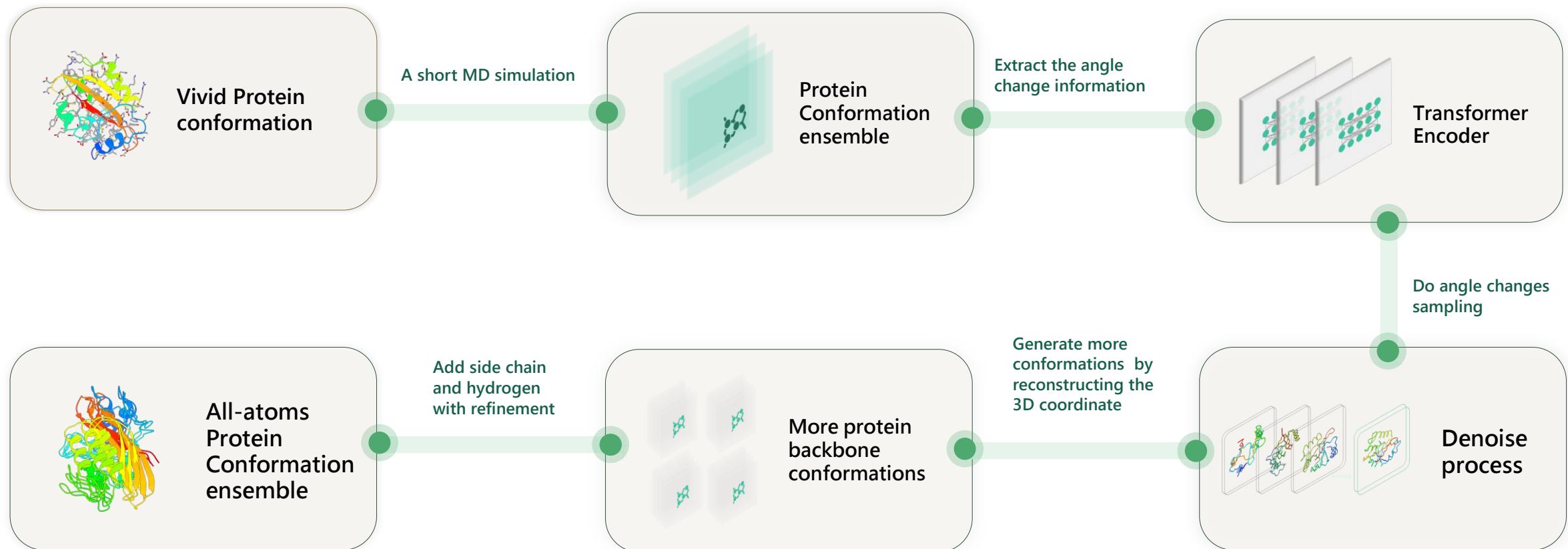
Backbone RMSD=1.9 Å

Light State



Backbone RMSD=1.8 Å

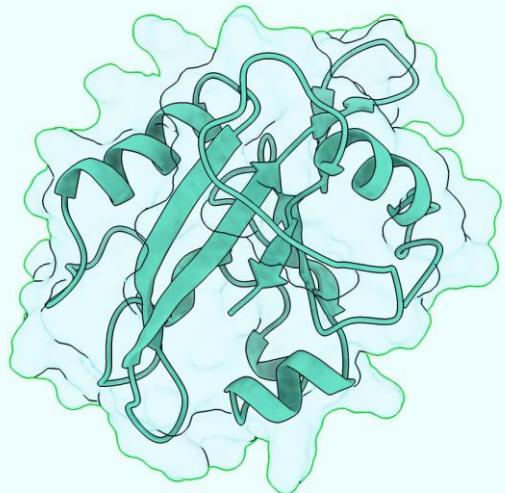
# CONCLUSION



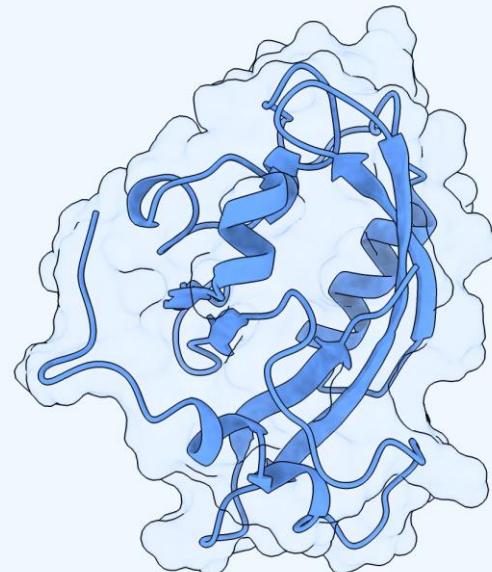
# Computational Results

- Angle vs angle deviation
- Compare with MD simulation
- Compare with AlphaFlow
- Limitation

