

Presentation | 2022-12-13

# Accelerated Diffusion Models for Protein Structure Generation

---

**Haowei Lin, Shenshen Li, Yuzhe Wang**

Peking University

{linhaowei, lssastar, wangyuzhe\_ccme}@pku.edu.cn

# Outline

---

- Motivation
  - *De novo* Protein Design
  - Diffusion-based Protein Structure Generation
  - Accelerated Sampling for Diffusion Models
- Method
  - Evaluation Metric of Sample Quality
  - Representation of Protein Structures
  - Algorithm Sketch
- Experiments and Conclusion
  - Results and Analysis
  - Comparison of Various Acceleration Methods
  - Conclusion

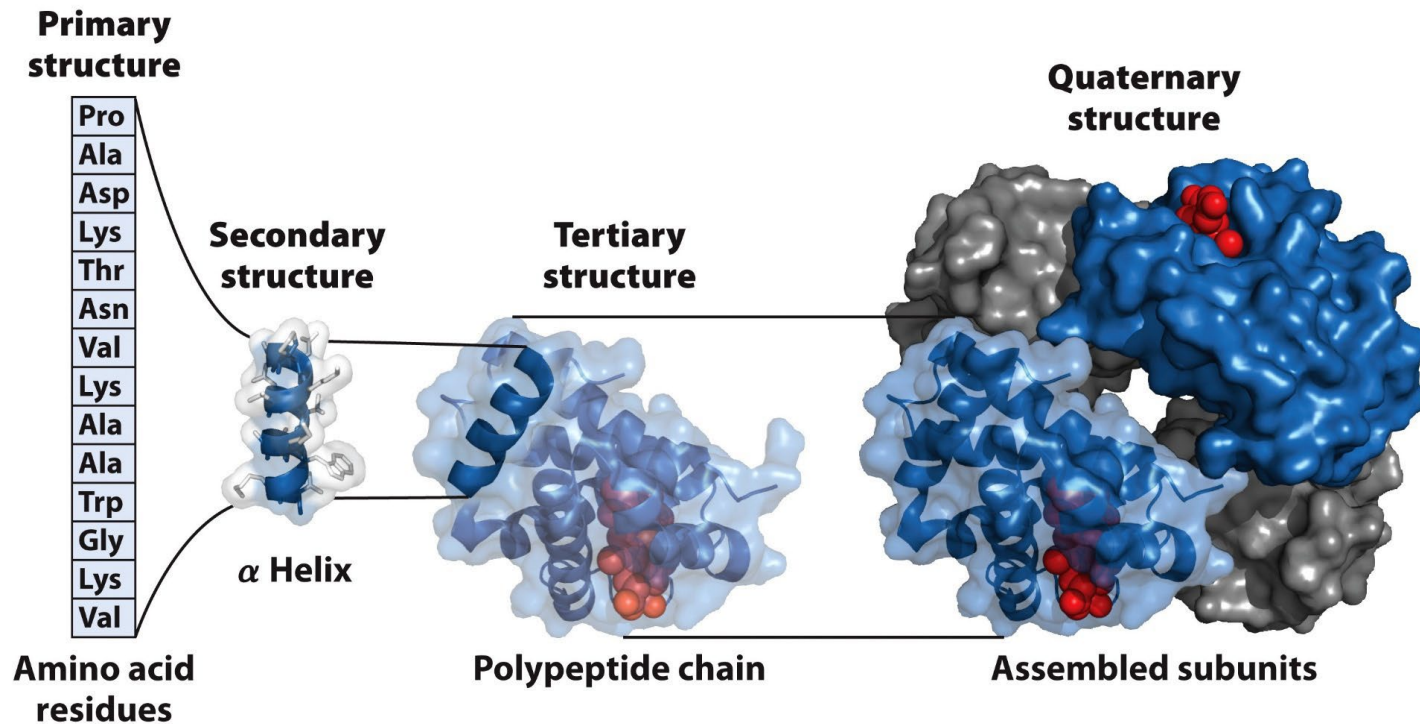
# Outline

---

- Motivation
  - *De novo* Protein Design
  - Diffusion-based Protein Structure Generation
  - Accelerated Sampling for Diffusion Models
- Method
  - Evaluation Metric of Sample Quality
  - Representation of Protein Structures
  - Algorithm Sketch
- Experiments and Conclusion
  - Results and Analysis
  - Comparison of Various Acceleration Methods
  - Conclusion

