



中國人民大學
RENMIN UNIVERSITY OF CHINA



高瓴人工智能学院
Gaoling School of Artificial Intelligence

Generalized Protein Pocket Generation with Prior-Informed Flow Matching

Zaixi Zhang, Marinka Zitnik, Qi Liu
(NeurIPS 2024 Spotlight)

Fanmeng Wang

2024-10-15

Outline

- Introduction
- Method
- Experiment
- Conclusion



➤ Introduction

➤ Method

➤ Experiment

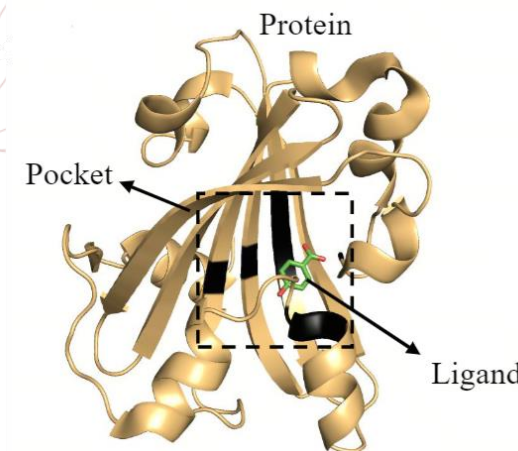
➤ Conclusion



Introduction

■ Protein Pocket Designing

- **Proteins** are the fundamental building blocks of living organisms, and they often interact with **ligands** (e.g., small molecules, nucleic acids, and peptides) to execute their functions.
- Therefore, one critical step in functional protein design is to **design protein pockets**, which refers to the protein interface binding with the ligand.
- However, the complexity of ligand-protein interactions, the variability of protein sidechains, and sequence-structure relationships pose great challenges for pocket design.





