

# Piloting Structure-Based Drug Design via Modality-Specific Optimal Schedule

Keyue Qiu<sup>\*12</sup> Yuxuan Song<sup>\*12</sup> Zhehuan Fan<sup>3</sup> Peidong Liu<sup>4</sup>  
Zhe Zhang<sup>2</sup> Mingyue Zheng<sup>3</sup> Hao Zhou<sup>1</sup> Wei-Ying Ma<sup>1</sup>

## Abstract

Structure-Based Drug Design (SBDD) is crucial for identifying bioactive molecules. Recent deep generative models are faced with challenges in geometric structure modeling. A major bottleneck lies in the twisted probability path of multi-modalities—continuous 3D positions and discrete 2D topologies—which jointly determine molecular geometries. By establishing the fact that noise schedules decide the Variational Lower Bound (VLB) for the twisted probability path, we propose VLB-Optimal Scheduling (VOS) strategy in this under-explored area, which optimizes VLB as a path integral for SBDD. Our model effectively enhances molecular geometries and interaction modeling, achieving a state-of-the-art PoseBusters passing rate of 95.9% on CrossDock, more than 10% improvement upon strong baselines, while maintaining high affinities and robust intramolecular validity evaluated on a held-out test set. Code is available at <https://github.com/AlgoMole/MolCRAFT>.

## 1. Introduction

Structure-Based Drug Design (SBDD) plays a pivotal role in the discovery of bioactive molecules, leveraging the knowledge of protein-ligand interactions to identify potential therapeutic compounds. At the core of SBDD is the accurate modeling of 3D protein-ligand geometries, as only when bioactive compounds can bind effectively to their target receptors can they elicit their therapeutic effects (Isert et al., 2023). Despite its importance, achieving high-fidelity interaction modeling remains a significant challenge, primarily due to the complexity of the underlying binding dynamics.

<sup>\*</sup>Equal contribution <sup>1</sup>Institute for AI Industry Research (AIR), Tsinghua University <sup>2</sup>Department of Computer Science and Technology, Tsinghua University <sup>3</sup>Shanghai Institute of Materia Medica, Chinese Academy of Sciences <sup>4</sup>Sichuan University. Correspondence to: Hao Zhou <zhouhao@air.tsinghua.edu.cn>.

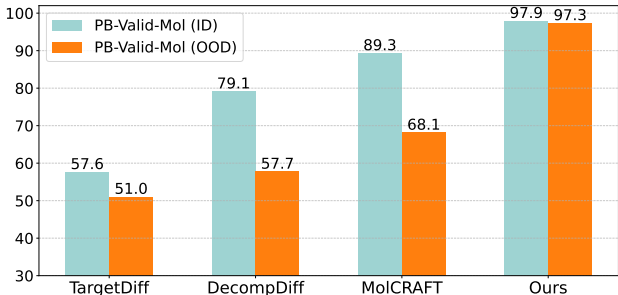


Figure 1. PoseBusters passing rates (%) for non-autoregressive models, where ours maintains the best ID and OOD performance. ID: in-distributional CrossDock test; OOD: out-of-distributional PoseBusters test. PB-Valid-Mol: intramolecular validity. Detailed results can be found in Appendix D.

Recent advances in geometric deep generative models have centered on non-autoregressive methods such as diffusion (Guan et al., 2023) and Bayesian Flow Network (BFN) (Qu et al., 2024), showing promise by capturing the structures at the global level. However, when evaluated with a rigorous out-of-distribution (OOD) test on PoseBusters benchmark (Buttenschoen et al., 2024), Fig. 1 shows a notable performance drop in intramolecular validity, suggesting that current global generation strategies may inadequately capture the fine-grained molecular geometries.

One key factor in geometric structure modeling is the intertwined probability path of different modalities: continuous atom positions and discrete molecular topologies. While these modalities jointly determine molecular geometries and protein-ligand interaction types, there lacks a systematic understanding in designing the twisted probability path for SBDD. Prior works (Peng et al., 2023; Vignac et al., 2023; Guan et al., 2023) adopt sophisticated noise schedules adapted for different modalities, and propose to generate 3D positions first. In our preliminary studies (Sec. 2), we identify the potential problem of prioritizing the 3D modality, where the model cannot effectively utilize 2D topology information in the generation, suggesting that the current noise schedule is suboptimal. Despite its importance, the design of optimal noise schedules for twisted probability path remains a largely underexplored area.

To address the gap left by previous methods, we aim at a systematic solution for the optimal schedule for the twisted

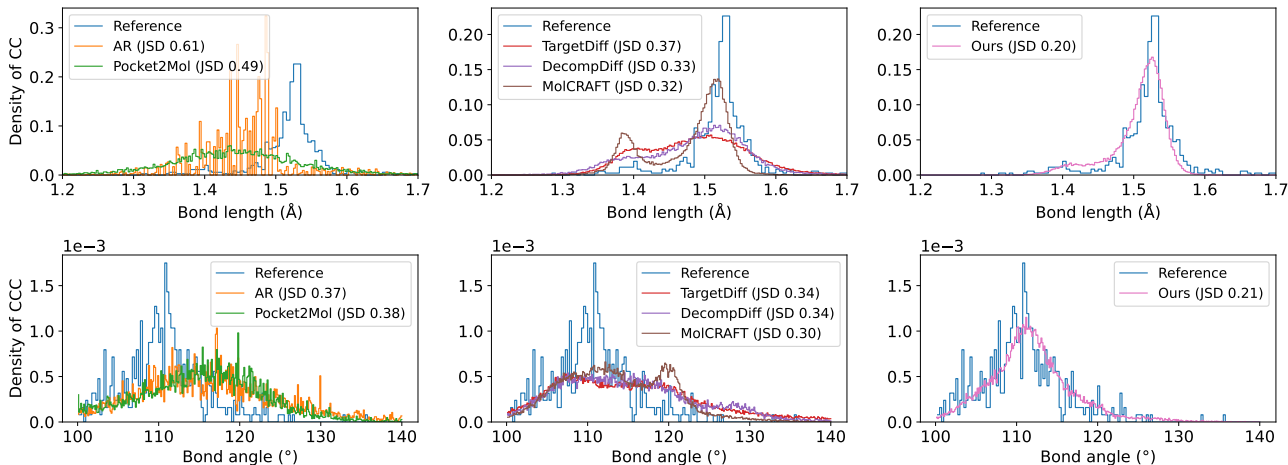


Figure 2. Visualization of the bond length and bond angle distribution for the most frequent bond types in CrossDock test. More types are shown in Appendix D.2, where our MolPilot consistently captures the molecular geometries most accurately.

probability path of different modalities, and emphasize the need for principled metrics relevant to generative modeling in order to evaluate schedule quality. EquiFM (Song et al., 2024b) introduces improved hybrid probability paths aligned with information-theoretic heuristics but lacks a theoretical foundation for optimality. To rigorously define optimality, we analyze the Variational Lower Bound (VLB) under varying noise schedules. Crucially, we prove that in the multi-modality generation, the VLB becomes a path-dependent integral; its value depends on the entire noise schedule, not just endpoints (Kingma et al., 2021).

This motivates our **VLB-Optimal Scheduling** (VOS), eliminating heuristic noise schedule design and directly linking it to the theoretical guarantees of the VLB landscape. Specifically, we combat the *combinatorial complexity in schedule design space* when navigating and optimizing for VLB. To make it tractable to exhaust the possibilities, we demonstrate that the design space of multi-modality noise schedules can be reduced to a two-dimensional plane. Building on this insight, we develop a generalized objective, showing that decoupling timesteps during training allows a single model to implicitly cover this space, generalizing beyond fixed schedules. This framework enables efficient interpolation and extrapolation of schedules at inference time, bypassing costly retraining for new design constraints. By advancing the underexplored area of modality-specific optimal scheduling, we address key shortcomings in SBDD models such as strained conformations and suboptimal interactions, substantially improving the geometric structure modeling.

Our contributions can be summarized as follows:

- We introduce VLB-Optimal Scheduling (VOS), a novel method for systematic noise schedule design in SBDD, achieving fine-grained control over multi-modality interdependence and improved interaction modeling.

- We establish the theoretical link between noise schedules and VLB in multi-modality probabilistic modeling. Unlike prior heuristic approaches, our principled method reveals path-dependent VLB dynamics and provides insights for modality-specific scheduling.
- Integrated with advanced frameworks, our proposed MolPilot achieves SOTA in de novo design with a remarkable PoseBusters passing rate of 95.9% on CrossDock, and competitive in local docking with 44.0% RMSD < 2Å on PoseBusters.

## 2. Issues with Current Probability Path

In this section, we introduce the formulation of Structure-Based Drug Design (SBDD) and demonstrate that the default noise schedule in current generative models leads to a 3D-dominant probability path. Through preliminary experiments, we demonstrate that the model trained along this path cannot adequately capture the interdependence between 2D and 3D modalities, which motivates the need for an optimal probability path facilitated by corresponding noise schedule.

### 2.1. Problem Formulation

Structure-based Drug Design (SBDD) involves modeling the conditional probability  $P(\mathbf{x}_M|\mathbf{x}_P)$ , where  $\mathbf{x}_M = (\mathbf{r}_M, \mathbf{h}_M, \mathbf{A}_M)$  represents the  $N$ -atom molecular geometry, and  $\mathbf{x}_P = (\mathbf{r}_P, \mathbf{h}_P)$  represents the protein target. Here,  $\mathbf{r} \in \mathbb{R}^{N \times 3}$  denotes continuous atom positions,  $\mathbf{h} \in \mathbb{R}^{N \times K_h}$  encodes discrete atom types, and  $\mathbf{A} \in \mathbb{R}^{N \times (N-1) \times K_A}$  encodes discrete bond types, with  $K_h$  and  $K_A$  denoting the number of atom types and bond types, respectively. In molecular generation,  $\mathbf{r}$  from the continuous modality is usually described as the 3D geometry, and  $(\mathbf{h}, \mathbf{A})$  from the discrete modality as the 2D topology of the molecular graph.

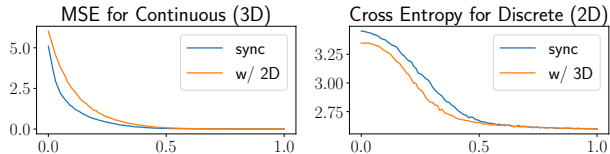


Figure 3. Validation loss curves on default schedule w.r.t. time in generation. *Sync*: Modalities at the same timestep. *w/ 2D*: Discrete modality at  $t = 1$ . *w/ 3D*: Continuous modality at  $t = 1$ .

## 2.2. 3D-driven Probability Path

Fig. 2 suggests that models have yet to accurately capture the molecular geometry, potentially resulting in 3D molecular structures that are incompatible with their 2D molecular graphs (Appendix D). To investigate this, we visualize the modality-specific validation losses for a vanilla BFN (Qu et al., 2024) trained with default noise schedule.

As shown in Fig. 3, the continuous modality loss decreases before the discrete modality, verifying that the twisted probability path is driven by 3D modality. However, the problem lies in that the model leverages cleaner 3D input to denoise the 2D topology effectively (*Right*), but fails the other way around (*Left*), as it performs worse when exposed to less noisy 2D input, suggesting that it is unable to benefit from cleaner 2D information. This highlights that the current noise schedule favors a twisted probability path driven by 3D modality, inducing a generative process functioning with lower-noise 3D conformation and higher-noise 2D topology.