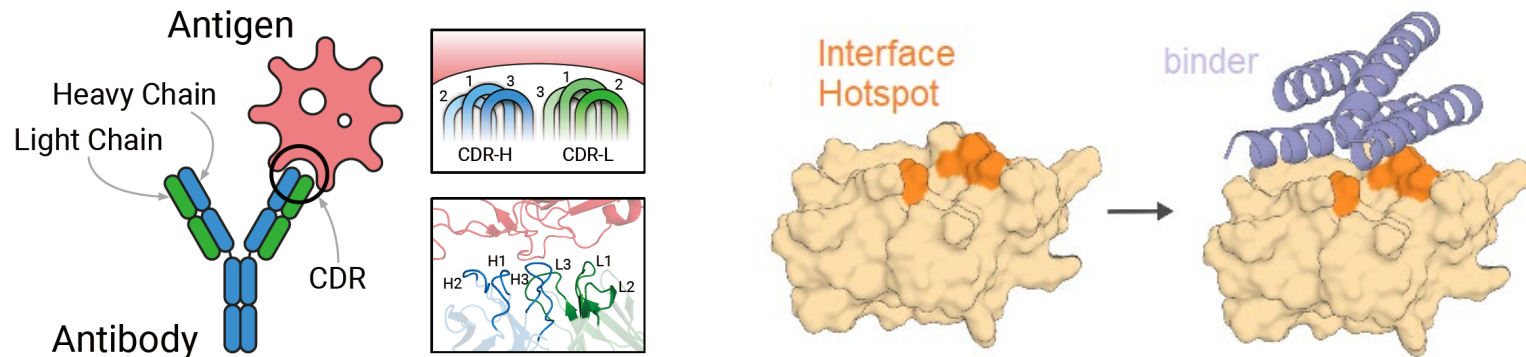
# De novo Protein Design

- Proteins: workhorse molecules of life
  - Protein functions: antibodies, enzymes, messengers...
  - **4 levels** of structure in proteins



# De novo Protein Design

- Protein structure generation
  - Generating novel protein structures for **antibody design** and **drug discovery**



- Challenge of computational *de novo* protein design

✓ AlphaFold2

**Structure prediction**  
Sequence known, structure unknown

Known amino-acid sequence

**Backbone sampling**

Guided by local native sequence

**Side-chain sampling**

Rotamers of native amino acids

Predicted structure

✓ Rosetta fixbb

**Fixed-backbone design**  
Sequence unknown, structure known

Known backbone structure

**Backbone sampling**

None

**Side-chain sampling**

Rotamers of all amino acids

Designed sequence

**De novo design**

Sequence unknown, structure unknown

Architecture definition

**Backbone sampling**

Sequence independent

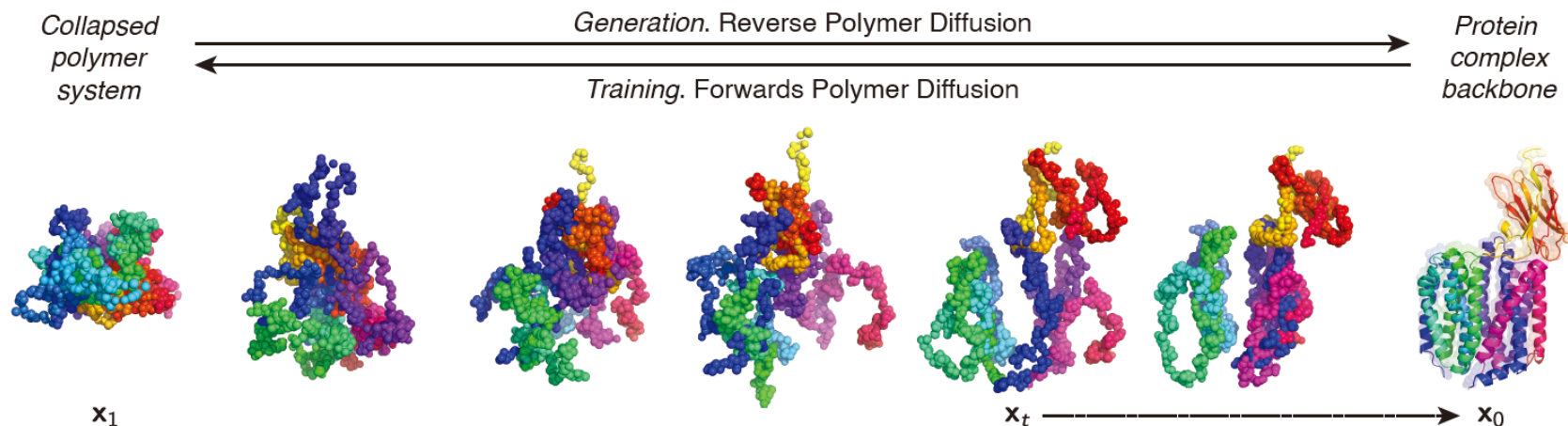
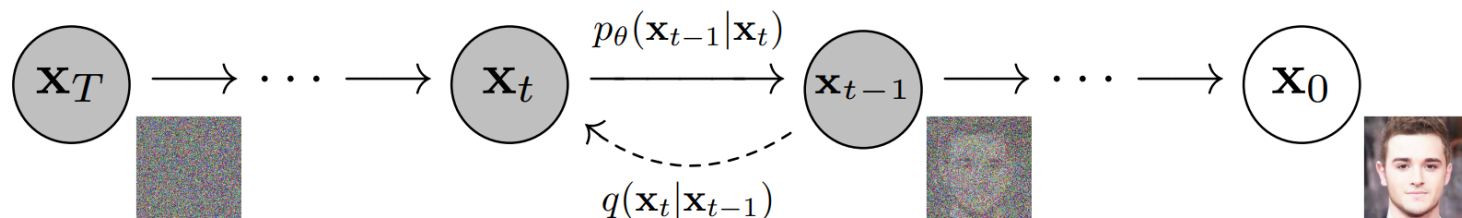
**Side-chain sampling**