Introduction

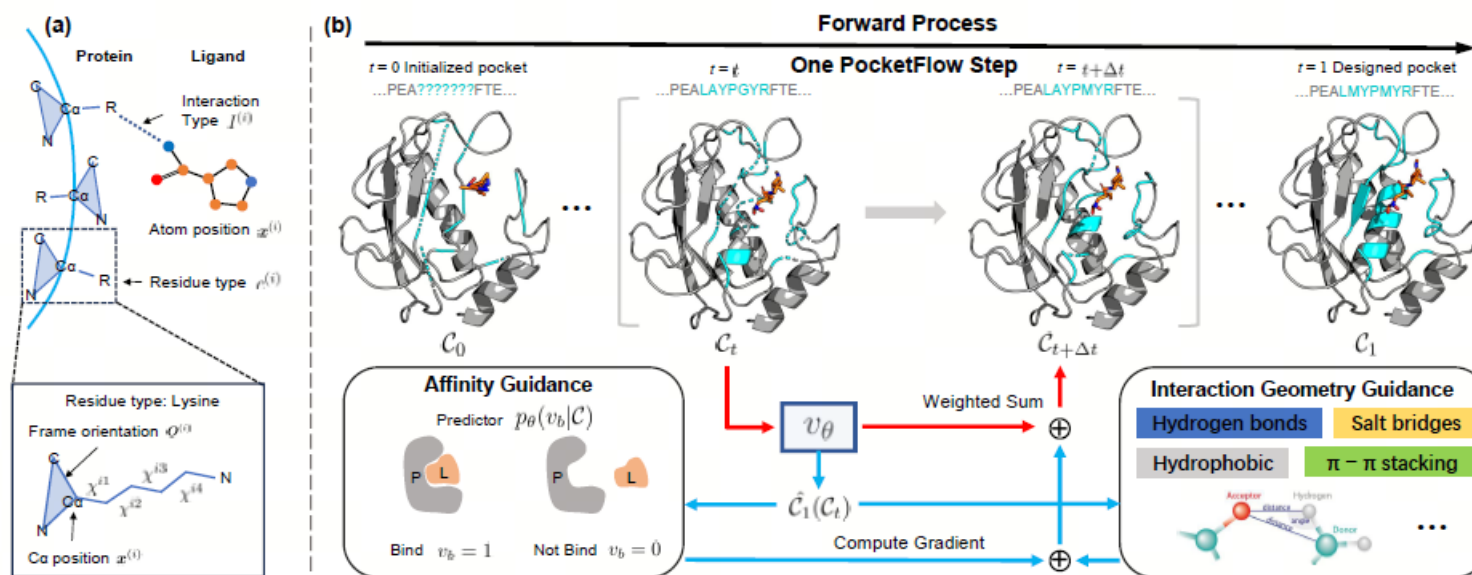
■ Existing works

- **Traditional methods** for protein pocket design mainly focus on physics modeling or template-matching. However, the involved physical energy calculation or substructure enumeration could be quite **time-consuming**.
- Recent advancements in pocket design have benefited a lot from **deep learning-based methods**.
 - On the one hand, these methods often **overlook essential domain knowledge**, such as the protein-ligand interactions and the geometric constraints governing them.
 - On the other hand, these methods **are restricted to small molecule ligands**, omitting other important ligand types such as nucleic acids and peptides.

Introduction

■ PocketFlow

- In this work, we propose **PocketFlow**, a **protein-ligand interaction prior-informed flow matching model** for protein pocket generation.



Generalized tasks
Strong performance

Figure 1: (a) Parameterization of protein-ligand complex. (b) Illustration of PocketFlow forward process. The affinity and interaction geometry guidance are proposed to improve affinity and structural validity. The red/blue lines denote the unconditional/guidance paths respectively.