Intuitively, a well-balanced probability path should allow the models to leverage cleaner input from either modality to inform the generative process effectively, thereby accurately capturing the molecular geometries required for effective drug design. The problem of current twisted probability path is strengthened by the observation of inaccurate geometric structure modeling, where SBDD baselines exhibit abnormal bond lengths and angles (Fig. 14-16), suggesting inconsistency between modalities. We hypothesize that the distorted molecular geometry originates from the gap between modalities in the twisted probability path, which inspires our search for an optimal scheduling.

We desire the optimal probability path defined by the optimal noise schedule to benefit from mutually informed modalities, with the potential to improve the generation of chemically valid and spatially accurate molecules, ultimately enhancing protein-ligand interaction modeling.

## 3. Background

Our work is built within the framework of Bayesian Flow Network (BFN) (Graves et al., 2023), the SOTA model for 3D molecular generation (Song et al., 2024a; Qu et al., 2024)

that shows superiority over diffusion counterparts. In this section, we introduce the noise schedule and define the corresponding noising process, paving the way for Variational Lower Bound (VLB) analysis w.r.t. noise schedules.

Similar to diffusion, BFN involves a noising process  $q$  for data  $\mathbf{x}$  to obtain a temporal sequence of observed latents  $\mathbf{y}_{1:n} := \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ , and optimizes the VLB:

$$\log p(\mathbf{x}) \geq \text{VLB} = \mathbb{E}_{q(\mathbf{y}_{1:n}|\mathbf{x})} [\log p(\mathbf{x} | \mathbf{y}_{1:n}) - D_{\text{KL}}(q(\mathbf{y}_{1:n}|\mathbf{x}) \| p(\mathbf{y}_{1:n}))] \quad (1)$$

where the noising process  $q$  is defined by the variational distribution  $q(\mathbf{y}_{1:n} | \mathbf{x}) = \prod_{i=1}^n q(\mathbf{y}_i | \mathbf{x})$  that depends on the noise schedule over timesteps:

$$q(\mathbf{y}_i | \mathbf{x}) = \begin{cases} \mathcal{N}(\mathbf{x}, \alpha_i^{-1} \mathbf{I}), & \text{continuous data} \\ \mathcal{N}(\alpha_i(K\mathbf{e}_x - \mathbf{1}), \alpha_i K \mathbf{I}), & \text{discrete one-hot} \end{cases} \quad (2)$$

where  $K$  is the number of classes,  $\mathbf{e}_x$  is Kronecker function, i.e., the projection from a class index  $\mathbf{x} = j$  to a one-hot vector  $\in \mathbb{R}^K$  with the  $j$ -th value equal to 1. For time step  $i$ , the discretized  $\alpha_i$  is obtained from modality-specific noise schedules, defined as monotonically increasing differentiable functions over  $t \in [0, 1]$ :

$$\beta_c(t) = \sigma_1^{-2t} - 1, \quad \beta_d(t) = \beta_1 t^2, \quad (3)$$

where  $\sigma_1 \in \mathbb{R}^+$  controls the noise for continuous atomic positions, and  $\beta_1 \in \mathbb{R}^+$  is the hyperparameter for discrete topology. The noise level  $\alpha(t) = \beta'(t) := d\beta(t)/dt$ .

BFN differs from diffusions in the generative process, where instead of noisy latents  $\mathbf{y}$ , the network is informed by a lower-variance posterior  $\theta$  given the observed latents, from the Bayesian flow distribution of Gaussian or Dirac  $\delta$ :

$$p_F(\theta | \mathbf{x}; t) = \begin{cases} \mathcal{N}(\gamma \mathbf{x}, \gamma(1 - \gamma) \mathbf{I}), & \text{continuous} \\ \delta(\theta - \text{softmax}(\mathbf{y})), & \text{discrete} \end{cases} \quad (4)$$

where  $\gamma := \frac{\beta_c(t)}{1 + \beta_c(t)}$ ,  $\mathbf{y}$  follows the discrete case in Eq. 2.

The network  $\tilde{\mathbf{x}}_\phi(\theta, t)$ <sup>1</sup> is trained to denoise  $\mathbf{x}$  given  $\theta$  instead of  $\mathbf{y}$ , and thus constructs a receiver distribution  $p_R(\mathbf{y} | \theta; t) = q(\mathbf{y} | \tilde{\mathbf{x}}_\phi(\theta, t))$ . We show in Appendix B.1 that the single-modality VLB objective simplifies to:

$$\mathcal{L}_{\text{VLB}}(\mathbf{x}) = -\mathbb{E}_{q(\mathbf{y}_{1:n}|\mathbf{x})} \log p(\mathbf{x} | \mathbf{y}_{1:n}) + \sum_{i=1}^n D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \| p(\mathbf{y}_i | \mathbf{y}_{1:i-1})) \quad (5)$$

and can be optimized by the continuous-time loss:

$$\begin{aligned} \mathcal{L}^\infty(\mathbf{x}) &= \mathbb{E}_{t \sim U(0,1), p_F(\theta|\mathbf{x};t)} D_{\text{KL}}(q(\mathbf{y} | \mathbf{x}; t) \| p_R(\mathbf{y} | \theta; t)) \\ &= \frac{1}{2} \mathbb{E}_{t \sim U(0,1), p_F(\theta|\mathbf{x};t)} \beta'(t) \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, t)\|^2 \end{aligned} \quad (6)$$

<sup>1</sup>The actual noise prediction model is specified as  $\tilde{\mathbf{x}}_\phi(\theta, t, \mathbf{x}_P)$  with the protein structure input. We omit  $\mathbf{x}_P$  for simplicity.

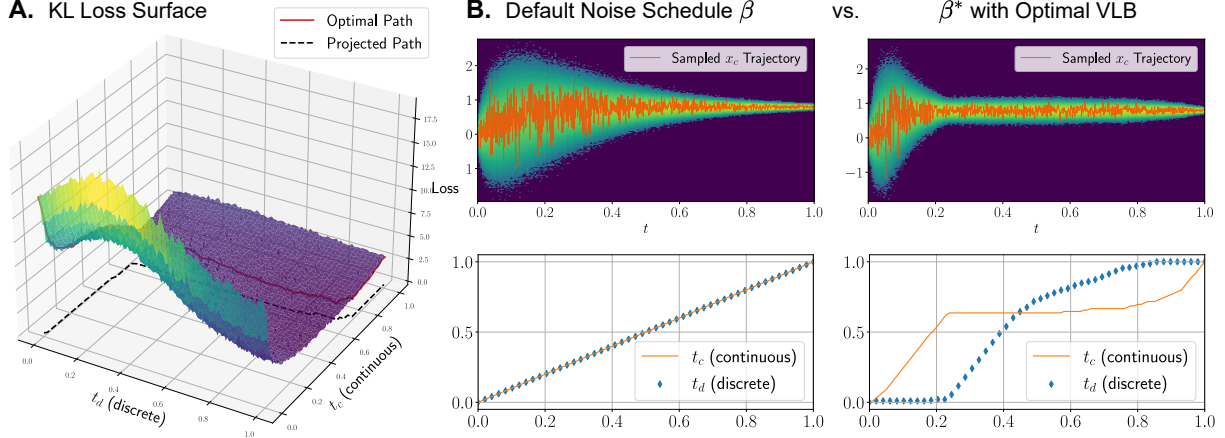


Figure 4. Our proposed VLB-Optimal Scheduling (VOS) that works by estimating the loss landscape and deriving an optimal noise schedule. **A.** Visualization of the loss surface over the function space. **B. Upper:** Visualization of the probability path of continuous data  $\mathbf{x}_c$ . **Lower:** The equivalent time-rescaling functions  $t_c \equiv f(t)$ ,  $t_d \equiv g(t)$ .

where we redefine  $\mathbf{x} := \sqrt{K}\mathbf{e}_\mathbf{x}$  for the discrete data for simplicity. Based on the monotonic and thus invertible function  $u \equiv \beta(t)$ , we perform a change of variables similarly to Kingma et al. (2021) as  $t = \beta^{-1}(u)$ , and rewrite the network as  $\tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)$ . Therefore, the loss in Eq. 6 is equivalent to the continuous-time loss expressed by noise schedule  $\beta(t) : [0, 1] \rightarrow \mathbb{R}^+$ :

$$\mathcal{L}^\infty(\mathbf{x}) = \frac{1}{2} \int_{\beta(0)}^{\beta(1)} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};u)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)\|^2 du. \quad (7)$$

In the single-modality case, this states that the VLB is invariant to the shape of noise schedule functions except for the endpoints (Kingma et al., 2021). However, we will see in the following section that this invariance no longer holds for multi-modality cases, where the VLB becomes:

$$\begin{aligned} \mathcal{L}_{\text{VLB}}(\mathbf{x}) = & \frac{1}{2} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};t)} \int_0^1 [\beta'_c(t) \|\mathbf{x}_c - \tilde{\mathbf{x}}_{\phi,c}(\boldsymbol{\theta}, \beta(t))\|^2 \\ & + \beta'_d(t) \|\mathbf{x}_d - \tilde{\mathbf{x}}_{\phi,d}(\boldsymbol{\theta}, \beta(t))\|^2] dt, \end{aligned} \quad (8)$$

which is simply summed along each modality due to the factorized nature of  $q$ , yet this results in path-dependent VLB which is closely related to the design of joint noise schedule  $\beta(t) : [0, 1] \rightarrow (\mathbb{R}^+)^2$ .

## 4. Methodology

We propose VLB-Optimal Scheduling (VOS), a principled methodology for optimal probability path in SBDD, by analyzing and identifying the optimal noise schedule for both discrete 2D topologies and continuous 3D atomic positions.

**Overview** (1) In Sec. 4.1, we theoretically analyze the path-dependent VLB for a given twisted probability path induced by specific noise schedule in that function space.

(2) In Sec. 4.2, we formalize the entire function space of possible noise schedules, and propose to navigate the space by time rescaling. (3) In Sec. 4.3, we propose a path-invariant generalized loss objective that enables efficient estimation of the VLB landscape over the entire function space of noise schedules. Then, we describe the method to search for the optimal schedule on the landscape that maximizes the VLB.

### 4.1. Path-dependent VLB for Joint Noise Schedule

In this section, we establish a key result for the twisted probability path, where the joint noise schedule  $\beta(t) : [0, 1] \rightarrow (\mathbb{R}^+)^2$  induces path-dependent VLB.

A key foundation in single modality is that the continuous-time loss, i.e., VLB, remains invariant to the shape of noise schedule function  $\beta(t)$ , except for the endpoints  $\beta(0), \beta(1)$  (Kingma et al., 2021), allowing the design of different schedules for efficient training. As a natural extension, we generalize this invariance to multi-modalities given  $\beta_c \neq \beta_d$ :

$$\mathcal{L}^\infty(\mathbf{x}) = \frac{1}{2} \int_{\beta_c(0)}^{\beta_c(1)} \int_{\beta_d(0)}^{\beta_d(1)} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};\boldsymbol{\beta})} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2 d\boldsymbol{\beta}. \quad (9)$$

This equation shows that the generalized loss is invariant to the decoupled schedules as a surface integral over the plane  $[\beta_c(0), \beta_c(1)] \times [\beta_d(0), \beta_d(1)]$ .