Light State

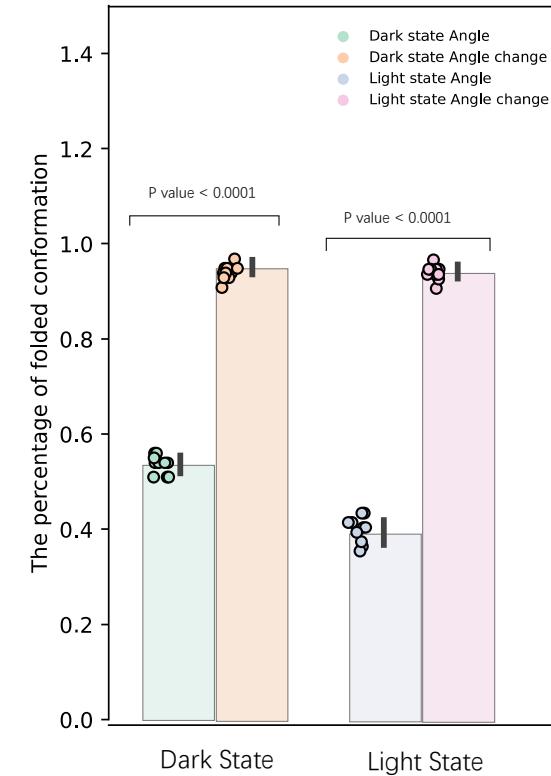
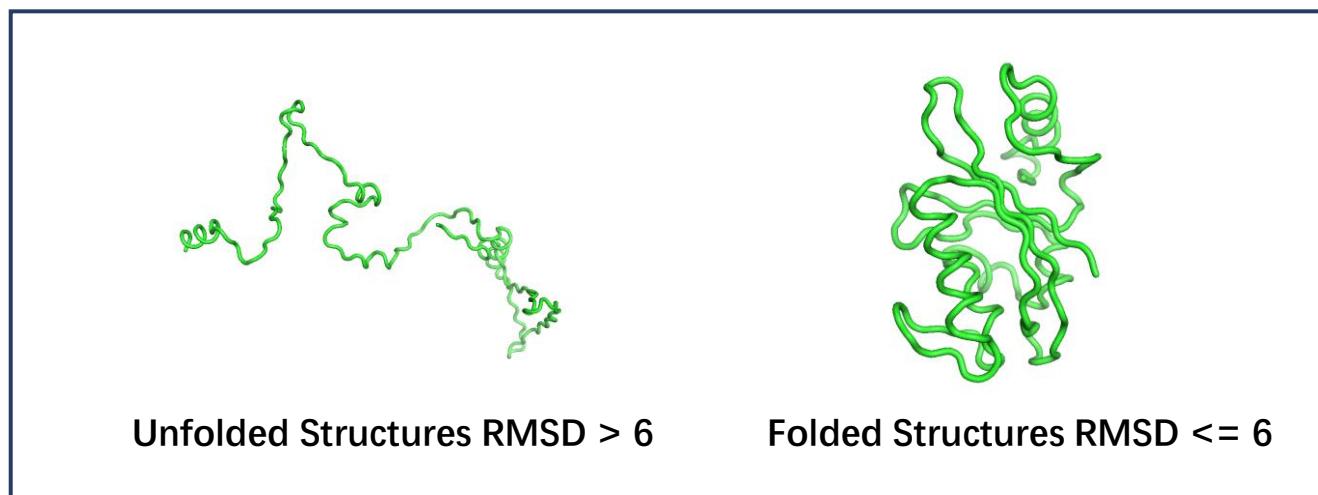
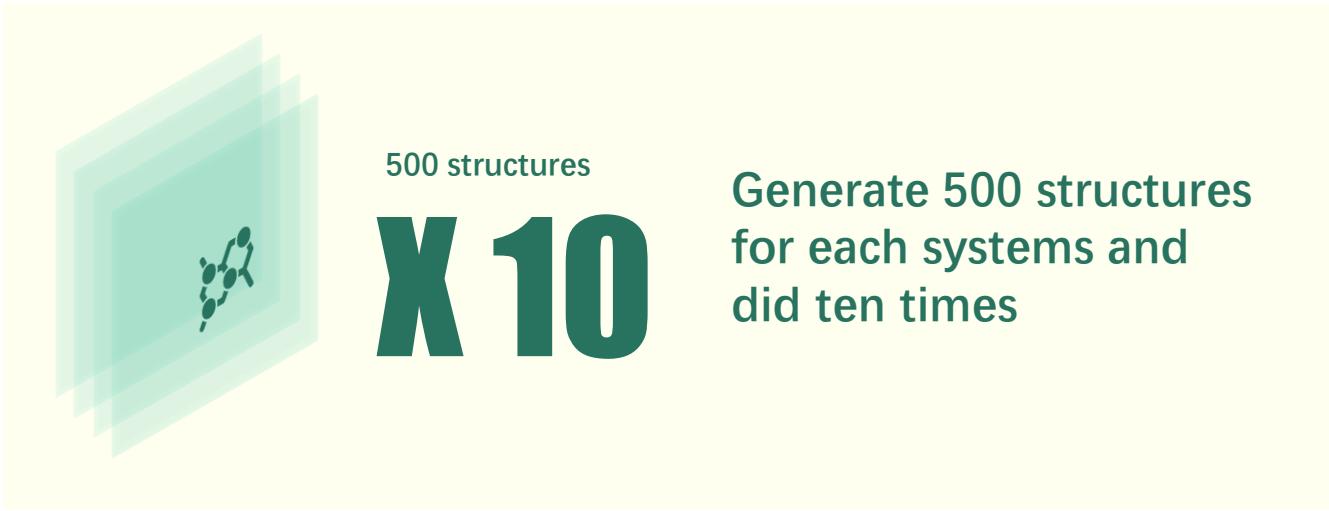