➤ Introduction

➤ **Method**

➤ Experiment

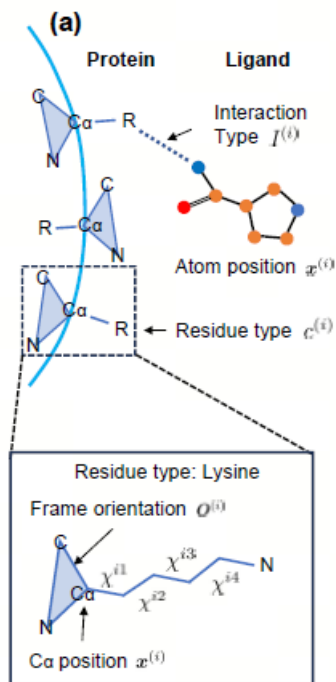
➤ Conclusion





Method

Notations

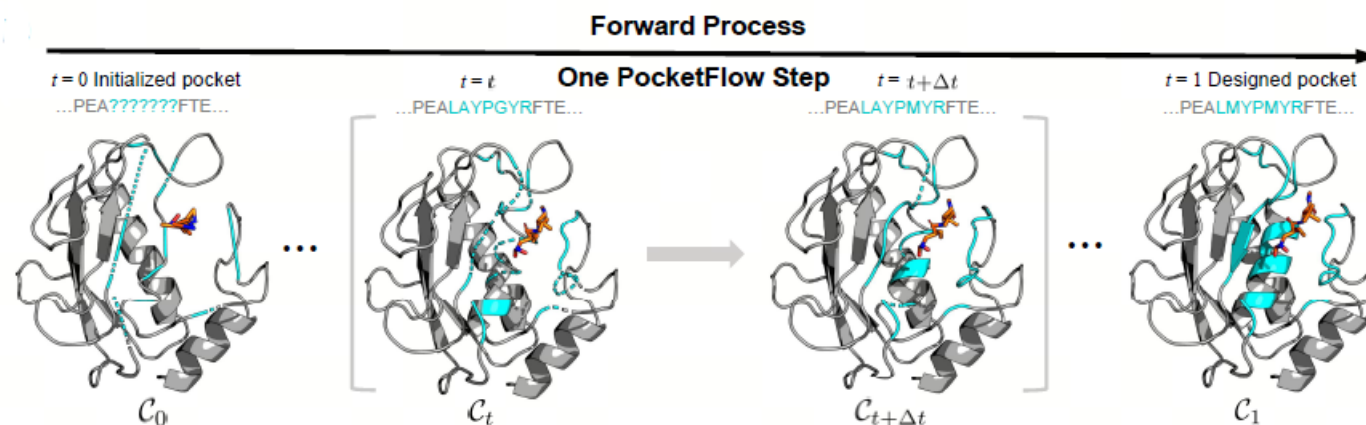


Notations. As shown in Figure 1(a), we model protein-ligand complex as $\mathcal{C} = \{\mathcal{P}, \mathcal{G}\}$ consisting of protein \mathcal{P} and ligand \mathcal{G} (small molecule as an example). Protein \mathcal{P} is composed of a sequence of residues (amino acids) with residue types denoted $c^{(i)} \in \mathbb{R}^{20}$. Consistent with [92, 87], the protein pocket $\mathcal{R} \subset \mathcal{P}$ is defined as the subset of residues closest to the ligand atoms under a threshold δ (e.g., 3.5 Å). In a residue, the backbone structure (consisting of C_α , N , C , O) is parameterized with C_α position $x^{(i)} \in \mathbb{R}^3$ and a frame orientation matrix $O^{(i)} \in SO(3)$ following [43, 84]. The sidechain is parameterized with maximal 4 torsion angles $\chi^{(i)} = \{\chi^{i1}, \chi^{i2}, \chi^{i3}, \chi^{i4}\} \in [0, 2\pi)^4$. Given these key parameters, the full atom protein structure can be derived with the ideal frame coordinates and the sidechain bond length/angles [43]. The protein-ligand interaction type for each residue is marked as $I^{(i)} \in \mathbb{R}^5$ (Hydrogen bond, Salt bridge, Hydrophobic, π - π stacking, no interaction). A pocket with N_r residues can be compactly represented as $\mathcal{R} = \{c^{(i)}, x^{(i)}, O^{(i)}, \chi^{(i)}, I^{(i)}\}_{i=1}^{N_r}$. As for the ligand, we use a generalized atom-level representation that accommodates various modalities including small molecules, peptides, and RNA. The atom types and bonding information between atoms are given and PocketFlow predicts the N_l ligand atom coordinates (also denoted as $x^{(i)}$ for conciseness).

Method

■ Problem Definition

- PocketFlow co-designs **residue types, 3D structures of the protein pocket, 3D structures of the ligand** conditioned on the ligand (could be small molecules, nucleic acids, peptides, etc.) and protein scaffold (the other parts of protein besides the pocket region).
- Here, each atom in the ligand is treated as an individual residue, but only coordinate (similar to C_α coordinates) need to be predicted.



Method

■ PocketFlow

- For conditional flow matching, the key point is to learn the **conditional vector field**, which can be calculated based on **prior distribution**, sample distribution and the **conditional flow**.

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1], p_1(x_1), p_t(x|x_1)} \|v_\theta(x, t) - u_t(x|x_1)\|_g^2$$

$$\frac{d}{dt} \psi_t(x) = u_t(\psi_t(x) | x_1)$$

- Here, considering that protein-ligand complex has **many different components**, we need to define PocketFlow for these different components (backbone, sidechain, and residue/interaction types).

$$\mathcal{R} = \{c^{(i)}, x^{(i)}, O^{(i)}, \chi^{(i)}, I^{(i)}\}_{i=1}^{N_r}$$

Define prior distribution and conditional flow for these different components

Method

■ PocketFlow on Backbone

➤ Prior Distribution:

- For \mathbf{C}_α coordinates: the **linear interpolation and extrapolation** based on the known coordinates of neighboring scaffold residues
- For **frame orientation matrix**: the **uniform distribution** on $\text{SO}(3)$

➤ Conditional flow:

$$\mathbf{x}_t^{(i)} = (1 - t)\mathbf{x}_0^{(i)} + t\mathbf{x}_1^{(i)} \quad \mathbf{O}_t^{(i)} = \exp_{\mathbf{O}_0^{(i)}}(t \log_{\mathbf{O}_0^{(i)}}(\mathbf{O}_1^{(i)}))$$