Rotamers of all amino acids

Designed backbone and designed sequence

# Diffusion-based Protein Structure Generation

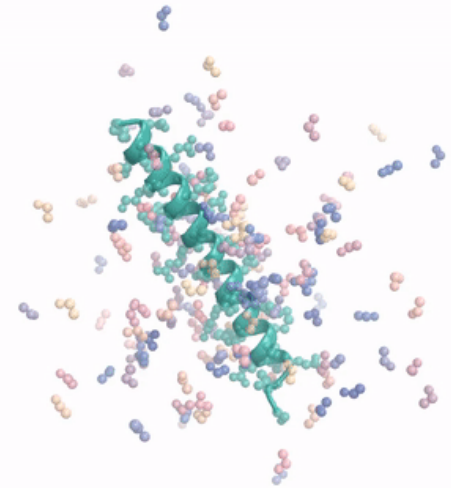
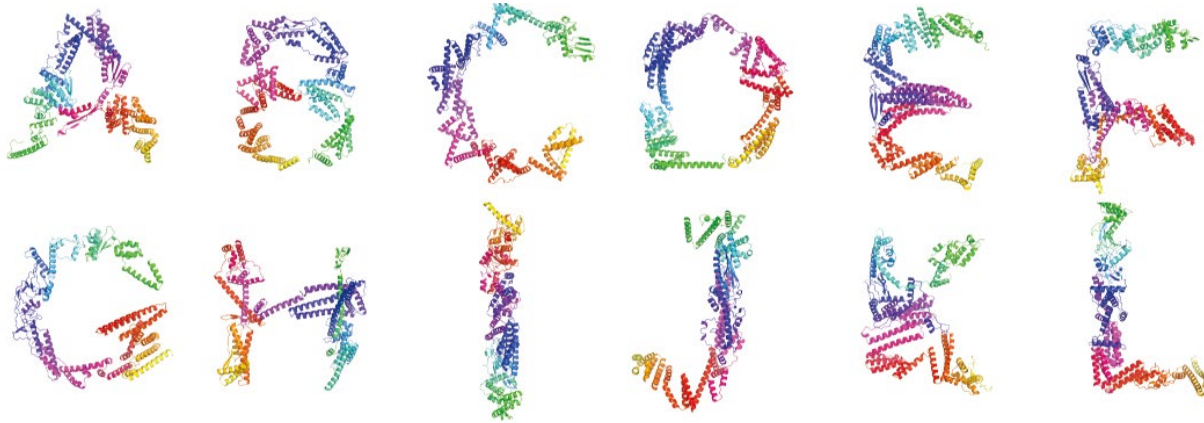
- Generative diffusion of protein structures





# Diffusion-based Protein Structure Generation

- Generative diffusion of protein structures



- Existing problem: **slow inference**
- Project aim: apply feasible **accelerated sampling** methods to **protein structure generation**
  - Improved efficiency of antibody design and drug discovery
  - Provide convenience for model refinement



Ingraham, J., Baranov, M., Costello, Z., Frappier, V., Ismail, A., Tie, S., ... & Grigoryan, G. (2022). Illuminating protein space with a programmable generative model. *bioRxiv*.

<https://www.bakerlab.org/2022/11/30/diffusion-model-for-protein-design/>



# Accelerated Sampling for Diffusion Models

---

- Accelerated sampling methods
  - Training-based: knowledge distillation, self-adaptive noise scheduling, Sample trajectory learning
  - Training-free: **DDIM**, Dynamic stepsize SDE solver, Analytic-DPM, **DPM-Solver**, **DPM-Solver++** (both used in Stable-Diffusion)
- Algorithm sketch of DPM-Solver/DPM-Solver++
  - Essence: fast ODE solver for the **probability flow ODE**

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t), \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T)$$

or equivalent the **diffusion ODE**

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{h}_{\theta}(\mathbf{x}_t, t) := f(t)\mathbf{x}_t + \frac{g^2(t)}{2\sigma_t}\epsilon_{\theta}(\mathbf{x}_t, t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I})$$

in which

$$f(t) = \frac{d \log \alpha_t}{dt}, \quad g^2(t) = \frac{d\sigma_t^2}{dt} - 2\frac{d \log \alpha_t}{dt}\sigma_t^2$$





# Accelerated Sampling for Diffusion Models

- Accelerated sampling methods
- Algorithm sketch of DPM-Solver/DPM-Solver++
  - Decouple linear parts and nonlinear parts of the diffusion ODE by **variation of constants** formula to formulate the exact solution

$$\mathbf{x}_t = e^{\int_t^s f(\tau) d\tau} \mathbf{x}_s + \int_s^t \left( e^{\int_t^\tau f(\tau) d\tau} \frac{g^2(\tau)}{2\sigma_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) \right) d\tau$$