Dark State



# Angles vs Angle Deviation

Using Angle Deviation is Superior to Using Angles

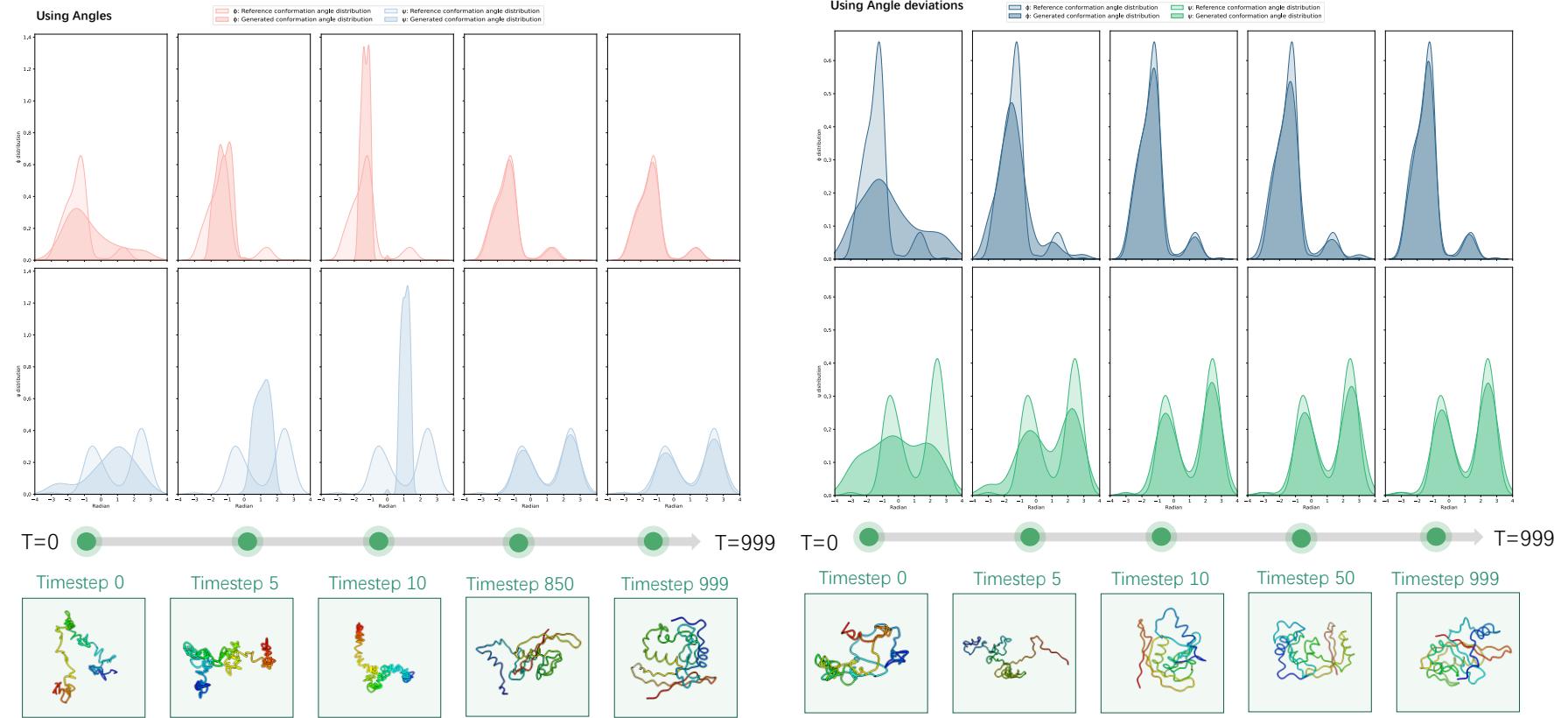


Comparison of the proportion of correctly folded structures generated using absolute angles versus angular deviation

# Angles vs Angle Deviation

Using Angle Deviation is Superior to Using Angles

Denoising process

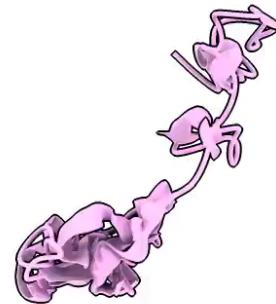


Denoising process of using angles and angle deviations for the Dark state system. Shifts in angle distributions, and conformations. (a) using angles (b) using angle deviations

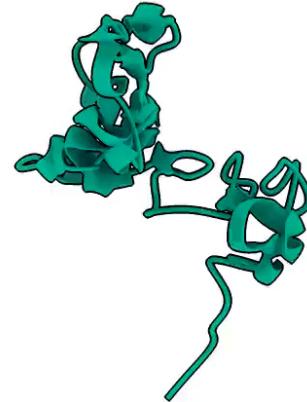
# Angles vs Angle Deviation

Using Angle Deviation is Superior to Using Angles

Using Angles



Using Angle Deviations



DENOISE PROCESS OF USING  
ANGLES AND ANGLE DEVIATIONS

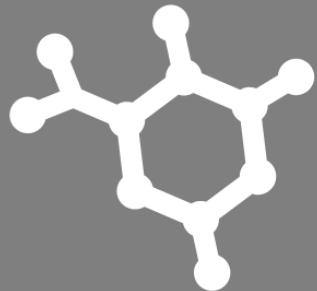
time step: 0

## Evaluations of conformations

### Using Angle Deviation

#### Our model

Generate 6000 conformations  
for each systems



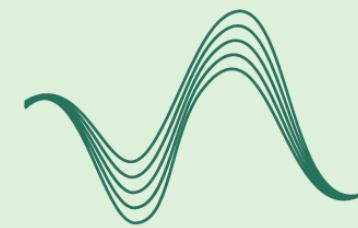
#### MD simulation

Collect 6001 conformations  
for each systems



#### Alpha Flow

Generate 6000 conformations  
for each systems only use  
Sequence



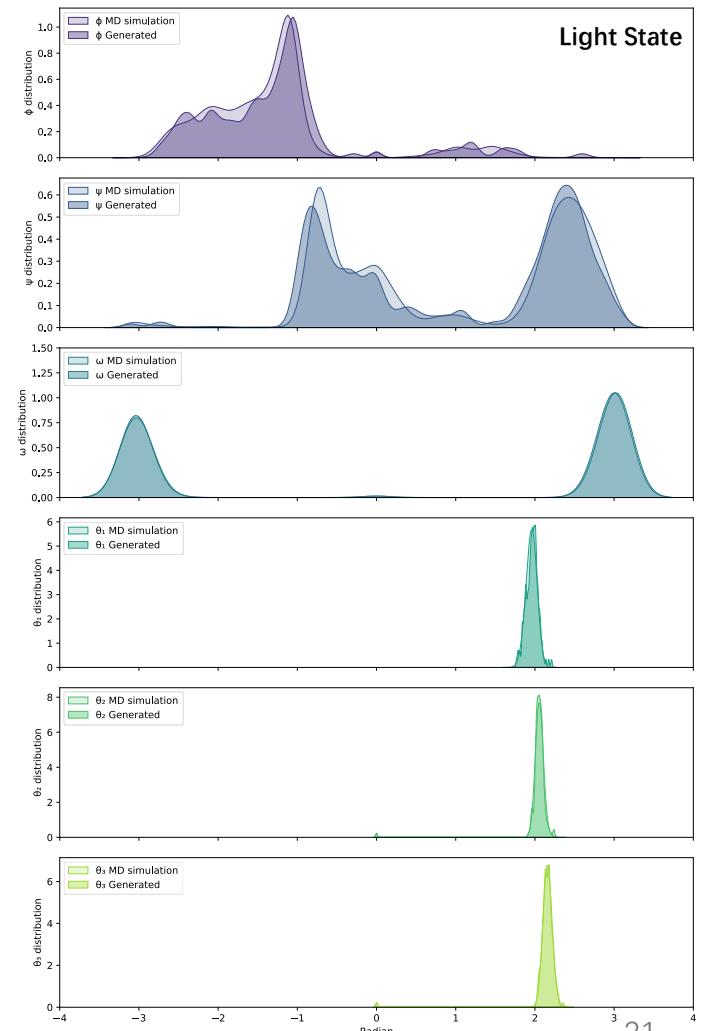
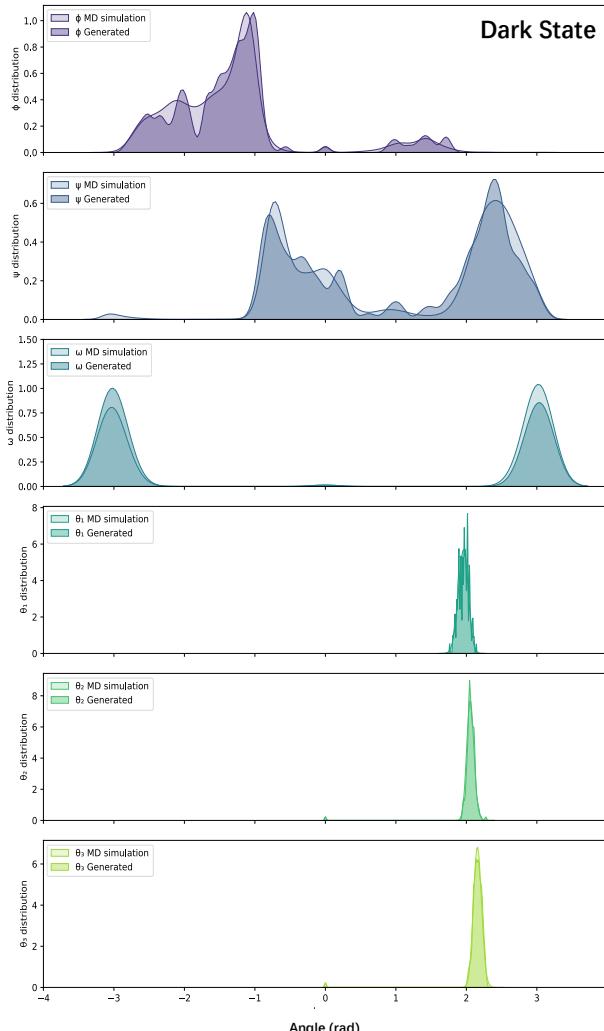
# Angles Distribution