➤ Loss:

$$\mathcal{L}_{coord}(\theta) = \mathbb{E}_{t, p_1(\mathbf{x}_1), p_0(\mathbf{x}_0), p_t(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_1)} \frac{1}{N_r + N_l} \sum_{i=1}^{N_r + N_l} \left\| v_\theta^{(i)}(\mathbf{x}_t^{(i)}, t) - \mathbf{x}_1^{(i)} + \mathbf{x}_0^{(i)} \right\|_2^2, \quad (1)$$

$$\mathcal{L}_{ori}(\theta) = \mathbb{E}_{t, p_1(\mathbf{O}_1), p_0(\mathbf{O}_0), p_t(\mathbf{O}_t | \mathbf{O}_0, \mathbf{O}_1)} \frac{1}{N_r} \sum_{i=1}^{N_r} \left\| v_\theta^{(i)}(\mathbf{O}_t^{(i)}, t) - \frac{\log_{\mathbf{O}_0^{(i)}}(\mathbf{O}_1^{(i)})}{1 - t} \right\|_{\text{SO}(3)}^2, \quad (2)$$

where we additionally consider N_l ligand atom coordinates in $\mathcal{L}_{coord}(\theta)$, for which we use Gaussian distribution at the center of ligand mass as the prior distribution.

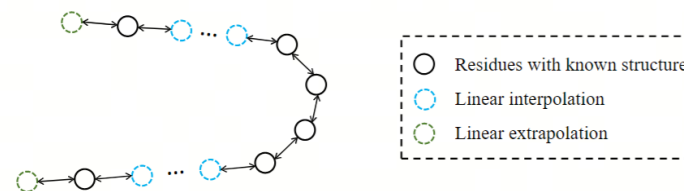


Figure 5: Structure initialization based on interpolation and extrapolation



Method

■ PocketFlow on Sidechain

As described in Sec. 3.1, the sidechain conformation of each residue can be represented as maximally four torsion angles $\chi^{(i)} = \{\chi^{i1}, \chi^{i2}, \chi^{i3}, \chi^{i4}\} \in [0, 2\pi)^4$. In a pocket with N_r residues, the sidechain torsion angles form a hypertorus \mathbb{T}^{4N_r} , which is the quotient space $\mathbb{R}^{4N_r} / 2\pi\mathbb{Z}^{4N_r}$ with the equivalence relation: $\chi = (\chi^1, \dots, \chi^{4N_r}) \sim (\chi^1 + 2\pi, \dots, \chi^{4N_r}) \sim (\chi^1, \dots, \chi^{4N_r} + 2\pi)$ [41, 90]. Following [42], the prior distribution is chosen as a uniform distribution over \mathbb{T}^{4N_r} . We regard the torsion angles as mutually independent and use interpolation paths as: $\chi_t = (1-t)\chi_0 + t(\chi'_1 - \chi_0)$ where $\chi'_1 = (\chi_1 - \chi_0 + \pi) \bmod (2\pi) - \pi + \chi_0$. The loss for the torsion angles is defined as:

$$\mathcal{L}_{tor}(\theta) = \mathbb{E}_{t, p_1(\chi_1), p_0(\chi_0), p_t(\chi_t | \chi_0, \chi_1)} \frac{1}{N_r} \sum_{i=1}^{N_r} \left\| v_{\theta}^{(i)}(\chi_t^{(i)}, t) - \chi_1'^{(i)} + \chi_0^{(i)} \right\|_2^2. \quad (3)$$

■ PocketFlow on Residue Types and Interaction Types

- The **prior distribution** is set as the uniform distribution and the **conditional flow** is defined as the Euclidean interpolation between initial data and sample data.