- Rewrite the solution by **change-of-variable** for  $\lambda_t := \log(\alpha_t/\sigma_t)$

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \hat{\epsilon}_\theta(\hat{\mathbf{x}}_\lambda, \lambda) d\lambda$$

- Apply  $(k-1)$ -th order Taylor expansion to the NN term and compute the integral **analytically** using integration-by-parts to obtain **DPM-Solver- $k$**

$$\begin{aligned} \mathbf{x}_{t_{i-1} \rightarrow t_i} = & \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \hat{\epsilon}_\theta^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda \\ & + \mathcal{O}((\lambda_{t_i} - \lambda_{t_{i-1}})^{k+1}) \end{aligned}$$

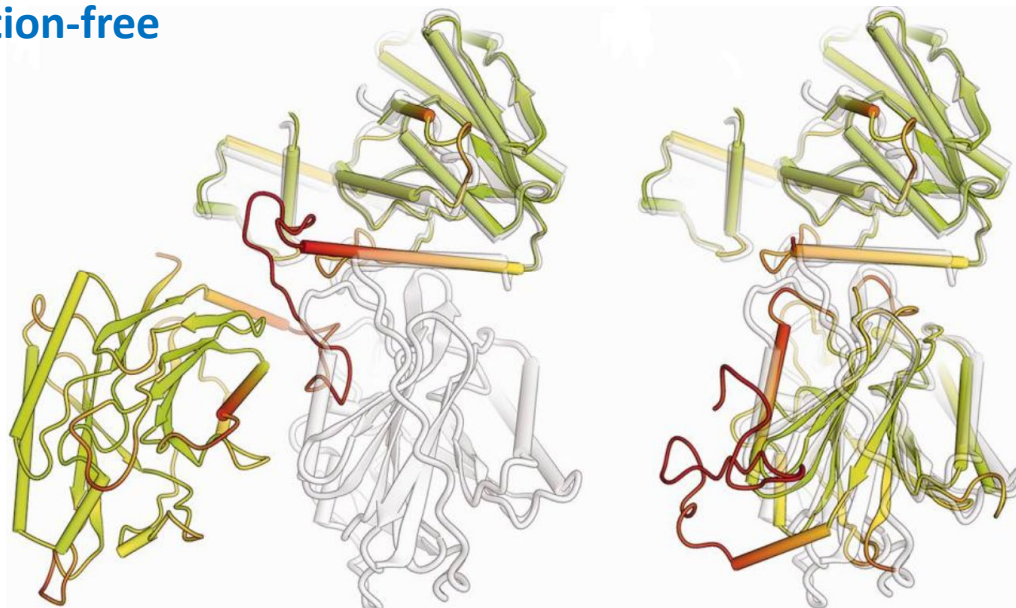
# Outline

---

- Motivation
  - *De novo* Protein Design
  - Diffusion-based Protein Structure Generation
  - Accelerated Sampling for Diffusion Models
- Method
  - Evaluation Metric of Sample Quality
  - Representation of Protein Structures
  - Algorithm Sketch
- Experiments and Conclusion
  - Results and Analysis
  - Comparison of Various Acceleration Methods
  - Conclusion