However, this generalized loss  $\hat{\mathcal{L}}$  no longer corresponds to the objective of generative modeling, where the VLB should be a path integral along a line on such plane:

$$\mathcal{L}^\infty(\mathbf{x}) = \frac{1}{2} \int_{\beta_c, \beta_d} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};\boldsymbol{\beta})} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2 d\boldsymbol{\beta}. \quad (10)$$

This curve is equivalent to a specific coupled noise schedule  $\boldsymbol{\beta} \in \mathcal{Z}$ . Therefore, we conclude that the VLB varies for different coupling even with the same endpoints  $\boldsymbol{\beta}(0) =$



**Algorithm 1** Deriving Optimal Schedule

**Input:** Multi-modality data  $\mathbf{x}$ , default noise schedule  $\beta = (\tilde{\beta}_c, \tilde{\beta}_d)$ , grid resolution  $M$ , step  $N$ , step scale  $K$ .  
**Output:** Optimal schedule  $\beta^*$ , generative model  $\tilde{\mathbf{x}}_\phi(\theta, t)$ .  
 Train  $\tilde{\mathbf{x}}_\phi(\theta, t)$  to minimize the generalized loss  $\hat{\mathcal{L}}$  (Eq. 14)  
 $\hat{\mathcal{L}}^\infty(\mathbf{x}) = \frac{1}{2} \int_0^1 \int_0^1 \mathbb{E}_{p_F(\theta|\mathbf{x};t)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, t)\|^2 dt_c dt_d$   
 Discretized grid  $g = (t_c, t_d)_{M \times M} \leftarrow [\frac{1}{M}, \dots, 1]^2$   
**foreach**  $(t_c, t_d)$  **in**  $g$  **do**  
     Compute cost  $C(t_c, t_d) \leftarrow \frac{1}{2} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, t)\|^2$   
 Optimal path  $\{t_c^* \equiv f(t), t_d^* \equiv g(t)\}, J^* \leftarrow \text{DP}(C, N, K)$   
 Optimal schedule  $\beta^*(t) \leftarrow (\tilde{\beta}_c(t_c^*), \tilde{\beta}_d(t_d^*))$

$(\beta_c(0), \beta_d(0))$  and  $\beta(1) = (\beta_c(1), \beta_d(1))$ , and it is no longer agnostic to the intermediate trajectory of  $\beta(t)$ , which is essentially different from the generative modeling within single modality.

The fundamental challenge, then, lies in identifying an optimal joint schedule  $\beta^*$  in the design space  $\mathcal{Z}$  with the best possible VLB integrated along that path:

$$\beta^* = \arg \min_{\beta \in \mathcal{Z}} \int_{\beta} \mathbb{E}_{p_F(\theta|\mathbf{x};\beta)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, \beta)\|^2 d\beta. \quad (11)$$

## 4.2. Design Space of Joint Noise Schedules

To formalize the design space of noise schedules, we define the space of monotonically increasing functions  $\mathcal{X} : [0, 1] \rightarrow \mathbb{R}^+$ , thus the coupled function space  $\mathcal{Z}$  is:

$$\{\beta(t) = (\beta_c(t), \beta_d(t)) \mid \forall \beta_c, \beta_d \in \mathcal{X}, \beta'_c > 0, \beta'_d > 0\}, \quad (12)$$

where  $\beta_c(t)$  and  $\beta_d(t)$  are monotonic schedules for the continuous and discrete modalities, respectively.

To explore arbitrary schedule configurations, we introduce time-rescaling functions  $t_c \equiv f(t)$ ,  $t_d \equiv g(t)$ , allowing for arbitrary  $\beta$  expressed in predefined  $\tilde{\beta}_c, \tilde{\beta}_d$  as in Eq. 3 and varying forms of implicit functions  $f$  and  $g$  by:

$$\beta(t) = (\tilde{\beta}_c(f(t)), \tilde{\beta}_d(g(t))) \quad (13)$$

for which we have the following theorem that guarantees this general form of  $\beta$  expressed in terms of  $t_c, t_d$  is sufficient to cover the entire function space  $\mathcal{Z}$ .

**Theorem 4.1.** *Suppose we have a monotonic function  $\beta_m(t) : [0, 1] \rightarrow \mathbb{R}^+$ , and let  $\beta_m(t)$  be any such monotonic function. Then there exists a time-rescaling function  $t_m \equiv f(t)$  such that  $\beta_m(t) = \tilde{\beta}_m(t_m)$ . In fact,  $f(t) = \tilde{\beta}_m^{-1}(\beta_m(t))$  has the same monotonicity as  $\beta_m(t), \tilde{\beta}_m(t)$ .*

**Remark 4.2.** It follows from the monotonicity of the schedule functions that we can obtain arbitrary combination of noise levels simply by setting  $t_c, t_d \in [0, 1]$  separately.

## 4.3. Navigating for the VLB-optimal Schedule

The first obstacle in identifying the optimal  $\beta^*$  is to obtain a model  $\tilde{\mathbf{x}}_\phi(\theta, \beta)$  that can be evaluated over all possible  $\beta \in \mathcal{Z}$ . We facilitate the VLB analysis of different joint noise schedules by training  $\tilde{\mathbf{x}}_\phi(\theta, t)$  to minimize the generalized loss  $\hat{\mathcal{L}}$  described in Eq. 9 through the change of variables  $t_c \equiv \beta_c^{-1}(t)$ ,  $t_d \equiv \beta_d^{-1}(t)$ , thereby the objective becomes:

$$\hat{\mathcal{L}}^\infty(\mathbf{x}) = \frac{1}{2} \int_0^1 \int_0^1 \mathbb{E}_{p_F(\theta|\mathbf{x};t)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, t)\|^2 dt_c dt_d, \quad (14)$$

for which the proposition holds (proof in Appendix B.4):

**Proposition 4.3.** *Suppose we have a model  $\tilde{\mathbf{x}}_\phi(\theta, t)$  trained by Eq. 14, and let  $\beta_c(t), \beta_d(t)$  be any monotonically increasing functions in  $\mathcal{X}$ . Then the line integral in Eq. 10 corresponds to the negative VLB for  $\beta(t) = (\beta_c(t), \beta_d(t))$ .*

Then, we present the method to navigate the function space of arbitrary  $\beta$ . Note that from Remark 4.2, we can discretize the function space through discretized combinations of  $t_c, t_d$ . Therefore, we can formulate finding the optimal  $\beta^*$  as a search problem for the solution with the minimal cumulative cost along the discretized  $N$ -step trajectory  $\{t_c, t_d\}$  from  $[0, 0]$  to  $[1, 1]$ .

Depicted in Fig. 4, we estimate the cost matrix across a grid of possible  $t_c$  and  $t_d$  values by evaluating the KL divergence over a batch of samples given the  $\hat{\mathcal{L}}$ -trained model  $\tilde{\mathbf{x}}_\phi(\theta, \beta)$ , and fit a smooth loss surface using B-spline interpolation.

**Definition 4.4 (Cost).** The cost at a given noise level  $t_c, t_d$  as implicit functions over  $t$  is defined as:

$$C(t_c, t_d) = \frac{1}{2} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\theta, \beta)\|^2 \quad (15)$$

We employ dynamic programming to solve the search problem for a minimal cumulative cost  $J$ :

$$J(t_c, t_d) = \min_{(\epsilon_c, \epsilon_d)} (J(t_c - \epsilon_c, t_d - \epsilon_d) + \alpha C(t_c, t_d)). \quad (16)$$

where  $\alpha := \beta'(t)$  is the accuracy level, and the valid ranges of  $(\epsilon_c, \epsilon_d)$  are approximated by the gradient of the smooth surface. It is guaranteed that the method yields a schedule with maximal VLB given the corresponding accuracy level:

**Remark 4.5.**  $J(1, 1)$  corresponds to an unbiased Monte-Carlo estimate of the optimal VLB among all generative models trained by arbitrary  $\beta \in \mathcal{Z}$ , and backtracking the path  $\{t_c = f(t), t_d = g(t)\}$  from  $[1, 1]$  to  $[0, 0]$  yields the optimal schedule  $\beta^*(t) = (\tilde{\beta}_c(t_c), \tilde{\beta}_d(t_d))$ .

Table 4 empirically validates the derived schedule with better VLB than default. The overall procedure is summarized in Algorithm 1, with function DP specified in Algorithm 2.

We visualize the derived optimal schedule in Fig. 4B. Intuitively, this time rescaling corresponds to a two-stage probability path: (1) *Shape-driven sketching*: First, the generation predominantly focuses on generating the continuous atom positions and largely ignores the discrete topology. Then, when  $t \in [0.3, 0.8]$ , the model starts to fit a possible 2D molecular graph into the rough shape of fixed 3D conformation. (2) *Topology-driven docking*: at the last stage  $t > 0.8$ , the generation enters into the docking stage, altering the conformation according to the discrete molecular topologies. By facilitating such a principled probability path, the optimal noise schedule effectively harnesses different modalities, achieving the best VLB and sample quality.

## 5. Experiments

### 5.1. Experimental Setup

We conduct two main experiments within the broader scope of SBDD: (1) de novo design in the in-distributional (ID) and out-of-distributional (OOD) settings, and (2) molecular docking. The same model checkpoint trained with generalized loss is evaluated throughout both main experiments.

**Dataset.** Following SBDD conventions (Luo et al., 2021), we adopt the same split of CrossDock (Francoeur et al., 2020) to train and validate our model, which consists of 100,000 training poses and 100 validation poses. (1) For de novo design, we evaluate on an OOD subset of PoseBusters (Buttenschoen et al., 2024) in addition to the ID CrossDock test set. We cluster all chains via MMseqs2 (Steinegger & Söding, 2017) and filter out test proteins with any chain  $> 30\%$  sequence identity w.r.t. CrossDock training sequences, obtaining 180 test proteins. This avoids possible information leakage from the original data splits (detailed discussion in Appendix C.1), and serves as a more reliable held-out test, where the structures of protein-ligand complexes are experimentally determined. (2) For molecular docking, after removing 10 proteins that cannot be processed for non-standard residues from the PoseBusters V2, we test on the remaining 298 protein-ligand complexes as ground-truth-based evaluation.

**Baselines.** We consider the following SBDD baselines: (1) Autoregressive models including AR (Luo et al., 2021), Pocket2Mol (Peng et al., 2022), (2) Diffusion-based Target-Diff (Guan et al., 2022), DecompDiff (Guan et al., 2023), and (3) BFN-based MolCRAFT (Qu et al., 2024). A detailed list of baselines can be referred to in Appendix C.3.

**Metrics.** We employ the following metrics: (1) For conformation quality, we report PoseBusters passing rate (**PB-Valid**), and **Strain Energy Passed** denotes the ratio of molecules passing PoseBusters internal energy check. (2)

For interaction modeling, we emphasize **RMSD**, the root-mean-squared distance between generated and ground truth pose in docking. We also report the self-consistency RMSD (**scRMSD**) in de novo design calculated between generated and Vina Dock pose as an indicator of binding mode consistency upon redocking. Affinity metrics are calculated by AutoDock Vina (Eberhardt et al., 2021), including binding affinities where **Vina Score** directly scores the generated molecular pose in the pocket, **Vina Min** quickly optimizes the pose in-place and scores the minimized pose, and **Vina Dock** exhaustively searches for optimal pose to obtain the lowest energy. (3) For molecular properties, we report drug-likeness (**QED**) and synthetic accessibility (**SA**), which are desired to fall within reasonable ranges. We additionally report the percentage of fully connected molecules (**Connected**) to show the performance of generation successes, together with the average number of atoms of those successfully generated molecules (**Size**).

### 5.2. De novo Design

We report the main results for ID and OOD design in Table 1 by sampling 100 molecules for each test protein.

**MolPilot achieves the most accurate molecular geometries.** For intramolecular validity, we achieve the highest passing rate of internal energy, demonstrating our conformation stability and showing robustness in both ID and OOD settings. We additionally report the bond length, angle, and torsion angle distributions in Fig. 14, 15, 16, further underscoring its superiority.