$$\mathcal{L}_{res} = \mathbb{E}_{t \sim \mathcal{U}(0,1), p_1(c_1), p_0(c_0), p_t(c | c_0, c_1)} \sum_{i=1}^{N_r} \text{CE} \left(c_t^{(i)} + (1-t)v_{\theta}^{(i)}(c_t^{(i)}, t), c_1^{(i)} \right),$$

$$\mathcal{L}_{inter} = \mathbb{E}_{t \sim \mathcal{U}(0,1), p_1(I_1), p_0(I_0), p_t(I | I_0, I_1)} \sum_{i=1}^{N_r} \text{CE} \left(I_t^{(i)} + (1-t)v_{\theta}^{(i)}(I_t^{(i)}, t), I_1^{(i)} \right).$$

Method

Model Architecture

- PocketFlow adopt the neural network architecture from the **FrameDiff**, which incorporates Invariant Point Attention from AF2 to encode spatial features combined with transformer layers to encode sequence-level features.

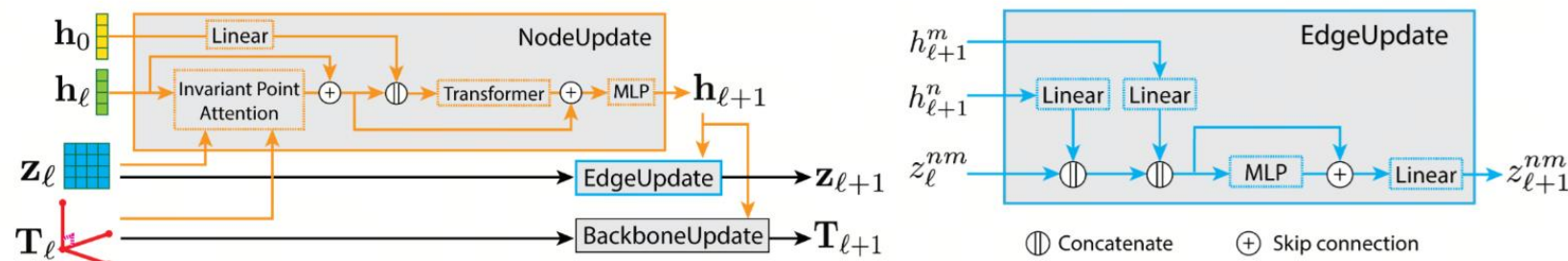


Figure 2: Framework of FrameDiff

Residue/Interaction Type and Torsion angle Prediction. We predict the residue/interaction types and sidechain torsion angles based on node embeddings.

$$\mathcal{R} = \{c^{(i)}, x^{(i)}, O^{(i)}, \chi^{(i)}, I^{(i)}\}_{i=1}^{N_r}$$

$$h_c = \text{MLP}(h^L), \quad h_I = \text{MLP}(h^L), \quad h_\chi = \text{MLP}(h^L), \quad (31)$$

$$c = \text{softmax}(\text{Linear}(h_c + h^L)), I = \text{softmax}(\text{Linear}(h_\psi + h^L)), \quad (32)$$

$$\chi = \text{Linear}(h_\chi + h^L) \bmod 2\pi \quad (33)$$



Method

unconditional
vector field guidance term

■ **Sampling** $\nabla_{\mathcal{C}_t} \log p(\mathcal{C}_t|y) = \nabla_{\mathcal{C}_t} \log p(\mathcal{C}_t) + \nabla_{\mathcal{C}_t} \log p(y|\mathcal{C}_t),$

➤ Generally, we use **classifier-guided sampling** and consider overall **binding affinity guidance** and **interaction geometry guidance**.

- For binding affinity guidance, we just train a separate lightweight affinity predictor for guidance.
- For interaction geometry guidance, we make the local geometries satisfy a series of distance/angle constraints.