The distributions of the six types of angles between the generated structures and the MD simulation.

(a) Dark state system (b) Light state system

Our model can successfully capture the angle distributions of the original MD simulation for both the dark state and the light state.

Angle distribution between the generated structures and the MD simulation.



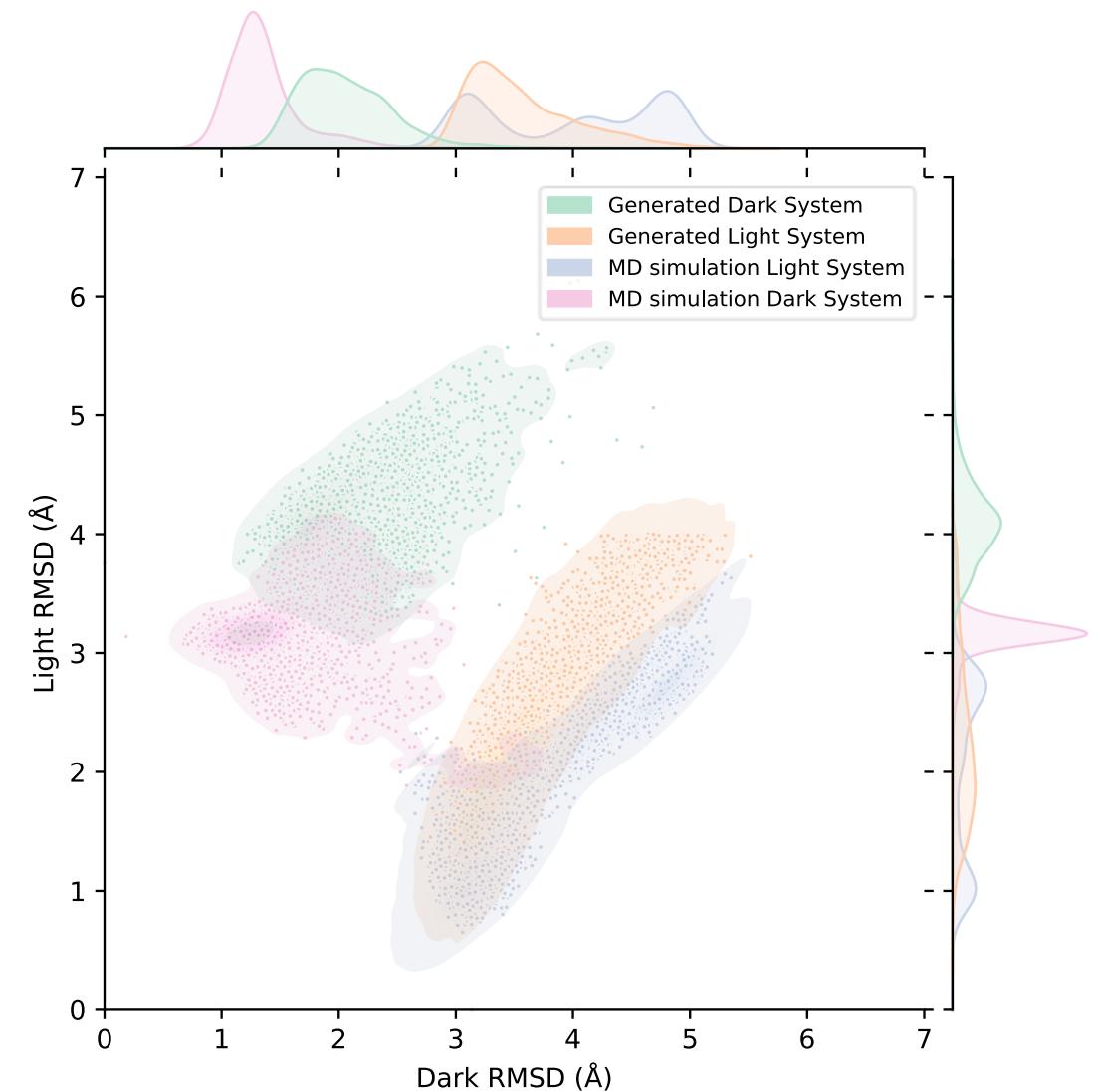
# RMSD Distribution

**RMSD distributions of the generated structures with the training set.**

The x-axis represents RMSD for the dark state while the y-axis represents RMSD for the light state.

**Our model can explore a more extensive conformational space**

RMSD distributions of the generated structures with the training set.

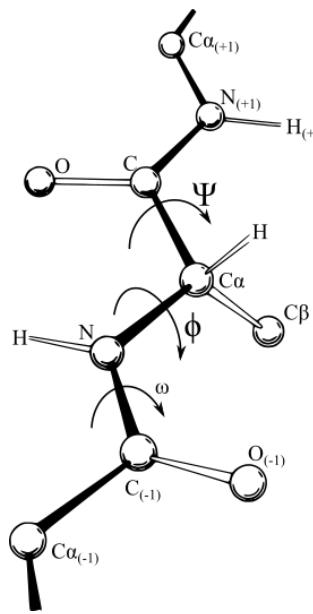


# Ramachandran Plot

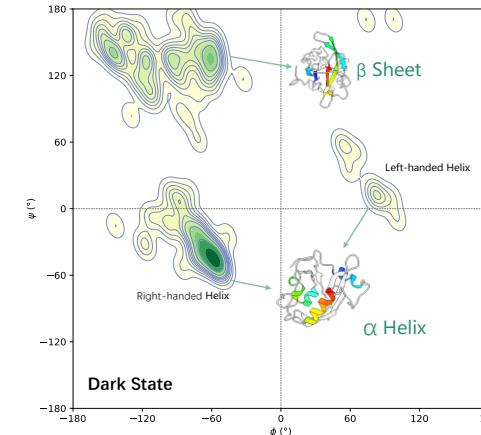
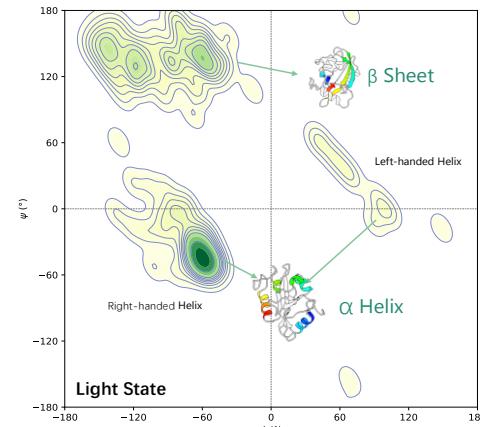
our generated structures possess reasonable second structures

The generated ensembles.