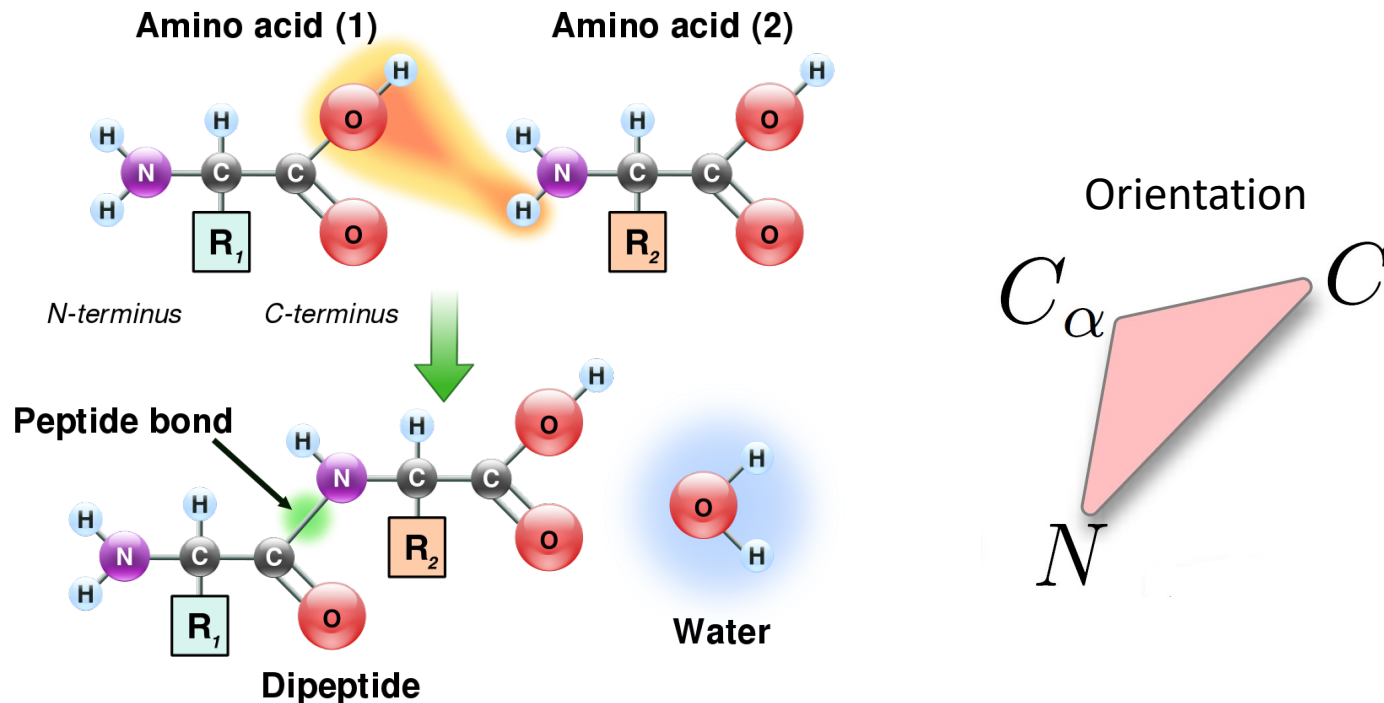
# Evaluation Metric of Sample Quality

- Sample quality assessment
  - Image/text generation: perplexity, FID, inception score etc.
  - Quality assessment of generated protein structures? → structure generation of **given amino acid sequence** (whose **ground-truth** structure is known)
- Evaluation metrics: **IDDT** (local **D**istance **D**ifference **T**est)
  - IDDT measures how well the environment in a reference structure is **reproduced** in a protein model (fraction of reproduced atom pairs, ↑)
  - **Superposition-free**



# Representation of Protein Structures

- Representation of an amino acid in a protein complex
  - Amino acid type:**  $s_i \in \{\text{ACDEFGHIKLMNPQRSTVWY}\}$
  - $\text{C}\alpha$  atom coordinate:**  $\mathbf{x}_i \in \mathbb{R}^3$
  - Orientation:**  $\mathbf{O}_i \in \text{SO}(3)$



# Algorithm Sketch

---

- Diffusion for C $\alpha$  coordinates
  - Forward diffusion

$$q(\mathbf{x}_j^t | \mathbf{x}_j^{t-1}) = \mathcal{N}(\mathbf{x}_j^t | \sqrt{1 - \beta^t} \cdot \mathbf{x}_j^{t-1}, \beta^t \mathbf{I})$$

$$q(\mathbf{x}_j^t | \mathbf{x}_j^0) = \mathcal{N}(\mathbf{x}_j^t | \sqrt{\bar{\alpha}^t} \cdot \mathbf{x}_j^0, (1 - \bar{\alpha}^t) \mathbf{I})$$

- Generative process

$$p(\mathbf{x}_j^{t-1} | \mathbf{x}^t) = \mathcal{N}(\mathbf{x}_j^{t-1} | \boldsymbol{\mu}(\mathbf{x}^t, t), \beta^t \mathbf{I})$$

$$\boldsymbol{\mu}(\mathbf{x}^t, t) = \frac{1}{\sqrt{\alpha^t}} \left( \mathbf{x}_j^t - \frac{\beta^t}{\sqrt{1 - \bar{\alpha}^t}} G(\mathbf{x}^t, t)[j] \right)$$

- Objective function

$$L^t = \mathbb{E} \left( \frac{1}{m} \sum_j \|\epsilon_j - G(\mathbf{x}^t, t)\|^2 \right)$$

# Algorithm Sketch

---

- Diffusion for amino acid orientations

- Forward diffusion

$$q(\mathbf{O}_j^t | \mathbf{O}_j^0) = \mathcal{IG}_{\text{SO}(3)}(\mathbf{O}_j^t | \text{ScaleRot}(\sqrt{\bar{\alpha}^t}, \mathbf{O}_j^0), 1 - \bar{\alpha}^t)$$

- $\mathcal{IG}_{\text{SO}(3)}$  denotes the isotropic Gaussian distribution on  $\text{SO}(3)$