Sampling. With the initialized data, the sampling process is the integration of the ODE $\frac{d\mathcal{C}_t}{dt} = v_\theta(\mathcal{C}_t, t)$ from $t = 0$ to $t = 1$ with an Euler solver [14]. γ, ξ_1, ξ_2 , and ξ_3 are set as 1 in the default setting. To apply the guidance, we use \tilde{v}_θ which is v_θ plus guidance terms (Equ. 7, 9, and 10):

$$\chi_{t+\Delta t}^{(i)} = \text{reg} \left(\chi_t^{(i)} + \tilde{v}_\theta(\chi_t^{(i)}, t) \Delta t \right); \quad (11)$$

$$\mathbf{x}_{t+\Delta t}^{(i)} = \mathbf{x}_t^{(i)} + \tilde{v}_\theta(\mathbf{x}_t^{(i)}, t) \Delta t; \quad \mathbf{O}_{t+\Delta t}^{(i)} = \mathbf{O}_t^{(i)} \exp \left(\tilde{v}_\theta(\mathbf{O}_t^{(i)}, t) \Delta t \right); \quad (12)$$

$$\mathbf{c}_{t+\Delta t}^{(i)} = \text{norm} \left(\mathbf{c}_t^{(i)} + \tilde{v}_\theta(\mathbf{c}_t^{(i)}, t) \Delta t \right); \quad \mathbf{I}_{t+\Delta t}^{(i)} = \text{norm} \left(\mathbf{I}_t^{(i)} + \tilde{v}_\theta(\mathbf{I}_t^{(i)}, t) \Delta t \right); \quad (13)$$

➤ Introduction

➤ Method

➤ **Experiment**

➤ Conclusion



Experiment

■ Datasets

- Following previous works, we consider two widely used **protein-small molecule binding datasets** for experimental evaluations:
 - **CrossDocked dataset:** This dataset is generated through crossdocking and is split with mmseqs2 at 30% sequence identity, leading to train/val/test set of 100k/100/100 complexes.
 - **Binding MOAD dataset:** This dataset contains experimentally determined protein-small molecule complexes and is split based on the proteins' enzyme commission number, leading to train/val/test set of 40k/100/100.
- Besides, to test the generalizability of PocketFlow to **other ligand modalities**, we further consider **PPDBench**, which contains 133 non-redundant complexes of **protein-peptides** and **PDBBind RNA**, which contains 56 **protein-RNA** pairs by filtering the PDBBind nucleic acid subset.

Experiment

■ Performance Metrics

- **Amino Acid Recovery (AAR):** the overlapping ratio between the predicted and ground truth residue types.
- **scRMSD:** the self-consistency Root Mean Squared Deviation between the generated and the predicted pocket's backbone atoms to reflect structural validity.
- **Binding affinity:** we choose different binding affinity metrics for different ligands.

predicted structures. To measure the binding affinity for protein-small molecule pairs, we calculate **Vina Score** with AutoDock Vina [78] following [64, 92]. For protein-peptide and protein-RNA pairs, we calculate **Rosetta $\Delta\Delta G$** [5] and **Rosetta-Vienna RNP- $\Delta\Delta G$** [44] respectively that measure the binding affinity change. The unit is kcal/mol and a lower Vina score/ $\Delta\Delta G$ indicates higher affinity.



Experiment

■ Small-molecule-binding Pocket Design

Table 1: Evaluation of different models on **small-molecule-binding** protein pocket design. We report the average and standard deviation values across three independent runs. We highlight the best two results with **bold text** and underlined text, respectively.