**MolPilot generates the best binding poses.** For protein-ligand intermolecular validity, we excel at the highest PB-Valid, ensuring reasonable binding poses. For binding affinities, our method outperforms all other models with an average Vina Score of -6.88 (ID) and -7.45 kcal/mol (OOD), very close to Vina Min and Dock.

**CrossDock benchmark may not be challenging enough.** In Table 1, we observe that our rate of PB-Valid and internal energy passed even outperforms the CrossDock test set, suggesting potential problems with this synthetic dataset (Qu et al., 2024). We dive deeper into the benchmarking statistics, and identify possible data leakage in Appendix C.1.

**Models underperform on the PoseBusters held-out test.** In Table 1, we select competitive baselines on CrossDock, and run the sampling given their released codebases and checkpoints. Sampling efficiency measured in the ratio of successfully generated molecules (**Connected**) shows that autoregressive models severely degrade on these challenging targets. Among non-autoregressive baselines, DecompDiff exhibits the most prominent degeneration in Vina affinities,

Table 1. Performance on CrossDock in an in-distribution (ID) setting and PoseBusters in an out-of-distribution (OOD) setting, where MolPilot shows robust results. ♡: results cited from Qu et al. (2024). †: results calculated by us using the official samples. ◇: results calculated by us using the official code. Top-2 results highlighted in **bold** and underlined, respectively.

Methods	PB-Valid†	Vina Score (↓)		Vina Min (↓)		Vina Dock (↓)		scRMSD	Energy	Connected†	QED	SA	Div	Size
	Avg. (†)	Avg.	Med.	Avg.	Med.	Avg.	Med.	<2 Å (†)	Passed† (†)	Avg. (†)	Avg.	Avg.	Avg.	Avg.
CrossDock (ID)♡	95.0%	-6.36	-6.46	-6.71	-6.49	-7.45	-7.26	34.0%	98.0%	-	0.48	0.73	-	22.8
AR♡	59.0%	-5.75	-5.64	-6.18	-5.88	-6.75	-6.62	36.5%	84.9%	93.5%	0.51	0.63	0.70	17.7
Pocket2Mol♡	72.3%	-5.14	-4.70	-6.42	-5.82	-7.15	-6.79	32.0%	<u>97.3%</u>	96.3%	0.57	0.76	0.69	17.7
FLAG♡	16.0%	45.85	36.52	9.71	-2.43	-4.84	-5.56	0.3%	83.4%	97.1%	0.61	0.63	0.70	16.7
TargetDiff♡	50.5%	-5.47	-6.30	-6.64	-6.83	-7.80	-7.91	37.1%	69.8%	90.4%	0.48	0.58	0.72	24.2
DiffSBDD†	37.6%	-1.44	-4.91	-4.52	-5.84	-7.14	-7.3	18.7%	74.0%	93.2%	0.47	0.58	0.73	24.4
DecompDiff♡	71.7%	-5.19	-5.27	-6.03	-6.00	-7.03	-7.16	24.2%	84.1%	82.9%	0.51	0.66	0.73	21.2
MolCRAFT♡	84.6%	-6.55	-6.95	-7.21	-7.14	-7.67	-7.82	<b>46.8%</b>	91.1%	96.7%	0.50	0.67	0.72	22.7
MolPilot (Ours)	<b>95.9%</b>	<b>-6.88</b>	<b>-7.03</b>	<b>-7.23</b>	<b>-7.27</b>	<b>-7.92</b>	<b>-7.92</b>	<u>41.1%</u>	<b>98.5%</b>	<b>97.9%</b>	0.56	0.74	0.69	22.6
PoseBusters (OOD)†	98.9%	-7.06	-7.05	-7.50	-7.41	-7.98	-7.82	59.4%	100%	-	0.40	0.72	-	25.7
AR◇	54.7%	-5.45	-5.17	-5.67	-5.38	-6.18	-5.94	35.3%	77.7%	39.1%	0.50	0.67	0.76	13.6
Pocket2Mol◇	<u>63.6%</u>	-5.39	-5.03	-6.64	-6.24	-7.40	-7.03	37.8%	97.4%	67.7%	0.57	0.74	0.73	17.4
TargetDiff◇	32.3%	-6.57	-6.78	-7.16	-7.31	<u>-8.18</u>	<b>-8.20</b>	32.3%	65.2%	81.3%	0.41	0.55	0.67	27.0
DecompDiff◇	40.2%	-3.14	-3.02	-4.03	-4.11	-5.06	-5.40	17.0%	80.1%	82.9%	0.47	0.66	0.81	19.3
MolCRAFT◇	57.8%	<u>-7.29</u>	<u>-7.11</u>	<u>-7.44</u>	<u>-7.22</u>	-7.95	-7.73	<u>46.4%</u>	71.6%	97.0%	0.41	0.65	0.65	23.8
MolPilot (Ours)	<b>79.1%</b>	<b>-7.59</b>	<b>-7.54</b>	<b>-7.74</b>	<b>-7.67</b>	<b>-8.20</b>	<b>-7.99</b>	<b>56.1%</b>	<b>98.1%</b>	<b>97.3%</b>	0.49	0.72	0.67	23.5

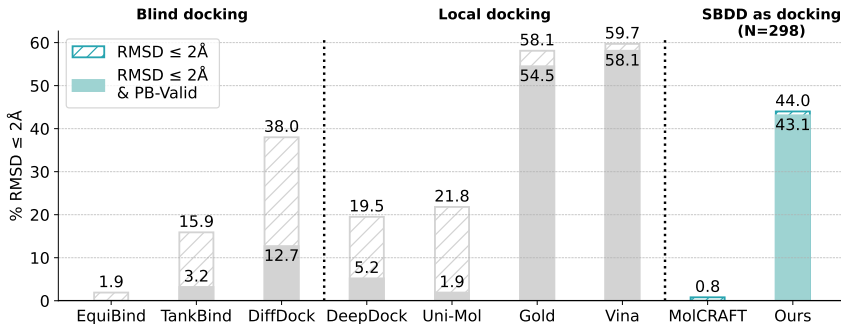


Figure 5. PoseBusters V2 structural accuracy (RMSD) and validity (PB-Valid) for docking methods (in gray, cited from Abramson et al. (2024)) and SBDD methods (in blue, calculated by us).

while TargetDiff displays a notable drop in PB-Valid and ranks the least, and MolCRAFT shows significantly more strained structures reflected by the tail distribution of Strain Energy in Fig. 8. This suggests that PoseBusters can serve as a challenging OOD test for de novo design.

**Our method is robust on the difficult PoseBusters test.** MolPilot reliably generates physically valid conformations within the protein pocket with an impressive 56.1% of cases with scRMSD < 2Å, comparable to the 59.4% obtained by redocking the co-crystallized ligands, which indicates its ability to capture the interaction patterns and maintain binding pose consistency. We additionally calculate the Tanimoto similarity of interaction profiles between generated and ground-truth molecules via ProLIF (Bouysset & Fiorucci, 2021). Table 2 shows that ours best matches the genuine interaction pattern of the co-crystallized structures.

Furthermore, considering that the accuracy of Vina docking tool evaluated on ground-truth complexes in fact places an upper bound, our result of 56.1% suggests that Vina

Table 2. Tanimoto similarity of interaction profiles between PoseBusters reference and generated PB-Valid molecules.

Methods	Sim. (†)
AR	0.221
Pocket2Mol	0.436
TargetDiff	0.458
DecompDiff	0.374
MolCRAFT	0.498
MolPilot (Ours)	<b>0.551</b>

might be reaching its limit in evaluating the pose generation accuracy. This brings us one step further towards ground truth-based evaluation, i.e., molecular docking.

### 5.3. Molecular Docking

We note that the model trained with  $\hat{\mathcal{L}}$  in Eq. 9 is also available for multimarginal generative modeling (Campbell et al., 2024), i.e., it can be used for molecular docking by an induced  $P(\mathbf{r} | \{\mathbf{h}, \mathbf{A}\}, \mathbf{x}_P)$ .

For a deeper understanding of our model’s ability in capturing genuine spatial interaction, we additionally evaluate its docking accuracy on PoseBusters V2 dataset with holo pocket residue structures specified, and report the RMSD results. While we did not specifically optimize our design w.r.t. molecular docking, Fig. 5 suggests that our method is a competitive SBDD model even compared with the local docking methods. This capability in recovering ground truth binding poses serves as an indicator of capturing true interactions, which is essential for designing bioactive molecules

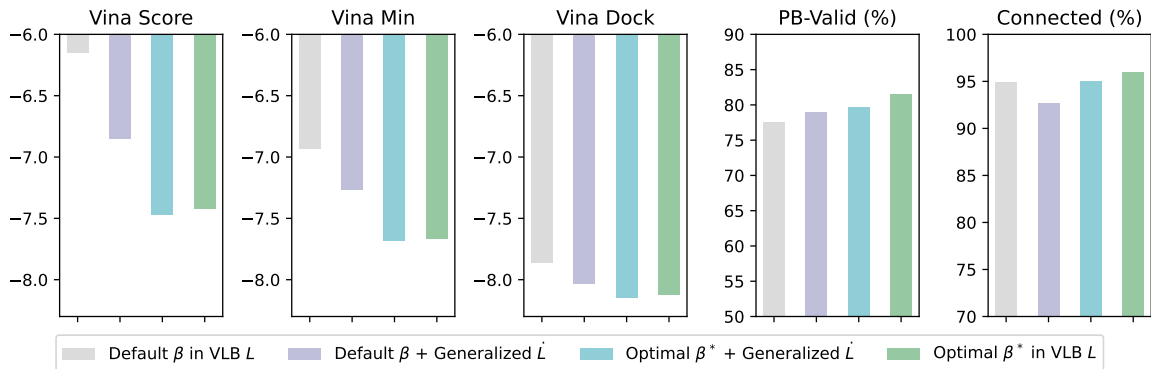


Figure 6. Ablation studies regarding loss objectives and noise schedules. Training with VLB (Eq. 10) as a line integral requires specific schedule  $\beta$  in training and sampling, while with generalized loss (Eq. 9) as a surface integral it only requires  $\beta$  for test time.

Table 3. Results of applying VOS to diffusion model (TargetDiff) on CrossDock. Top-1 results highlighted in **bold**.

Methods	PB-Valid $\dagger$	Vina Score ( $\downarrow$ )		Vina Min ( $\downarrow$ )		Vina Dock ( $\downarrow$ )		scRMSD	Energy	Connected $\dagger$	QED	SA
	Avg. ( $\uparrow$ )	Avg.	Med.	Avg.	Med.	Avg.	Med.	<2 Å ( $\uparrow$ )	Passed $\dagger$ ( $\uparrow$ )	Avg. ( $\uparrow$ )	Avg.	Avg.
TargetDiff	50.5%	-5.47	-6.30	-6.64	-6.83	-7.80	-7.91	37.1%	69.8%	90.4%	0.48	0.58
TargetDiff + $\dot{L}$	53.7%	-6.27	-6.31	-6.82	-6.78	-7.87	-7.90	36.8%	70.3%	89.2%	0.50	0.62
TargetDiff + $\dot{L}$ + $\beta^*$ (VOS)	<b>58.1%</b>	<b>-6.46</b>	<b>-6.53</b>	<b>-7.04</b>	<b>-7.09</b>	<b>-8.04</b>	<b>-8.12</b>	<b>40.2%</b>	<b>73.2%</b>	<b>93.4%</b>	0.49	0.59

with potential biological efficacy. We believe that incorporating molecular docking as a subtask of SBDD also sheds light on the model’s ability of interaction modeling.