- Generative process

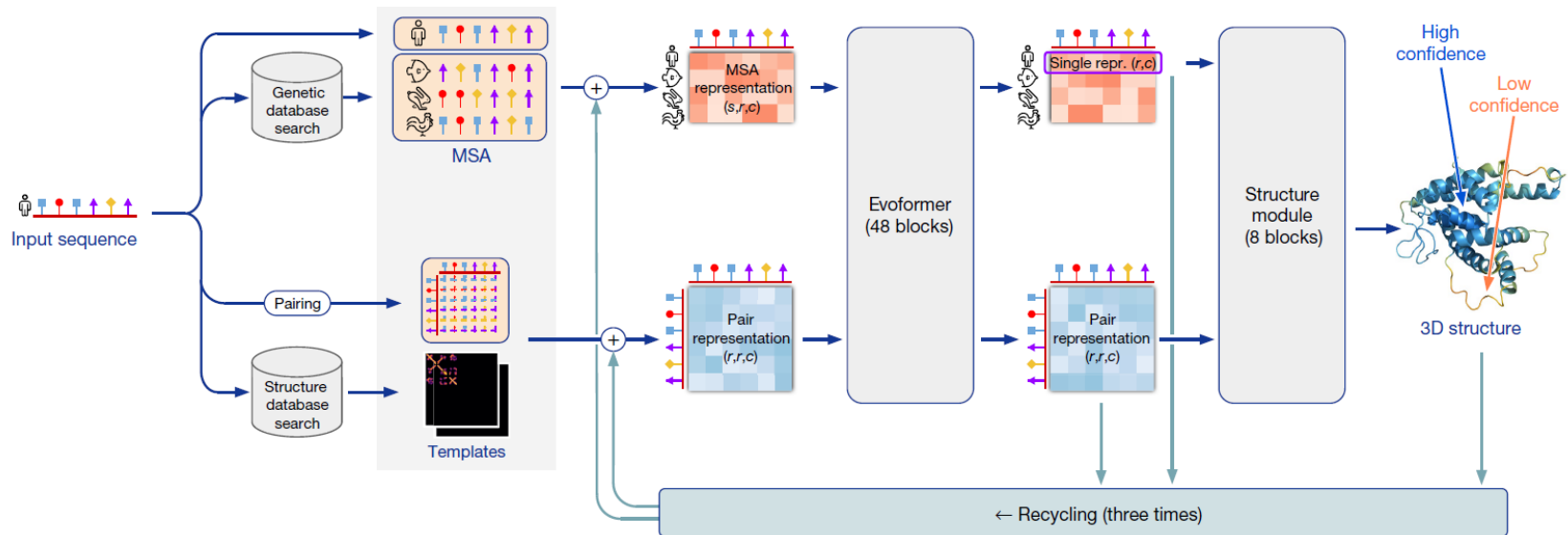
$$p(\mathbf{O}_j^{t-1} | \mathbf{O}^t) = \mathcal{IG}_{\text{SO}(3)}(\mathbf{O}_j^{t-1} | H(\mathbf{O}^t, t)[j], \beta^t)$$

- Objective function

$$L^t = \mathbb{E} \left( \frac{1}{m} \sum_j \| (\mathbf{O}_j^0)^T H(\mathbf{O}_j^t, t)[j] - \mathbf{I} \|^2 \right)$$