Model	CrossDocked			Binding MOAD		
	AAR (\uparrow)	scRMSD (\downarrow)	Vina (\downarrow)	AAR (\uparrow)	scRMSD (\downarrow)	Vina (\downarrow)
Test set	-	0.65	-7.016	-	0.67	-8.076
DEPACT	31.52 \pm 3.26%	0.73 \pm 0.06	-6.632 \pm 0.18	35.30 \pm 2.19%	0.77 \pm 0.08	-7.571 \pm 0.15
dyMEAN	38.71 \pm 2.16%	0.79 \pm 0.09	-6.855 \pm 0.06	41.22 \pm 1.40%	0.80 \pm 0.12	-7.675 \pm 0.09
FAIR	40.16 \pm 1.17%	0.75 \pm 0.03	-7.015 \pm 0.12	43.68 \pm 0.92%	0.72 \pm 0.04	-7.930 \pm 0.15
RFDiffusionAA	50.85 \pm 1.85%	0.68 \pm 0.07	-7.012 \pm 0.09	49.09 \pm 2.49%	0.70 \pm 0.04	-8.020 \pm 0.11
PocketFlow	52.19\pm1.34%	0.67 \pm 0.04	-8.236\pm0.16	54.30\pm1.70%	<u>0.68\pm0.03</u>	-9.370\pm0.24
w/o Aff Guide	50.94 \pm 1.37%	0.65\pm0.04	-7.375 \pm 0.10	51.43 \pm 1.52%	0.64\pm0.04	-8.380 \pm 0.19
w/o Geo Guide	49.80 \pm 1.41%	0.68 \pm 0.03	-8.120 \pm 0.14	53.49 \pm 1.53%	0.71 \pm 0.05	-9.197 \pm 0.22
w/o Geo & Aff Guide	48.50 \pm 1.66%	0.71 \pm 0.06	-7.135 \pm 0.13	49.71 \pm 1.68%	0.69 \pm 0.03	-8.241 \pm 0.18
w/o Inter Learning	50.72 \pm 1.20%	<u>0.66\pm0.03</u>	-7.968 \pm 0.15	52.25 \pm 1.74%	<u>0.68\pm0.05</u>	-9.031 \pm 0.17

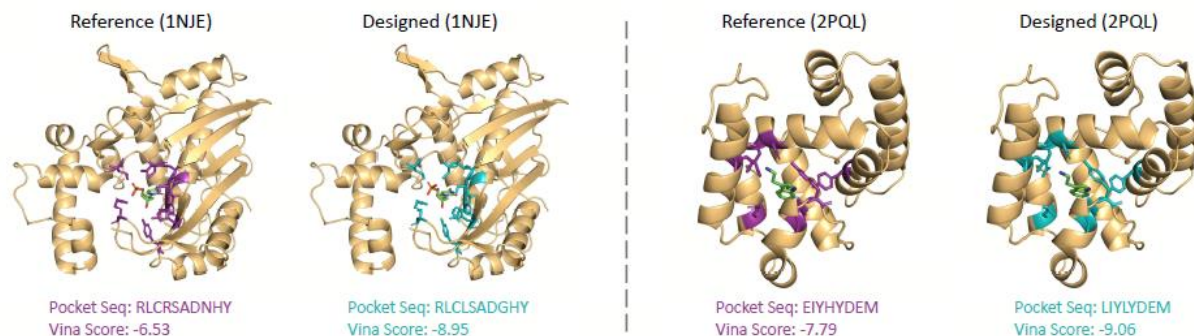
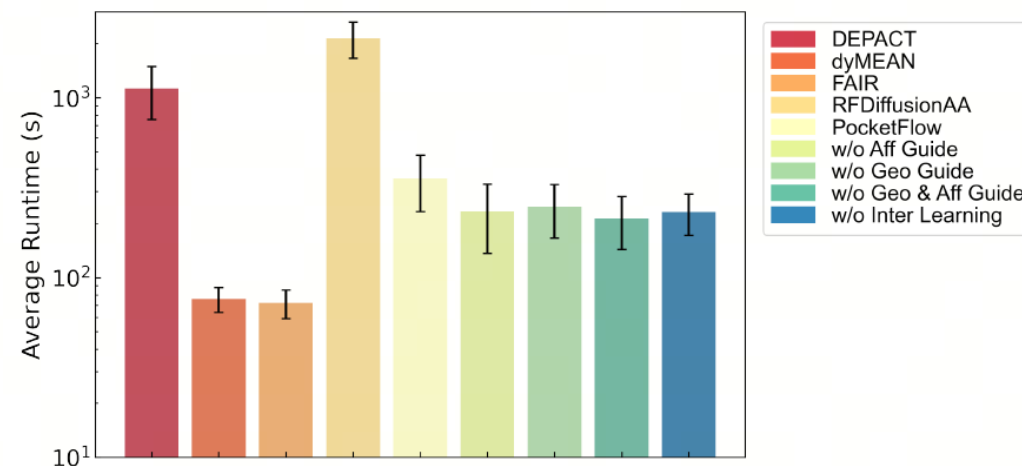


Figure 2: Case studies of small-molecule-binding protein pocket design. We show the reference and designed structures/sequences of two protein pockets from the CrossDocked (PDB ID: 1NJE) and Binding MOAD (PDB ID: 2PQL) datasets respectively.