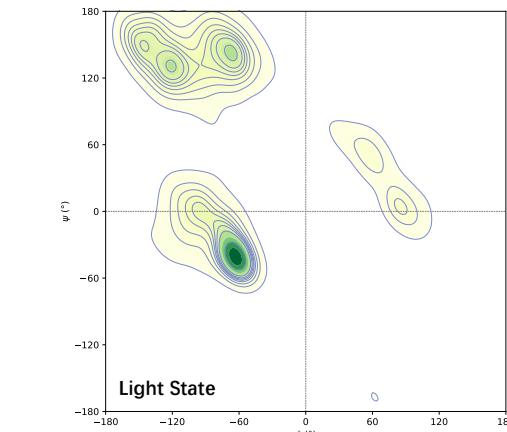
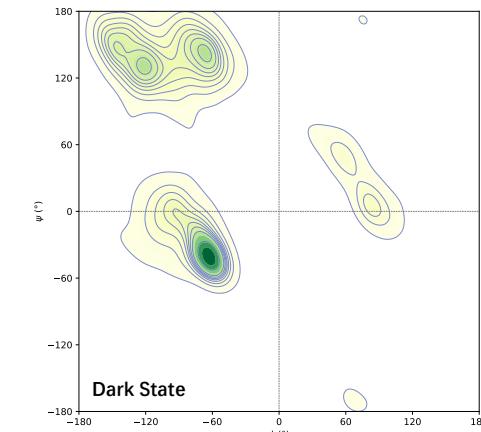
Ramachandran plot [ $\phi, \psi$ ] plot



Our model can successfully capture secondary structure features of protein structures, including alpha helix and beta sheet

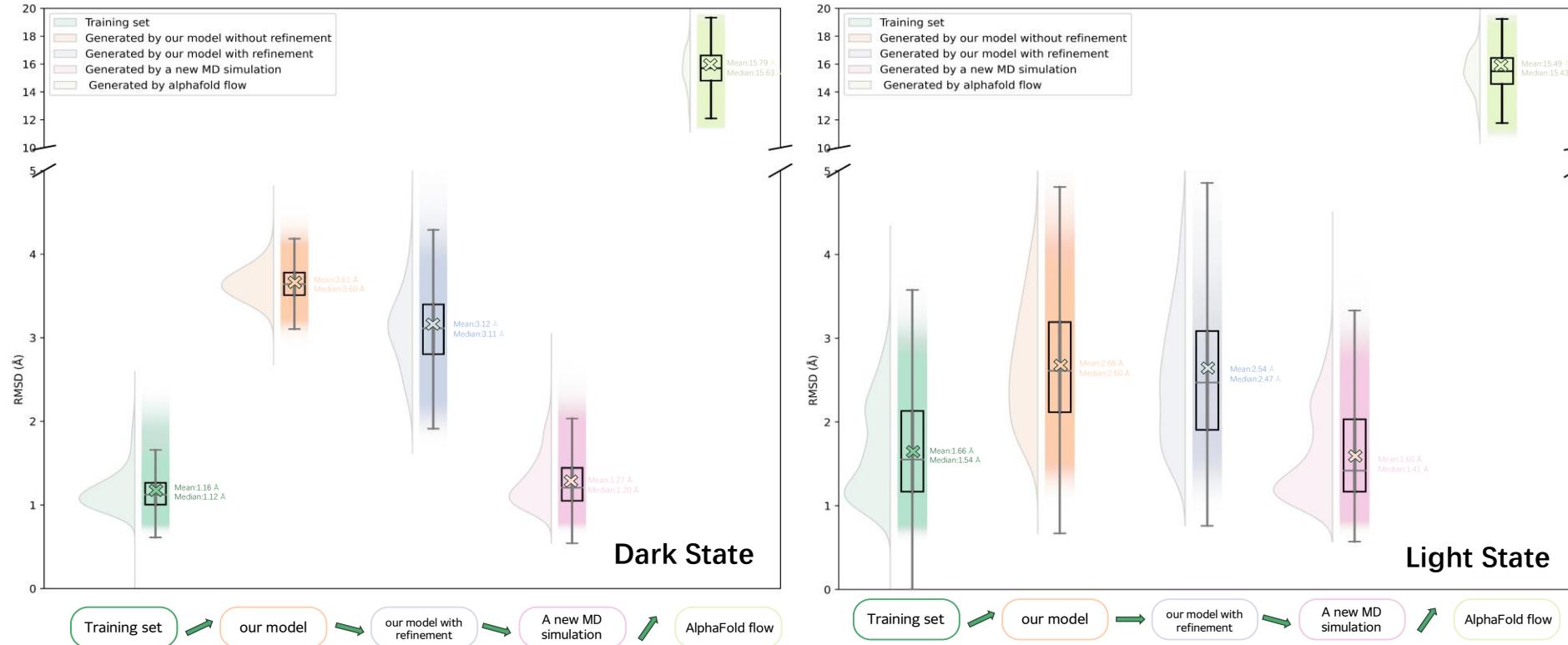


MD simulation ensembles



# Pairwise RMSD Distribution

Compare with Alphaflow



our model, particularly with refinement, produces structures closer to the original MD ensemble than those generated by AlphaFold Flow,

After refinement, there is no bond crash or break

## Physical Reasonable

### Bond Clash:

The distance between any two  $\alpha$ -carbon is **less than** two times the van der Waals radius of  $\alpha$ -carbon with an overlap tolerance:  $\delta = 2 \times 1.7 - 0.4 = 3.0 \text{ \AA}$ .