#### 5.4. Ablation Studies

We conduct ablation studies by sampling 10 molecules for each of 180 proteins in PoseBusters OOD test, showing the significance of our VOS method in Fig. 6.

**Effect of generalized loss  $\dot{L}$**  The generalized training objective boosts the generative model’s ability in capturing spatial interactions, indicated by the improved Vina affinities especially the Vina Score that directly scores the generated poses. We attribute this gain to the fact that this loss objective forces the model to learn better-balanced probability path for multi-modality molecular data as shown in Fig. 11, especially the previously missing ability to denoise continuous 3D positions with cleaner discrete 2D information on the region of higher-noise 3D and lower-noise 2D data.

**Effect of optimal schedule  $\beta^*$**  When training under the generalized objective  $\dot{L}$  and sampling under different noise schedules  $\beta$ , it can be seen that optimal schedule  $\beta^*$  consistently achieves better performance than default noise schedule. Moreover, the optimal noise schedule can be further adopted in training under the VLB objective, yielding a slightly improved performance. This is due to the fact that VLB remains invariant as a line integral along the path determined by the optimal noise schedule function, yet the

denoising model effectively allocates more capacity in optimizing the VLB, leading to significantly faster convergence in about half the training time with generalized objective  $\dot{L}$ .

#### 5.5. Generality of VOS

To demonstrate VOS’s broader applicability, we have integrated it with the diffusion-based framework TargetDiff. We train for 140k steps following the default configuration with our generalized objective  $\dot{L}$ , and then derive the optimal schedule that proves to resemble the shape in Fig. 4B.

We sample 10 molecules per target on CrossDock and report the results in Table 3, showing that VOS successfully enhances conformation quality for diffusion models as well, with generated poses achieving Vina Scores closely matching Vina Min values, indicating near-optimal realistic poses.

### 6. Related Works

**Structure-based Drug Design (SBDD)** SBDD generative models focus on the joint generation of discrete molecular topology and continuous conformation conditioned on protein-binding pockets. Recent approaches have centered on capturing the critical protein-ligand interactions for the generation of high-affinity molecules. Autoregressive models (Luo et al., 2021; Peng et al., 2022; Liu et al., 2022) generate molecules atom-by-atom while preserving geometric equivariance but are computationally expensive. Fragment-based methods (Powers et al., 2022; ZHANG et al., 2023; Lin et al., 2023) improve efficiency by generating motifs



instead of atoms, though they often require post-processing to mitigate error accumulation. Non-autoregressive models such as diffusion-based approaches (Schneuing et al., 2022; Guan et al., 2022; 2023) and Bayesian Flow Networks (BFNs) (Qu et al., 2024) focus on full-atom generation, enabling scalability and improved controllability. Another line of works incorporates interaction context or property guidance signals such as binding affinity, to enhance the profile of generated ligands (Huang et al., 2024a;b; Qiu et al., 2025). While significant progress has been made, challenges in conformation quality and interaction modeling remain central to advancing structure-based drug design (Harris et al., 2023).

**Molecular Docking** In the context of SBDD, molecular docking concerns predicting the 3D conformation given 2D molecular topology and the protein pocket structures. Traditional search-based local-docking approaches include AutoDock Vina (Eberhardt et al., 2021), Gold (Jones et al., 1997). Deep learning-based methods are divided into blind docking methods with holo protein structures specified, such as DiffDock (Corso et al., 2023), and local docking methods such as DeepDock (Méndez-Lucio et al., 2021) and UniMol (Zhou et al., 2023). These methods additionally rely on the RDKit-initialized molecular conformation as input.

## 7. Conclusion

This work addresses a critical challenge in Structure-Based Drug Design (SBDD) by introducing better-balanced dynamics in the generative process of multi-modalities. Equipped with an optimal noise schedule in terms of Variational Lower Bound (VLB), MolPilot achieves a SOTA PoseBusters passing rate of 95.9% on CrossDock with improved molecular geometry and interaction modeling, offering a promising step forward in the field of drug discovery.

## Impact Statement

This work makes a contribution to the field of computational drug discovery by improving the modeling of protein-ligand interactions, which are crucial for identifying therapeutic compounds. Our approach can be generalized to other tasks involving multimodal generative modeling such as material design. Ultimately, we hope this research could accelerate the discovery of bioactive compounds by enabling more accurate and reliable in silico design of protein-binding ligands, leading to more effective therapeutic interventions.

## Acknowledgments

This work is supported by the Natural Science Foundation of China (Grant No. 62376133) and sponsored by Beijing Nova Program (20240484682) and the Wuxi Re-

search Institute of Applied Technologies, Tsinghua University (20242001120). The authors would like to thank Jingjing Gong, Hanlin Wu and Zhilong Zhang for their helpful comments on this work.

## References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pp. 1–3, 2024.
- Bouysset, C. and Fiorucci, S. Prolif: a library to encode molecular interactions as fingerprints. *Journal of cheminformatics*, 13(1):72, 2021.
- Buttenschoen, M., Morris, G. M., and Deane, C. M. Posebusters: Ai-based docking methods fail to generate physically valid poses or generalise to novel sequences. *Chemical Science*, 15(9):3130–3139, 2024.
- Campbell, A., Yim, J., Barzilay, R., Rainforth, T., and Jaakkola, T. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 5453–5512. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/campbell124a.html>.
- Corso, G., Stärk, H., Jing, B., Barzilay, R., and Jaakkola, T. Diffdock: Diffusion steps, twists, and turns for molecular docking. *International Conference on Learning Representations (ICLR)*, 2023.
- Dieleman, S., Sartran, L., Roshannai, A., Savinov, N., Ganin, Y., Richemond, P. H., Doucet, A., Strudel, R., Dyer, C., Durkan, C., et al. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- Eberhardt, J., Santos-Martins, D., Tillack, A. F., and Forli, S. Autodock vina 1.2. 0: New docking methods, expanded force field, and python bindings. *Journal of chemical information and modeling*, 61(8):3891–3898, 2021.
- Francoeur, P. G., Masuda, T., Sunseri, J., Jia, A., Iovanisci, R. B., Snyder, I., and Koes, D. R. Three-dimensional convolutional neural networks and a cross-docked data set for structure-based drug design. *Journal of chemical information and modeling*, 60(9):4200–4215, 2020.
- Graves, A., Srivastava, R. K., Atkinson, T., and Gomez, F. Bayesian flow networks. *arXiv preprint arXiv:2308.07037*, 2023.

- Guan, J., Qian, W. W., Peng, X., Su, Y., Peng, J., and Ma, J. 3d equivariant diffusion for target-aware molecule generation and affinity prediction. In *The Eleventh International Conference on Learning Representations*, 2022.
- Guan, J., Zhou, X., Yang, Y., Bao, Y., Peng, J., Ma, J., Liu, Q., Wang, L., and Gu, Q. DecompDiff: Diffusion models with decomposed priors for structure-based drug design. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 11827–11846. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/guan23a.html>.
- Harris, C., Didi, K., Jamasb, A. R., Joshi, C. K., Mathis, S. V., Lio, P., and Blundell, T. Benchmarking generated poses: How rational is structure-based drug design with generative models? *arXiv preprint arXiv:2308.07413*, 2023.
- Huang, Z., Yang, L., Zhou, X., Qin, C., Yu, Y., Zheng, X., Zhou, Z., Zhang, W., Wang, Y., and Yang, W. Interaction-based retrieval-augmented diffusion models for protein-specific 3d molecule generation. In *Forty-first International Conference on Machine Learning*, 2024a. URL <https://openreview.net/forum?id=eejhd9FCp3>.
- Huang, Z., Yang, L., Zhou, X., Zhang, Z., Zhang, W., Zheng, X., Chen, J., Wang, Y., Bin, C., and Yang, W. Protein-ligand interaction prior for binding-aware 3d molecule diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Isert, C., Atz, K., and Schneider, G. Structure-based drug design with geometric deep learning. *Current Opinion in Structural Biology*, 79:102548, April 2023. ISSN 0959440X. doi: 10.1016/j.sbi.2023.102548. URL <https://linkinghub.elsevier.com/retrieve/pii/S0959440X23000222>.
- Jones, G., Willett, P., Glen, R. C., Leach, A. R., and Taylor, R. Development and validation of a genetic algorithm for flexible docking. *Journal of molecular biology*, 267(3): 727–748, 1997.
- Kingma, D., Salimans, T., Poole, B., and Ho, J. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Lin, H., Huang, Y., Zhang, O., Liu, Y., Wu, L., Li, S., Chen, Z., and Li, S. Z. Functional-group-based diffusion for pocket-specific molecule generation and elaboration. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 34603–34626. Curran Associates, Inc., 2023. URL <https://arxiv.org/abs/2306.13769>.
- Liu, M., Luo, Y., Uchino, K., Maruhashi, K., and Ji, S. Generating 3d molecules for target protein binding. In *International Conference on Machine Learning*, 2022.
- Luo, S., Guan, J., Ma, J., and Peng, J. A 3D Generative Model for Structure-Based Drug Design. *Advances in Neural Information Processing Systems*, 34: 6229–6239, 2021. URL <http://arxiv.org/abs/2203.10446>.
- Méndez-Lucio, O., Ahmad, M., del Rio-Chanona, E. A., and Wegner, J. K. A geometric deep learning approach to predict binding conformations of bioactive molecules. *Nature Machine Intelligence*, 3(12):1033–1039, 2021.
- Peng, X., Luo, S., Guan, J., Xie, Q., Peng, J., and Ma, J. Pocket2Mol: Efficient molecular sampling based on 3D protein pockets. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 17644–17655. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/peng22b.html>.
- Peng, X., Guan, J., Liu, Q., and Ma, J. Moldiff: Addressing the atom-bond inconsistency problem in 3d molecule diffusion generation. In *International Conference on Machine Learning*, pp. 27611–27629. PMLR, 2023.
- Powers, A. S., Yu, H. H., Suriana, P. A., and Dror, R. O. Fragment-based ligand generation guided by geometric deep learning on protein-ligand structures. In *ICLR2022 Machine Learning for Drug Discovery*, 2022. URL <https://openreview.net/forum?id=192L9cr-8HU>.
- Qiu, K., Song, Y., Yu, J., Ma, H., Cao, Z., Zhang, Z., Wu, Y., Zheng, M., Zhou, H., and Ma, W.-Y. Empower structure-based molecule optimization with gradient guidance. *arXiv e-prints*, pp. arXiv–2411, 2025.
- Qu, Y., Qiu, K., Song, Y., Gong, J., Han, J., Zheng, M., Zhou, H., and Ma, W.-Y. MolCRAFT: Structure-based drug design in continuous parameter space. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=KaAQu5rNU1>.
- Schneuing, A., Du, Y., Harris, C., Jamasb, A., Igashov, I., Du, W., Blundell, T., Lió, P., Gomes, C., Welling, M., Bronstein, M., and Correia, B. Structure-based

- Drug Design with Equivariant Diffusion Models, October 2022. URL <http://arxiv.org/abs/2210.13695>. arXiv:2210.13695 [cs, q-bio].
- Song, Y., Gong, J., Qu, Y., Zhou, H., Zheng, M., Liu, J., and Ma, W.-Y. Unified generative modeling of 3d molecules via bayesian flow networks. *arXiv preprint arXiv:2403.15441*, 2024a.
- Song, Y., Gong, J., Xu, M., Cao, Z., Lan, Y., Ermon, S., Zhou, H., and Ma, W.-Y. Equivariant flow matching with hybrid probability transport for 3d molecule generation. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Steinegger, M. and Söding, J. Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature biotechnology*, 35(11):1026–1028, 2017.
- Vignac, C., Osman, N., Toni, L., and Frossard, P. Midi: Mixed graph and 3d denoising diffusion for molecule generation, 2023. URL <https://arxiv.org/abs/2302.09048>.
- ZHANG, Z., Zheng, S., Min, Y., and Liu, Q. Molecule generation for target protein binding with structural motifs. In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Rq13idF0F73>.
- Zhou, G., Gao, Z., Ding, Q., Zheng, H., Xu, H., Wei, Z., Zhang, L., and Ke, G. Uni-mol: A universal 3d molecular representation learning framework. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=6K2RM6wVqKu>.