State-of-the-art performance with acceptable generation efficiency

Experiment

■ Generalization to Other Ligand Domains

- It explores whether the pretrained PocketFlow on the combination of CrossDocked and Binding MOAD can generalize to peptide and RNA binding pocket design.

Table 2: Evaluation of different approaches on the **peptide** and **RNA** datasets. DEPACT is not reported here because it is specially designed for small molecules. dyMEAN, FAIR, and PocketFlow are pretrained on protein-small molecule datasets and we use the checkpoint of RFDiffusionAA [1].

Model	PPDBench			PDBBind RNA		
	AAR (\uparrow)	scRMSD (\downarrow)	$\Delta\Delta G$ (\downarrow)	AAR (\uparrow)	scRMSD (\downarrow)	$\Delta\Delta G$ (\downarrow)
Test set	-	0.64	-	-	0.59	-
dyMEAN	26.29 \pm 1.05%	0.71 \pm 0.05	-0.23 \pm 0.04	25.90 \pm 1.22%	0.71 \pm 0.04	-0.18 \pm 0.03
FAIR	32.53 \pm 0.89%	0.86 \pm 0.04	0.05 \pm 0.07	24.90 \pm 0.92%	0.80 \pm 0.05	0.13 \pm 0.05
RFDiffusionAA	46.85 \pm 1.45%	0.65\pm0.06	-0.62 \pm 0.05	44.69\pm1.90%	0.65\pm0.03	-0.45 \pm 0.07
PocketFlow	48.19\pm1.34%	0.67 \pm 0.04	-1.06\pm0.04	44.34 \pm 1.16%	0.69 \pm 0.01	-0.78\pm0.07
w/o Aff Guide	47.78 \pm 1.18%	0.70 \pm 0.02	-0.47 \pm 0.10	42.15 \pm 1.56%	0.68 \pm 0.04	-0.35 \pm 0.11
w/o Geo Guide	47.30 \pm 1.94%	0.72 \pm 0.05	-0.96 \pm 0.08	41.73 \pm 2.34%	0.77 \pm 0.09	-0.65 \pm 0.15
w/o Geo & Aff Guide	44.63 \pm 1.79%	0.78 \pm 0.05	-0.31 \pm 0.05	39.70 \pm 1.24%	0.78 \pm 0.06	-0.26 \pm 0.08
w/o Inter Learning	36.41 \pm 1.38%	0.74 \pm 0.06	-0.34 \pm 0.05	36.27 \pm 1.47%	0.82 \pm 0.13	-0.23 \pm 0.06

Comparable performance
to the state-of-the-art baseline

➤ Introduction

➤ Method

➤ Experiment

➤ **Conclusion**



Conclusion

- In this paper, we proposed PocketFlow, a **protein-ligand interaction prior-informed flow matching model** for **protein pocket generation**.
- We define **multimodal flow matching** for **protein backbone frames, sidechain torsion angles**, and **residue/interaction types** to appropriately represent the protein-ligand complex.
- The **binding affinity and interaction geometry guidance** effectively improve the validity and affinity of the generated pockets.
- Moreover, PocketFlow offers a **unified framework** covering small-molecule, nucleic acids, and peptides-binding protein pocket generation.



中國人民大學
RENMIN UNIVERSITY OF CHINA



高瓴人工智能学院
Gaoling School of Artificial Intelligence

Thank You for listening!

Fanmeng Wang

2024-10-15