5 %

### Bond break

Maximum adjacent  $\alpha$ -carbon  
 $> 4.19 \text{ \AA}$

0 %

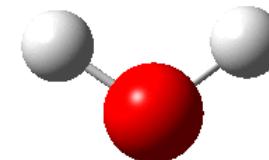
## Future improvement

Combine force field into the model

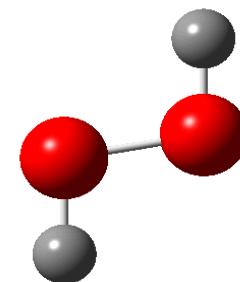
Force fields: bonded + nonbonded



$$\sum_{bond} K_b (b - b_0)^2$$



$$\sum_{angle} K_\theta (\theta - \theta_0)^2$$

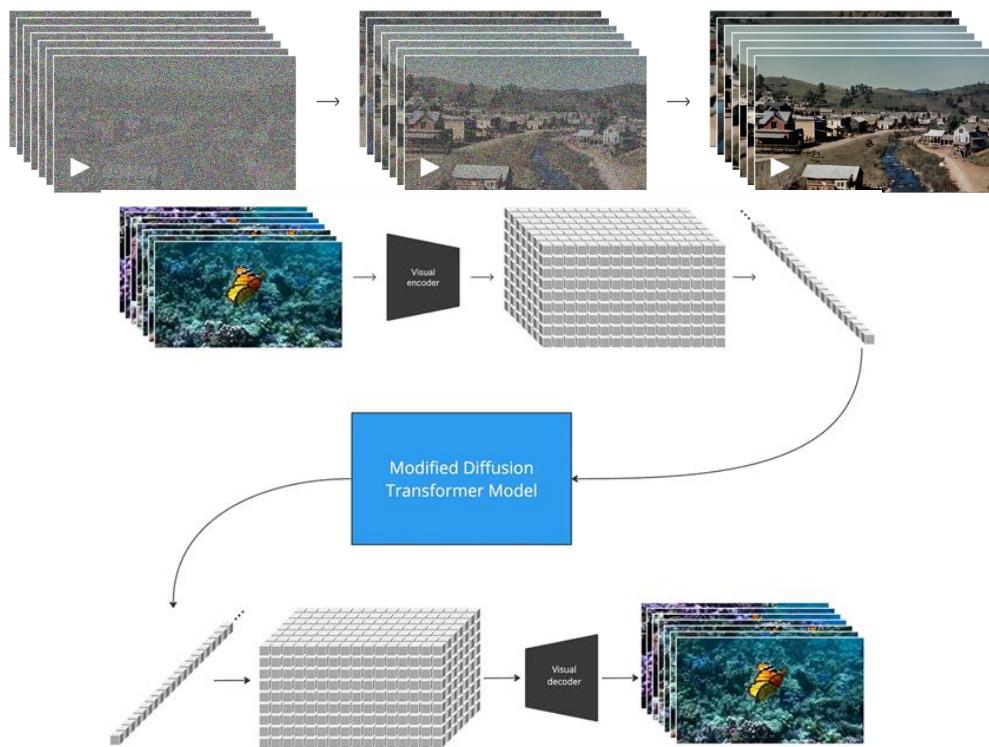


$$\sum_{dihedral} K_\varphi (1 + \cos(n\varphi - \theta))^2$$

## Future improvement

## Video Generations

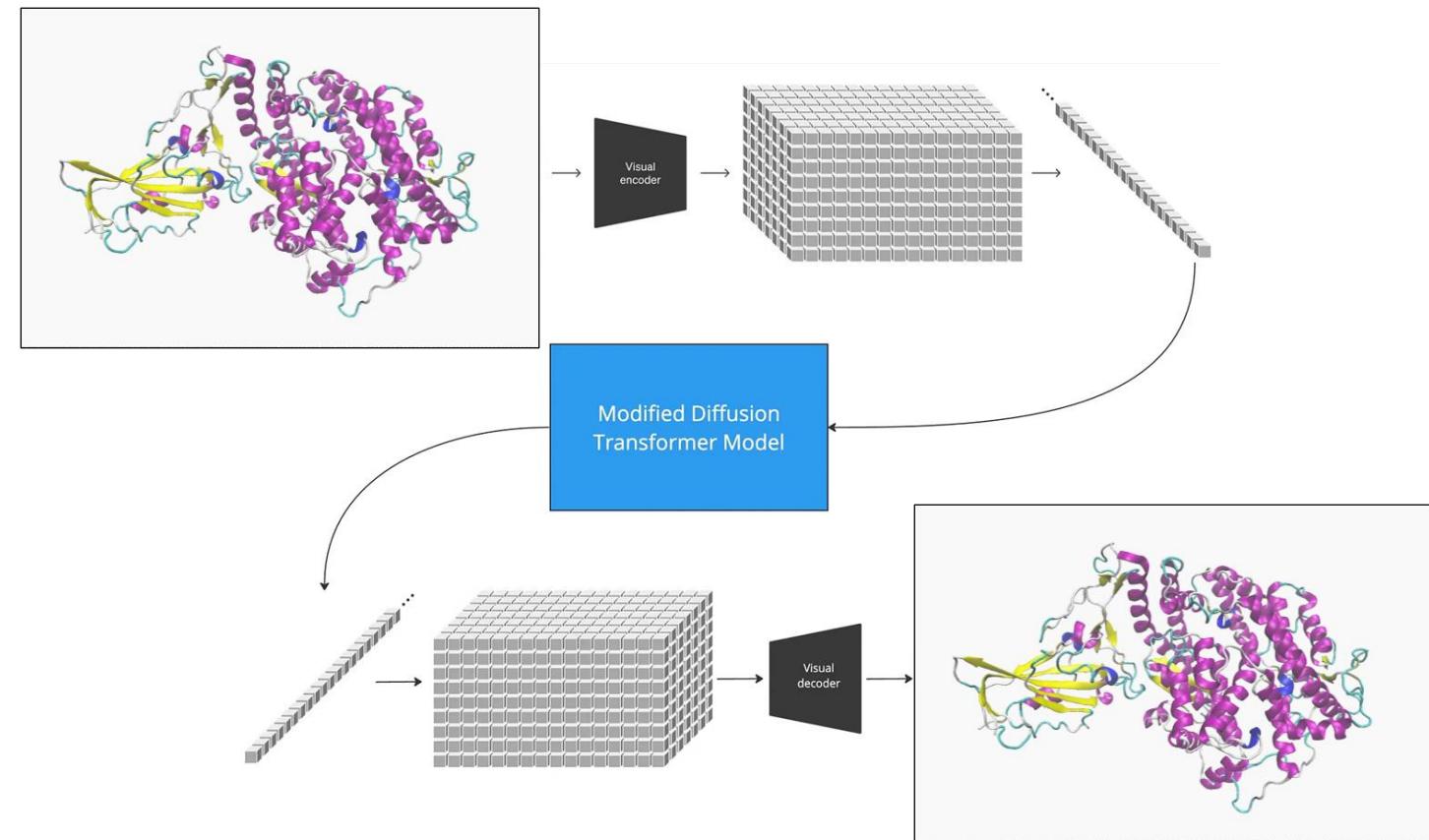
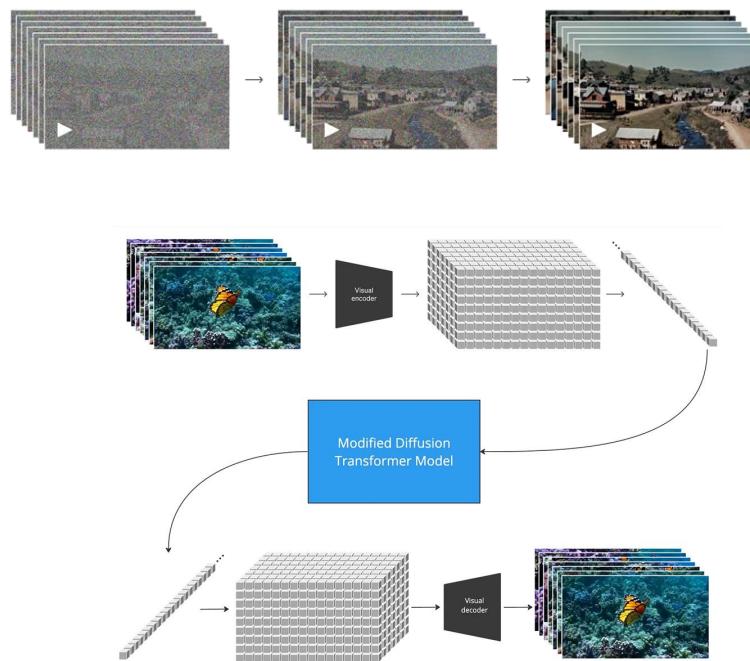
## Sora: Creating video from text