## A. Details about the Proposed Method

### A.1. Dynamic Programming for Optimal Schedules

---

**Algorithm 2** Dynamic Programming for Optimal Path
 

---

**Input:** Cost matrix  $C \in \mathbb{R}^{M \times M \times 2}$ , step budget  $L$ , step scale  $K$ , default noise schedule  $\beta(t)$

**Output:** Optimal *path*, minimal cumulative cost  $J^*$

**Initialization:**

Fit a smooth loss surface  $\tilde{C}$  using B-spline interpolation.

$J[:, M, : M, : N + 1] \leftarrow \infty, J[0, 0, 0] \leftarrow 0$

$\text{prev}[:, M, : M, : N + 1, : 2] \leftarrow -1$

$\text{active\_states} \leftarrow \{(0, 0, 0)\}$

$\text{size} \leftarrow \frac{K}{1 + \|\nabla \tilde{C}\|}$ , where  $\|\nabla \tilde{C}\| = \sqrt{(\frac{\partial \tilde{C}}{\partial x})^2 + (\frac{\partial \tilde{C}}{\partial y})^2}$

**Dynamic Programming:**

**for**  $l \leftarrow 0$  **to**  $L - 1$  **do**

$\text{states} \leftarrow \emptyset$

**foreach**  $(x, y, l) \in \text{active\_states}$  **do**

**foreach**  $(\epsilon_x, \epsilon_y) \in \text{valid\_steps}(\text{size}[x, y])$  **do**

$x_t, y_t \leftarrow x + \epsilon_x, y + \epsilon_y$

$\alpha \leftarrow [\frac{\beta_c(x_t) - \beta_c(x)}{\epsilon_x}, \frac{\beta_d(y_t) - \beta_d(y)}{\epsilon_y}]$

$\text{cost} \leftarrow J[x, y, l] + \alpha \cdot C[x_t, y_t]$

**if**  $J[x_t, y_t, l + 1] > \text{cost}$  **then**

$J[x_t, y_t, l + 1] \leftarrow \text{cost}$

$\text{prev}[x_t, y_t, l + 1] \leftarrow (x, y)$

$\text{states} \leftarrow \text{states} \cup \{(x_t, y_t, l + 1)\}$

$\text{active\_states} \leftarrow \text{states}$

$x, y, \text{path} \leftarrow M - 1, M - 1, \emptyset$

$l \leftarrow \arg \min J[x, y, :]$

**while**  $l \geq 0$  **do**

$\text{path.append}((x, y))$

$(x, y) \leftarrow \text{prev}[x, y, l], l \leftarrow l - 1$

**return** *path*,  $\min J[M - 1, M - 1, :]$

---

**Convergence guarantee** As the grid resolution  $N \rightarrow \infty$ , the discrete solution converges to the continuous VLB-optimal path  $\beta^*(t)$ . Since the VLB for a joint schedule  $\beta(t)$  is a line integral over the model’s loss field in the 2D noise space (Appendix B.4), and training with the generalized loss  $\hat{\mathcal{L}}^\infty$  ensures the model’s predictions are accurate everywhere in this space, enabling the evaluation of VLBs for arbitrary paths  $\beta(t)$  (Appendix B.5), the dynamic programming solution recovers the optimal path asymptotically, maximizing the VLB over the function space  $\mathcal{Z}$ .

**Time complexity** The computational complexity is  $O(NM^2K)$ , which does not add too much computational overhead and can be solved within a few minutes on a single CPU. Once solved, the optimal schedule can be adopted for test time, enabling the VLB-optimal generative process, where the multi-modalities effectively inform each other at the optimal accuracy level. However, in order to estimate the cost matrix over  $M \times M$  discretized grid, the same number of evaluations of the generative model  $\tilde{x}_\phi$  is required, which scales quadratically to the discretization  $M$  and involves considerable GPU computation. Empirically, we found that the loss surface is smooth enough to allow for effective interpolation, therefore  $M = 20$  would yield an accurate enough estimate, taking less than 20 minutes.

Table 4. Validation loss as the sum of scaled KL divergence at sample time for different scheduling.

Schedule	Default	Optimal	Learned
Validation Loss	5.51	4.23	5.39



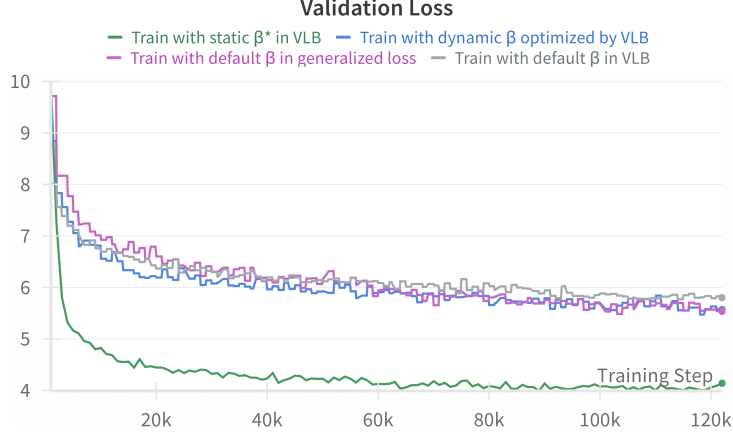


Figure 7. Validation losses for different configurations over the course of training.

### A.2. Why Learning for Optimal Schedule Might Fail

Given the obtained optimal schedule  $\beta^*(t)$ , we can train the generative model with invariant optimal VLB (Eq. 10) more effectively, allocating the model capacity towards learning to denoise only along the selected path. We found this results in earlier convergence, taking around half the previous training time for generalized objective (Eq. 9) with similar performance.

We would like to take one step further to see if we can directly optimize the joint noise schedule  $\beta(t)$  in generative training with the negative VLB objective in Eq. 10, i.e., optimizing the generative model and the noise schedule simultaneously through the same loss objective.

We parametrize the learnable time-rescaling function by

$$f(t) = t + (1 - t) \cdot \text{sigmoid}(h_\psi(t)) \quad (17)$$

where  $h_\psi(x)$  is a monotonic neural network which we choose to be a three-layered MLP with Softplus activation before the final output. This rescaling function is strictly monotonically increasing with endpoint constraints, i.e.,  $f(0) = 0, f(1) = 1$ .

However, we found that the learned noise schedule does not change much from its initialized shape, and the corresponding VLB proves to be suboptimal in Table 4. We hypothesize that there are several challenges in learning for the optimal schedule during training: (1) Large combinatorial space with limited exposure to possible noise schedules in the process of active learning. Unlike single-modality models, where the range of noise levels is easily explored within the modality (Dieleman et al., 2022), multi-modality generative models can suffer from exposure to only a subset of the possible noise combinations. This incomplete exploration makes it challenging to optimize directly for the best schedule, as large regions of the design space remain unvisited. It impacts the VLB estimation since the model  $\tilde{x}_\phi$  is still learning to adapt to the scaled noisy inputs during training, and is not invariant w.r.t. noisy input at different scales. (2) Optimization difficulties. Directing taking the gradient w.r.t.  $t_c, t_d$  introduces numerical instability, particularly when  $\tilde{\beta}_c$  involves exponentiation, which can grow or shrink very rapidly for certain values of  $t_c$ , leading to vanishing or exploding gradients. Moreover, the VLB optimization is sensitive to the trajectory defined by the noise schedule, and the optimization landscape of may have many local minima, especially due to the coupling of noise schedules, which can make it difficult to find the global optimum.

Addressing these might require more sophisticated methods such as bilevel optimization, which we leave for future work.

### A.3. Sensitivity to Optimal Path Choice

**Effect of Interpolation** We conduct additional experiments interpolating with a coefficient  $c$  between identity time function  $t_i$  ( $c = 0$ ) and our derived optimal time-rescaling functions  $t_o$  ( $c = 1$ ), by setting time functions  $t = c \cdot t_i + (1 - c)t_o$ . Our findings show a clear trend of improving performance as we move toward the optimal schedule.

Table 5 shows a generally monotonic improvement in binding pose quality on both datasets as we move from identity to optimal scheduling functions, with the most significant gains occurring when  $c \in [0.75, 1]$ . Additionally, the case with

Table 5. Results of interpolating between identity time function ( $c = 0$ ) and optimal time-rescaling functions ( $c = 1$ ), where for each protein target 10 molecules are sampled.

$c$	Vina Score ( $\downarrow$ )		Vina Min ( $\downarrow$ )		QED	SA	Connected	PB-Valid
	Avg.	Med.	Avg.	Med.	Avg.	Avg.	Avg. ( $\uparrow$ )	Avg. ( $\uparrow$ )
CrossDock (ID)	-6.36	-6.46	-6.71	-6.49	0.48	0.73	-	0.95
0 (Identity)	-6.63	-6.94	-7.06	-7.06	0.55	0.76	0.95	0.95
0.25	-6.66	-6.88	-7.02	-7.01	0.55	0.76	0.94	0.96
0.5	-6.70	-6.82	-7.04	-6.94	0.55	0.76	0.95	<b>0.97</b>
0.75	-6.87	-6.97	-7.21	-7.11	0.55	0.76	<b>0.96</b>	0.96
1 (Optimal)	<b>-6.92</b>	<b>-7.02</b>	<b>-7.23</b>	<b>-7.18</b>	0.55	0.75	<b>0.96</b>	0.95
PoseBusters (OOD)	-7.06	-7.05	-7.50	-7.41	0.40	0.72	-	0.99
0 (Identity)	-7.35	-7.38	-7.64	-7.48	0.48	0.73	0.94	0.79
0.25	-7.42	-7.37	-7.66	-7.50	0.48	0.73	0.94	0.79
0.5	-7.20	-7.36	-7.56	-7.51	0.48	0.73	<b>0.95</b>	0.78
0.75	-7.44	-7.43	-7.75	-7.54	0.48	0.73	<b>0.95</b>	0.79
1 (Optimal)	<b>-7.52</b>	<b>-7.52</b>	<b>-7.79</b>	<b>-7.65</b>	0.49	0.73	<b>0.95</b>	<b>0.80</b>

$c = 0$  also demonstrates the performance boost brought by our generalized training objective compared with the default objective (not shown in the table, please refer to Fig. 6).

**Effect of Accuracy Level  $\alpha$  as Scaling Factor** The desired scaling factor  $\alpha$  ought to be the derivative of joint schedule function  $\beta$ , which is not known before solving the search problem. To determine the appropriate accuracy level in practice, we have empirically found that we can actually approximate the accuracy by taking the derivative directly from the default schedule function  $\tilde{\beta}(t) = (\tilde{\beta}'_c(t), \tilde{\beta}'_d(t))$ , instead of taking the slope  $\frac{\beta(t) - \beta(t-\epsilon)}{\epsilon}$ , and this applies both to the experiments with BFN and diffusion models like TargetDiff. Note that by setting the accuracy level as such, the optimization objective no longer corresponds to the exact likelihood, but a rescaled sum of KL divergence terms that put more weight on the continuous variable. We hypothesize that this suggests a gap between likelihood estimation and sample quality, where the latter is more influenced by the 3D structure part.