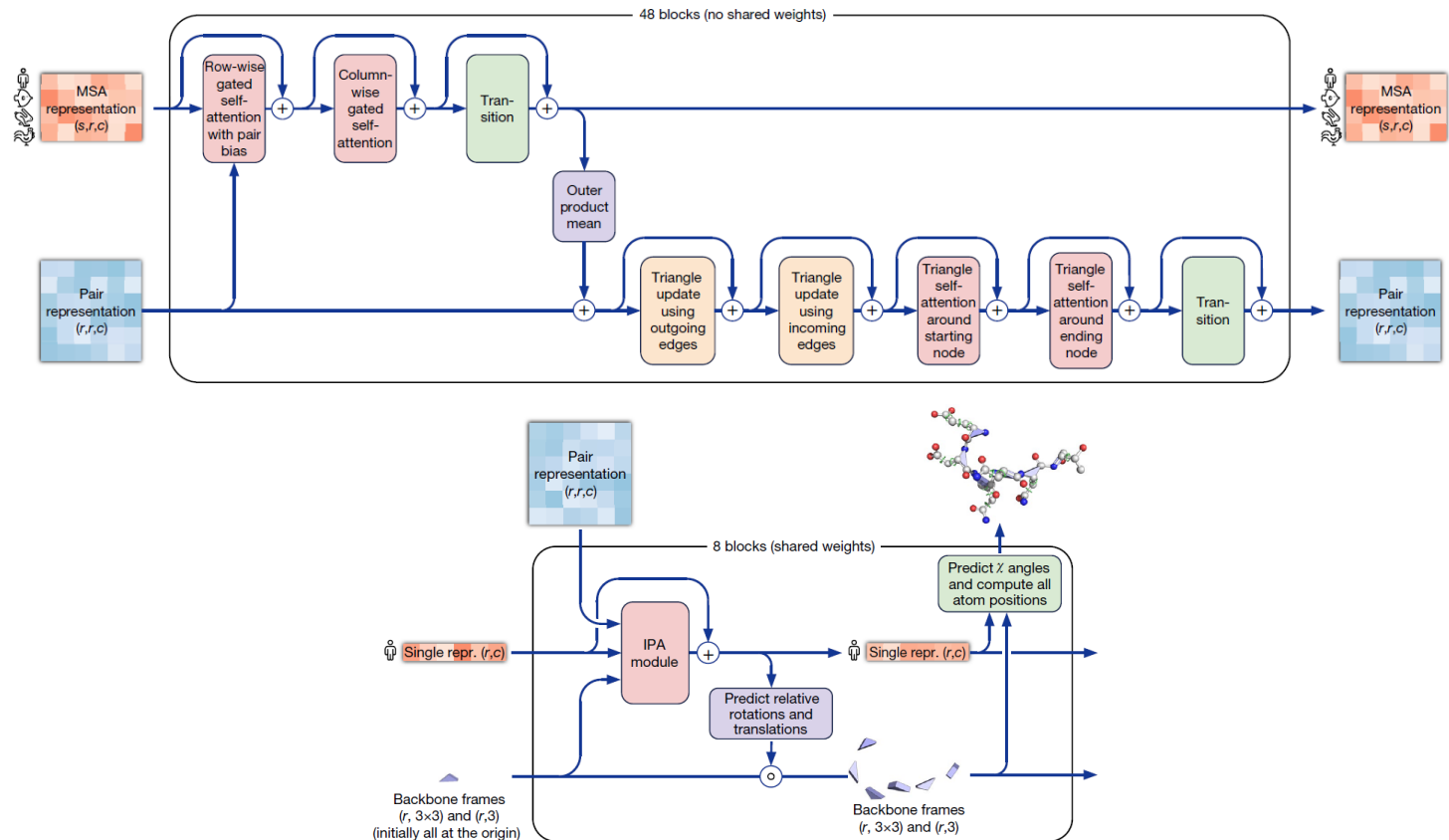
# Algorithm Sketch

- Sequence condition
  - Latent feature of protein sequences: pretrained **AlphaFold2 Evoformer**
  - Model architecture backbone: pretrained **AlphaFold2 structure module**



# Algorithm Sketch

- Sequence condition
  - Latent feature of protein sequences: pretrained **AlphaFold2 Evoformer**
  - Model architecture backbone: pretrained **AlphaFold2 structure module**





# Outline

---

- Motivation
  - *De novo* Protein Design
  - Diffusion-based Protein Structure Generation
  - Accelerated Sampling for Diffusion Models
- Method
  - Evaluation Metric of Sample Quality
  - Representation of Protein Structures
  - Algorithm Sketch
- Experiments and Conclusion
  - Results and Analysis
  - Comparison of Various Acceleration Methods
  - Conclusion

# Results and Analysis

---

- Baseline performances of the vanilla diffusion model
  - Test set: **CAMEO dataset** which consists 146 of the most recent single-chain proteins
  - Case study: protein **PDB\_ID = 7BI4\_A** (on which all accelerated sampling algorithms are tested)

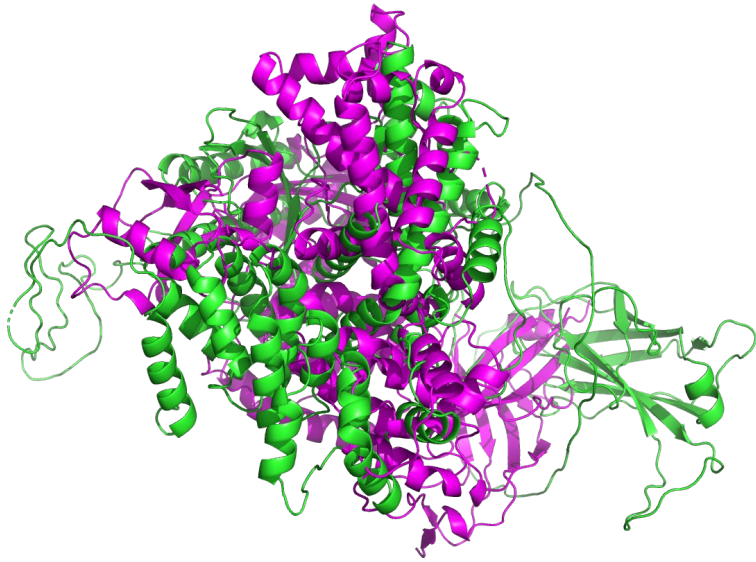
vanilla	1000 NFE	500 NFE	100 NFE	50 NFE	25 NFE	1 NFE
lddt(average)	0.824	0.820	0.818	0.8	0.74	0.30
lddt(7BI4_A)	0.70	0.685	0.676	0.65	0.45	0.18
time(CAMEO)	42331	21888	5356	3279	2321	1273
time(7BI4_A)	1257.3	648.2	158.7	96.9	69.8	40.1