# Future improvement

## Video Generations

### SORA



# Angular-Deviation-Diffuser

pip install Angular-  
Deviation-Diffuser

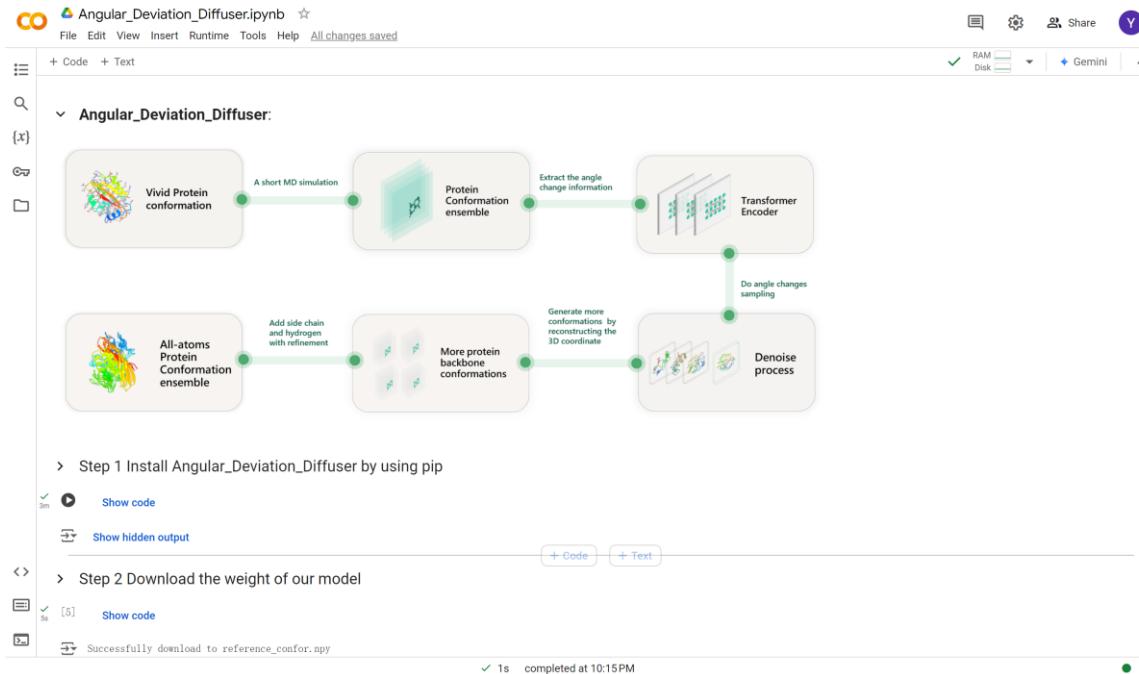


The screenshot shows the PyPI project page for "Angular-Deviation-Diffuser 1.0.5". The page has a blue header with the Python logo and a search bar. Below the header, the project name "Angular-Deviation-Diffuser 1.0.5" is displayed, along with a "Latest version" button and a release date of "Released: Nov 12, 2024". A yellow banner at the top encourages users to join the Python Developers Survey 2024. The main content area includes sections for "Navigation" (Project description, Release history, Download files), "Project description" (Angular Deviation Diffuser), "Overview", and "Background". It also features a "Verified details" section with a green checkmark and a "Maintainers" section showing "yiyangalan".

<https://pypi.org/project/Angular-Deviation-Diffuser/>



# Angular-Deviation-Diffuser



Google  
colab

```
#@title Step 1 Install Angular_Deviation_Diffuser by using pip
!pip install Angular-Deviation-Diffuser
!pip -q install git+https://github.com/sokrrypton/ColabDesign.git@v1.1.1
!python -c 'import pyrosetta_installer; pyrosetta_installer.install_pyrosetta()'

# Step 2 Download the weight of our model
# Step 3 Generate conformations
```

Successfully download to reference\_confor.npy  
Successfully download to model\_para.pth

```
core.pack.interaction_graph.interaction_graph_factory: Instantiating DensePDIInteractionGraph
protocols.relax.FastRelax: CMD: repack 2442.41 0.853644 0.853644 0.14575
protocols.relax.FastRelax: CMD: scale:fa.rep 2592.37 0.853644 0.853644 0.154
protocols.relax.FastRelax: CMD: min 109.719 1.2922 1.2922 0.154
protocols.relax.FastRelax: CMD: coord:cst_weight -101.439 1.2922 1.2922 0.154
protocols.relax.FastRelax: CMD: scale:fa.rep 121.721 1.2922 1.2922 0.30745
core.pack.task: Packer task: initialize from command line()
core.pack.pack_rotamers: built 1770 rotamers at 147 positions.
core.pack.interaction_graph.interaction_graph_factory: Instantiating DensePDIInteractionGraph
protocols.relax.FastRelax: CMD: monoval -78.1738 -1.2922 1.2922 0.30745
```