#### A.4. Model Architecture

We employ the equivariant Graph Transformer architecture similar to Guan et al. (2023), with a few modifications to attention score calculations designed to save memory. Denote node positions as  $\mathbf{r} \in \mathbb{R}^{N \times 3}$ , one-hot atom types as  $\mathbf{h} \in \mathbb{R}^{N \times K_h}$ , and one-hot bond types as  $\mathbf{A} \in \mathbb{R}^{N \times (N-1) \times K_A}$ , where  $K_h$  and  $K_A$  are the numbers of atom and bond types, respectively. We brief the model architecture below.

**Heterogeneous Message Passing for Node Update** The heterogeneous update for node representation consists of two parts: (1) a  $K$ -nearest neighbors graph  $\mathcal{G}_K$  built for interactions within the protein-ligand complex, and (2) a fully connected ligand graph  $\mathcal{G}_L$  within ligand atoms for ligand bond-based interactions.

$$\Delta \mathbf{h}_{K,i} \leftarrow \sum_{j \in \mathcal{N}_K(i)} \phi_{h_K}(\mathbf{X}_{kv}^K, \mathbf{X}_q^K, \text{concat}(\mathbf{d}_{ij}, \mathbf{h}_{ij}^{\text{edge}})), \quad \Delta \mathbf{h}_{L,i} \leftarrow \sum_{j \in \mathcal{N}_L(i)} \phi_{h_L}(\mathbf{X}_{kv}^L, \mathbf{X}_q^L, \mathbf{h}_{ij}^{\text{bond}}),$$

where the input features for complex graph  $\mathbf{X}_{kv}^K = \text{concat}(\mathbf{h}_i, \mathbf{h}_j, \mathbf{h}^{\text{edge}})$ ,  $\mathbf{X}_q^K = \text{concat}(\mathbf{h}_i, \mathbf{h}^{\text{edge}})$ ,  $\mathbf{d}_{ij} = \|\mathbf{r}_i - \mathbf{r}_j\|$  is the pairwise distance, and  $\mathbf{h}^{\text{edge}}$  is the one-hot encoding for edge types between protein-ligand atoms (protein-protein, protein-ligand, ligand-protein or ligand-ligand). For ligand graph, the input features is defined as  $\mathbf{X}_{kv}^L = \text{concat}(\mathbf{h}_i, \mathbf{h}_j, \mathbf{h}_{\text{bond}})$ ,  $\mathbf{X}_q^L = \text{concat}(\mathbf{h}_i, \mathbf{h}_{\text{bond}})$ , where  $\mathbf{h}_{\text{bond}}$  denotes the flattened bond embeddings.

The node-level representation is updated by aggregating the heterogeneous messages:

$$\mathbf{h}_i \leftarrow \mathbf{h}_i + \text{FFN}(\mathbf{h}_i + (\Delta \mathbf{h}_{K,i} + \Delta \mathbf{h}_{L,i}) \mathbf{W}_h).$$

**Intra-Ligand Message Passing for Bond Update** The bond representation is updated by a directional message passing:

$$\mathbf{h}_{ji}^{\text{bond}} \leftarrow \sum_{k \in \mathcal{N}_L(j) \setminus \{i\}} \phi_E(\mathbf{X}_{kv}^E, \mathbf{X}_q^E, \sigma(\mathbf{d}_{ji})).$$

where the input features are defined as  $\mathbf{X}_{kv}^E = \text{concat}(\mathbf{h}_j, \mathbf{h}_k, \mathbf{m}_{kj})$ ,  $\mathbf{X}_q^E = \text{concat}(\mathbf{h}_{kj}^{\text{bond}}, \sigma(\mathbf{d}_{kj}), \sigma(\mathbf{d}_{ji}), \mathbf{a}_{ijk}, \mathbf{h}_k, \mathbf{h}_j)$ , and  $\mathbf{X}_q^E = \text{concat}(\mathbf{h}_{ji}^{\text{bond}}, \mathbf{h}_i)$ .  $\sigma$  is the Gaussian smearing function applied to the pairwise distance, and  $\mathbf{a}$  is the angular encoding for the bond triplets.

**Heterogeneous Message Passing for Position Updates** Similarly, the positions are updated from heterogeneous messages:

$$\Delta \mathbf{r}_{K,i} \leftarrow \sum_{j \in \mathcal{N}_K(i)} (\mathbf{r}_j - \mathbf{r}_i) \odot \phi_{r_K}(\mathbf{X}_{kv}^K, \mathbf{X}_q^K, \text{concat}(\mathbf{d}_{ij}, \mathbf{h}_{ij}^{\text{edge}})), \quad \Delta \mathbf{r}_{L,i} \leftarrow \sum_{j \in \mathcal{N}_L(i)} (\mathbf{r}_j - \mathbf{r}_i) \phi_{r_L}(\mathbf{X}_{kv}^L, \mathbf{X}_q^L, \mathbf{h}_{ij}^{\text{bond}})$$

given the updated node features  $\mathbf{h}_{ij}$  and bond features  $\mathbf{h}_{ij}^{\text{bond}}$ . The final position is updated by

$$\mathbf{r}_i \leftarrow \mathbf{r}_i + (\Delta \mathbf{r}_{K,i} + \Delta \mathbf{r}_{L,i}) \cdot \mathbf{1}_{\text{mol}}$$

where  $\mathbf{1}_{\text{mol}}$  is the boolean mask for ligand atoms, as the protein atoms need to be fixed to ensure equivariance.

**Attention Mechanism** The function  $\phi$  is implemented as a multi-headed attention mechanism. The query, key, value are obtained by passing input features through linear layer without bias to save memory:

$$\mathbf{K} = \mathbf{X}_{kv} \mathbf{W}_k, \quad \mathbf{V} = \mathbf{X}_{kv} \mathbf{W}_v, \quad \mathbf{Q} = \mathbf{X}_q \mathbf{W}_q.$$

The attention scores are computed using:

$$\alpha_{ij} = \frac{\mathbf{q}_i(\mathbf{k}_j \odot \mathbf{w}_{ij})}{\sqrt{d_h}}, \quad \mathbf{w}_{ij} = \tanh(\mathbf{e}_{ij} \mathbf{W}_{\text{edge}}),$$

where  $d_h$  is the head dimension,  $\mathbf{w}_{ij}$  is the edge-specific weights modulating the contribution of distant neighbors, and  $\odot$  denotes Hadamard product.

## A.5. Implementation Details

**Training and Inference** We use Adam optimizer with learning rate 5e-4, batch size of 16, and fit the model with 3.1 million parameters on one NVIDIA 80GB A100 GPU. We set  $\beta_1 = 1.5$  for discrete atom types and bond types,  $\sigma_1 = 0.05$  for atom coordinates. The training converges in 200K steps (around 24 hours). For inference, we use the exponential moving average of the weights from training that is updated at every optimization step with a decay factor of 0.999. We run inference with 100 sampling steps with the same variance reduction sampling strategy as Qu et al. (2024).

**Hyperparameters for Network** We set the network to be kNN graphs with  $k = 32$ ,  $N = 9$  layers with  $d = 128$  hidden dimension, 16-headed attention, and dropout rate 0.1.

**Featurization** We describe our atom-level and edge-level featurization for the protein-ligand complexes. At the atom-level, each protein atom is represented with a one-hot element indicator (H, C, N, O, S, Se), a 20-dimensional one-hot amino acid type indicator, and a 1-dimensional backbone flag. Ligand atoms are featurized with a one-hot element indicator (C, N, O, F, P, S, Cl) coupled with an aromaticity flag. For edge-level featurization in the heterogeneous protein-ligand graph, edge types are encoded as a 4-dimensional one-hot vector indicating whether the edge is between ligand atoms, protein atoms, ligand-protein atoms, or protein-ligand atoms. For the ligand graph, bonds are represented with a 4-dimensional one-hot bond type vector (non-bond, single, double, triple), where aromatic bonds are kekulized using RDKit.

## B. Proof

### B.1. Derivation of Eq. 6 as single-modality VLB (Eq. 5)

We begin with the negative Variational Lower Bound (VLB) in Eq. 5:

$$\mathcal{L}_{\text{VLB}}(\mathbf{x}) = -\mathbb{E}_{q(\mathbf{y}_{1:n}|\mathbf{x})}[\log p(\mathbf{x} | \mathbf{y}_{1:n}) + \sum_{i=1}^n D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel p(\mathbf{y}_i | \mathbf{y}_{1:i-1}))].$$

In further analysis, we focus on the treatment of the second term on the right-hand side. For convenience in the subsequent derivation, we define:

$$\begin{aligned}\mathcal{L}^n(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{y}_{1:n}|\mathbf{x})} \left[ \sum_{i=1}^n D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x})) \right] \\ &= n \mathbb{E}_{i \sim U(1,n), q(\mathbf{y}_{1:n}|\mathbf{x})} \left[ D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x})) \right]\end{aligned}\quad (18)$$

Next, to specify the conditional distribution  $p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x})$ , we utilize  $p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x}) = \frac{p(\mathbf{y}_{1:i}|\mathbf{x})}{p(\mathbf{y}_{1:i-1}|\mathbf{x})}$ , and incorporate the following decomposition:

$$p(\mathbf{y}_{1:i} | \mathbf{x}) = \prod_{j=1}^i \mathbb{E}_{p_F(\boldsymbol{\theta}_{j-1}|\mathbf{x}, t_{j-1})} [p_R(\mathbf{y}_j | \boldsymbol{\theta}_{j-1}, t_{j-1})] \quad (19)$$

It follows that:

$$p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x}) = \mathbb{E}_{p_F(\boldsymbol{\theta}_{i-1}|\mathbf{x}, t_{i-1})} [p_R(\mathbf{y}_i | \boldsymbol{\theta}_{i-1}, t_{i-1})] \quad (20)$$

Substituting this result into Eq. 18, we obtain:

$$\mathcal{L}^n(\mathbf{x}) = n \mathbb{E}_{i \sim U(1,n), \boldsymbol{\theta} \sim p_F(\boldsymbol{\theta}|\mathbf{x}, t_{i-1})} \left[ D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel p_R(\mathbf{y}_i | \boldsymbol{\theta}, t_{i-1})) \right] \quad (21)$$

Given our parametrization with neural network, we can further replace  $p_R(\mathbf{y}_i | \boldsymbol{\theta}, t_{i-1})$  with  $q(\mathbf{y}_i | \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t))$ :

$$\mathcal{L}^n(\mathbf{x}) = n \mathbb{E}_{i \sim U(1,n), \boldsymbol{\theta} \sim p_F(\boldsymbol{\theta}|\mathbf{x}, t_{i-1})} \left[ D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel q(\mathbf{y}_i | \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t))) \right] \quad (22)$$

Following Graves et al. (2023), different noising processes are specified for  $q$  for different modalities. For continuous data,  $q(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{x}, \alpha^{-1} \mathbf{I})$ . For discrete data,  $q(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\alpha(K\mathbf{e}_\mathbf{x} - \mathbf{1}), \alpha K \mathbf{I})$ , where we define  $\mathbf{x} := \sqrt{K} \mathbf{e}_\mathbf{x}$ . Denoting  $\mathcal{L}^\infty(\mathbf{x}) = \lim_{n \rightarrow \infty} \mathcal{L}^n(\mathbf{x})$ ,  $\alpha(t, \epsilon) = \beta(t) - \beta(t - \epsilon)$ , we have

$$\begin{aligned}\mathcal{L}^\infty(\mathbf{x}) &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \mathbb{E}_{i \sim U(1,n), p_F(\boldsymbol{\theta}|\mathbf{x}, t_{i-1})} D_{\text{KL}}(q(\mathbf{y}_i | \mathbf{x}) \parallel q(\mathbf{y}_i | \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t))) \\ &= \lim_{\epsilon \rightarrow 0} \mathbb{E}_{i \sim U(1,n), p_F(\boldsymbol{\theta}|\mathbf{x}, t_{i-1})} \frac{\alpha(t, \epsilon)}{2\epsilon} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)\|^2 \\ &= \frac{1}{2} \mathbb{E}_{i \sim U(1,n), p_F(\boldsymbol{\theta}|\mathbf{x}, t_{i-1})} \beta'(t) \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)\|^2\end{aligned}\quad (23)$$

which establishes the fact that Eq. 6 corresponds to the single-modality VLB.