The test results for the vanilla diffusion model

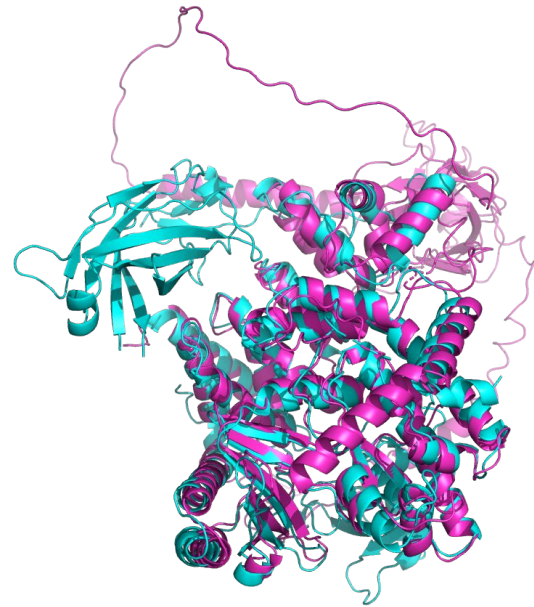
# Results and Analysis

---

- The 3D structure of **Real** and predicted (**50 NFE** / **1000 NFE**) 7Bl4\_A



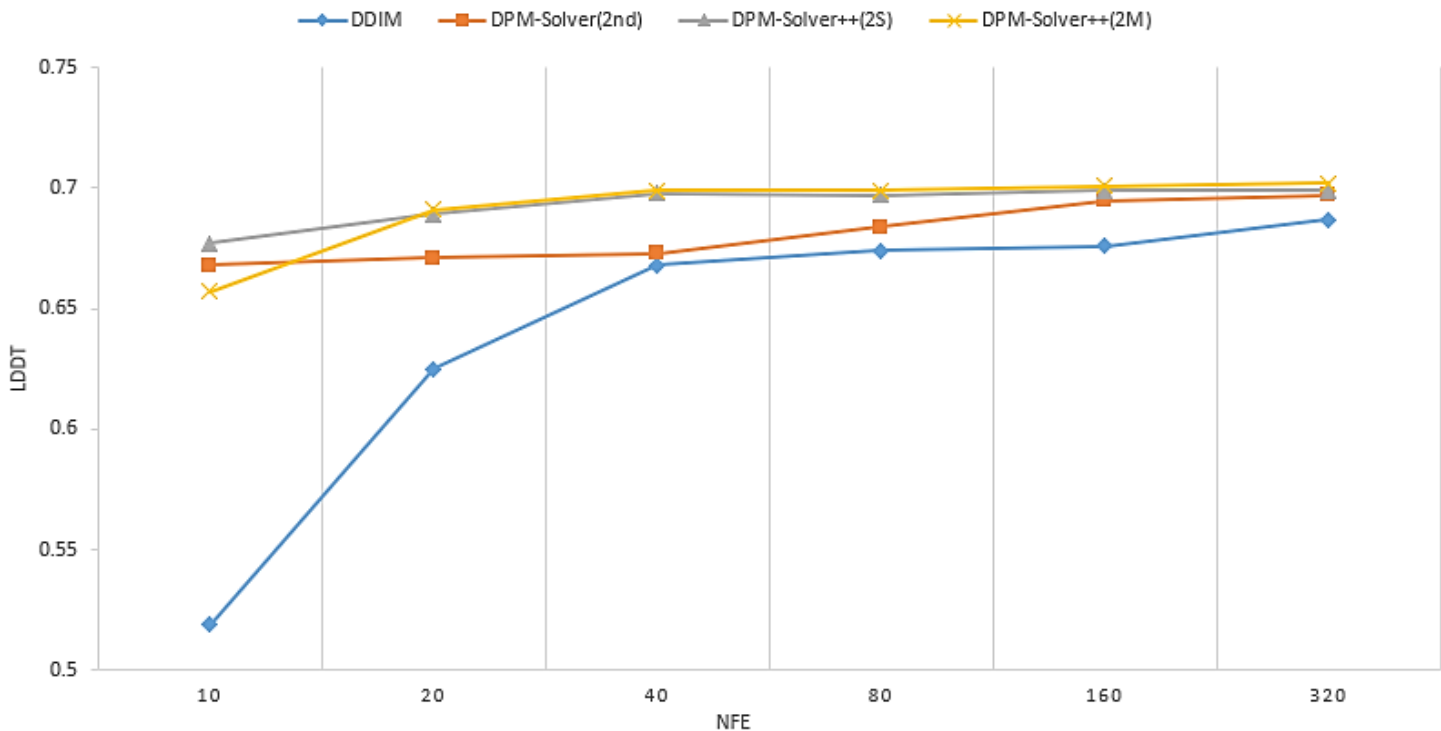
50 NFE, lddt=0.65, time=96.9 s



1000 NFE, lddt=0.70, time=1257.3 s

# Comparison of Various Acceleration Methods

- The results of 7Bl4\_A
  - **DPM-Solver++** converges at about 40 NFE and outperforms all other algorithms
  - DPM-Solver performs slightly worse than DPM-Solver++ but beats DDIM



# Conclusion

---

- For 7Bl4\_A:
  - The acceleration methods significantly improve the convergence speed of the model by a large margin (**1000 NFE → 40 NFE, 1200 s → 60 s**)
  - **DPM-Solver++** performs best when applied to diffusion-based protein structure generation
  - DDIM performs the worst (no convergence at 320 NFE)

# Thank you

---

**Haowei Lin, Shenshen Li, Yuzhe Wang**

Peking University

{linhaowei, lssastar, wangyuzhe\_ccme}@pku.edu.cn