<https://colab.research.google.com/drive/1paTyFVRMzD4b75DeFjYyXIMV38d1eE1I?usp=sharing>

## Angular-Deviation-Diffuser



Angular\_Deviation\_Diffuser.ipynb

File Edit View Insert Runtime Tools Help All changes saved

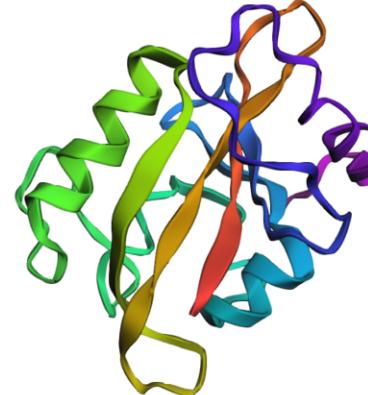
+ Code + Text

> Step 4 Display the conformation's 3D structure

animate: interactive  
color: rainbow  
dpi: 200

Show code

1s completed at 10:15 PM



This screenshot shows a Google Colab notebook titled "Angular\_Deviation\_Diffuser.ipynb". The notebook interface includes a menu bar with File, Edit, View, Insert, Runtime, Tools, Help, and a status bar indicating "All changes saved". A sidebar on the left contains code and text sections, with the text section currently expanded to show the instruction "Step 4 Display the conformation's 3D structure". Below this, there are three configuration options: "animate: interactive", "color: rainbow", and "dpi: 200". A "Show code" link is also present. The main workspace displays a 3D ribbon model of a complex molecule, likely a protein or nucleic acid, rendered in a rainbow color gradient. The bottom of the screen shows a progress bar indicating "1s completed at 10:15 PM".

<https://colab.research.google.com/drive/1paTyFVRMzD4b75DeFjYyXIMV38d1eE1I?usp=sharing>

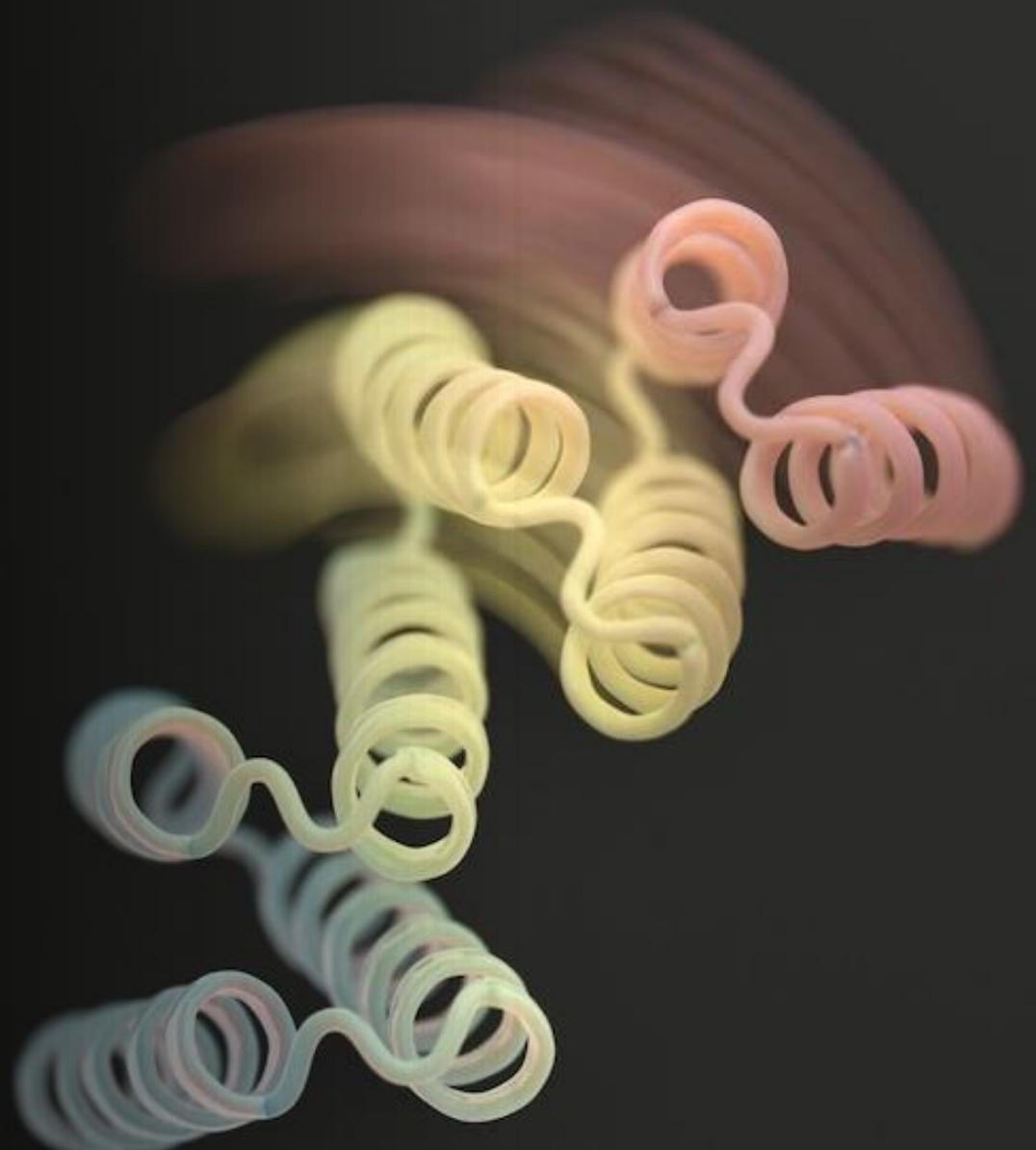
## Part IV

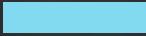
# CONCLUSION

- A transformer-based diffusion model that employs angular deviations to generate protein backbone structures with greater accuracy and efficiency.
- our model significantly outperforms traditional approaches that rely on absolute angle inputs, the use of angular deviations resulted in a more stable and reliable folding process.
- The computational experiments demonstrated that Our model, particularly when refined, has a better performance to Alpha Flow in reproducing VVD protein conformations ,also align closely with MD-generated structures.



SMU<sup>®</sup>





# Acknowledgements

---

- To all group members in Tao Group!
- To my friends and family 
- To HPC in SMU!
- Thanks for attending my Defense!



SMU.<sup>®</sup>

Thank you for your attention!

Questions and suggestions?