## B.2. Equivalence between Eq. 6 and Eq. 7

Rewrite the expectation in the above equation into an integral form and make the substitution  $u = \beta(t)$ . This substitution can be transformed as follows:  $u = \beta(t) \Rightarrow du = \beta'(t)dt \Rightarrow dt = \frac{du}{\beta'(t)}$ , and we directly rewrite  $\tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)$  as  $\tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)$ .

$$\begin{aligned}\mathcal{L}^\infty(\mathbf{x}) &= \frac{1}{2} \int_0^1 \int_{p_F(\boldsymbol{\theta}|\mathbf{x}, t)} \beta'(t) \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)\|^2 d\boldsymbol{\theta} dt \\ &= \frac{1}{2} \int_{\beta(0)}^{\beta(1)} \int_{p_F(\boldsymbol{\theta}|\mathbf{x}, u)} \beta'(\beta^{-1}(u)) \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)\|^2 \frac{du}{\beta'(\beta^{-1}(u))} d\boldsymbol{\theta} \\ &= \frac{1}{2} \int_{\beta(0)}^{\beta(1)} \int_{p_F(\boldsymbol{\theta}|\mathbf{x}, u)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)\|^2 du d\boldsymbol{\theta} \\ &= \frac{1}{2} \int_{\beta(0)}^{\beta(1)} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x}, u)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, u)\|^2 du\end{aligned}\quad (24)$$

which gives the form of Eq. 7.



### B.3. Equivalence between Eq. 9 and 14

Recall that the generalized loss over  $\beta \in \mathcal{Z}$  is:

$$\mathcal{L}^\infty(\mathbf{x}) = \frac{1}{2} \int_{\beta_c, \beta_d} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};\beta)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, \beta)\|^2 d\beta.$$

and the function space  $\mathcal{Z}$  is reparameterized into a product of uniform distributions over  $t_c, t_d$ , yielding

$$\dot{\mathcal{L}}^\infty(\mathbf{x}) = \frac{1}{2} \int_0^1 \int_0^1 \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};t)} \|\mathbf{x} - \tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)\|^2 dt_c dt_d,$$

this is equivalent to implicitly sampling  $f(t), g(t)$ , covering all possible time-rescaling functions as stated in Remark 4.2.

The joint distribution over  $(\beta_c, \beta_d)$  is thus

$$p(\beta_c, \beta_d) = p(f, g) = U(0, 1) \times U(0, 1), \quad (25)$$

where  $\beta_c(t) = \tilde{\beta}_c(f(t))$  and  $\beta_d(t) = \tilde{\beta}_d(g(t))$ .

Therefore, we see that the time-rescaling functions  $f(t), g(t)$  act as latent variables defining the coupling of  $\beta_c, \beta_d$ . Integrating over  $(t_c, t_d)$  marginalizes over all possible  $\beta$ , ensuring the model is trained to denoise all possible combinations of noise levels. By demonstrating this equivalence, we reframe the invariant objective described in Eq. 9 as Eq. 14.

### B.4. Derivation of Eq. 10 as multi-modality VLB

To rigorously prove that the VLB corresponds to a line integral along a 1D submanifold (path) in the 2D noise schedule space, recall that the joint noise schedule be a parameterized path in the joint function space  $\mathcal{Z}$ :

$$\beta(t) = (\beta_c(t), \beta_d(t)), \quad t \in [0, 1],$$

where  $\beta_c(t), \beta_d(t)$  are monotonically increasing functions with fixed endpoints:

$$\beta_c(0) = \beta_{c,0}, \quad \beta_c(1) = \beta_{c,1}, \quad \beta_d(0) = \beta_{d,0}, \quad \beta_d(1) = \beta_{d,1}.$$

Following Appendix B.1 and the factorized nature of  $q$ , the VLB for a joint schedule  $\beta(t)$  and multi-modality  $\mathbf{x}$  is:

$$\mathcal{L}_{\text{VLB}}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{p_F(\boldsymbol{\theta}|\mathbf{x};t)} \int_0^1 [\beta'_c(t) \|\mathbf{x}_c - \tilde{\mathbf{x}}_{\phi,c}(\boldsymbol{\theta}, \beta(t))\|^2 + \beta'_d(t) \|\mathbf{x}_d - \tilde{\mathbf{x}}_{\phi,d}(\boldsymbol{\theta}, \beta(t))\|^2] dt,$$

where  $\tilde{\mathbf{x}}_\phi = (\tilde{\mathbf{x}}_{\phi,c}, \tilde{\mathbf{x}}_{\phi,d})$  is the denoising model output for continuous ( $c$ ) and discrete ( $d$ ) modalities.

The VLB can be interpreted as a line integral over the trajectory  $\beta(t) \in \mathcal{Z}$ . Define the vector field  $\mathbf{F}(\beta_c, \beta_d)$  as:

$$\mathbf{F}(\beta_c, \beta_d) = \begin{pmatrix} \|\mathbf{x}_c - \tilde{\mathbf{x}}_{\phi,c}(\boldsymbol{\theta}, \beta)\|^2 \\ \|\mathbf{x}_d - \tilde{\mathbf{x}}_{\phi,d}(\boldsymbol{\theta}, \beta)\|^2 \end{pmatrix}.$$

The VLB then becomes:

$$\mathcal{L}_{\text{VLB}}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{\boldsymbol{\theta}} \int_{\beta(t)} \mathbf{F}(\beta_c, \beta_d) \cdot d\beta,$$

where  $d\beta = (\beta'_c(t)dt, \beta'_d(t)dt)$  is the differential vector along the path  $\beta(t)$ , and  $\cdot$  denotes inner product. By redefining  $\tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, \beta) := \mathbf{F}(\beta_c, \beta_d)$ , we obtain Eq. 10.

### B.5. Derivation of Proposition 4.3

*Proof.* Since the model  $\tilde{\mathbf{x}}_\phi(\boldsymbol{\theta}, t)$  is trained by the generalized loss  $\dot{\mathcal{L}}^\infty$ , which is equivalent to integrating the scalar field  $\|\mathbf{x} - \tilde{\mathbf{x}}_\phi\|^2$  over the entire function space  $\mathcal{Z}$ , it is exposed to all pairs  $(\beta_c, \beta_d) \in \mathcal{Z}$  within the region  $[\tilde{\beta}_c(0), \tilde{\beta}_c(1)] \times [\tilde{\beta}_d(0), \tilde{\beta}_d(1)]$ , with  $\tilde{\beta}_c, \tilde{\beta}_d$  defined in Eq. 3. Therefore, the model achieves minimal prediction error everywhere in the joint noise space, allowing accurate computation for any specific path  $\beta(t)$ .

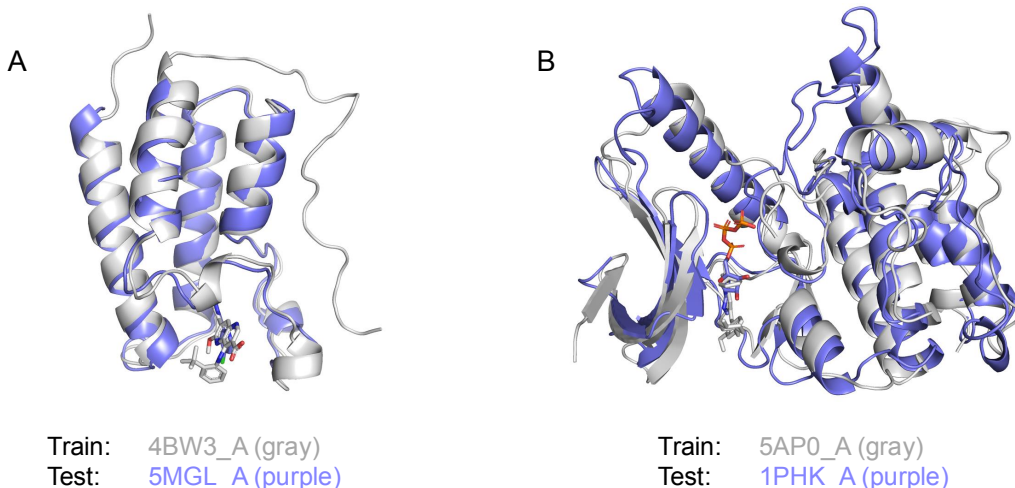


Figure 8. Structural alignment of randomly selected cases where training sequence and test sequence exhibit  $> 30\%$  identity, showing nearly identical protein structures as well as ligand binding pockets.

Following Appendix B.4, the line integral  $\mathcal{L}^\infty(\mathbf{x})$  corresponds to the VLB for a given path  $\beta(t)$ . Since the integral is equivalent to restricting the generalized loss  $\hat{\mathcal{L}}^\infty(\mathbf{x})$  to the 1D manifold defined by  $\beta(t)$ , and the model’s predictions are accurate along any submanifold within the space, evaluating  $\mathcal{L}^\infty(\mathbf{x})$  for a specific  $\beta(t)$  uses the model’s pre-optimized predictions at each point on the path, thereby yielding a valid VLB estimate. This establishes that the generalized training enables estimating VLBs for arbitrary joint schedules  $\beta \in \mathcal{Z}$ .  $\square$

## C. Details on Benchmarking

### C.1. Problems with CrossDock

Table 6. Test protein chains that exhibit a sequence identity higher than 30% to any chain of training sequences. For those with more than 10 overlapping training chains, we randomly select 10 items to display, and report the total number of unique overlapping chains. Each item is characterized by its PDB ID and chain ID.

Test PDB	#Overlap	Training PDB(s) with Sequence Identity $> 30\%$
1PHK_A	104	2WQO_A, 3TKU_B, 3WYY_A, 3WYX_A, 4C4E_A, 4BHZ_A, 2XKD_A, 4ZEG_A, 4CV8_A, 5E18_A
1UMD_B	1	1NI4_B
5NGZ_A	5	5F6Y_A, 5F6W_A, 5F6D_A, 5F6E_A, 5F6U_A
5L1V_A	5	4DNJ_A, 1Z8Q_A, 1Z8O_A, 1EGY_A, 1JIP_A
1FMC_B	8	5L7T_A, 5ICS_C, 5ICM_A, 5L7W_A, 5HS6_A, 5EN4_A, 5L7Y_A, 5JS6_A
4PXZ_A	4	4Z34_A, 4Z35_A, 5WIV_A, 4Z36_A
3HY9_B	1	1R55_A
4AZF_A	6	5J1V_B, 5X8I_B, 2B9J_A, 1Z57_A, 2VAG_A, 5J1W_C
3PDH_A	4	1YC2_D, 1YC2_B, 1S7G_B, 1S7G_C
5MGL_A	127	5UVW_A, 5D0C_A, 4O7B_A, 5KU3_A, 5DX4_A, 4A9F_A, 5CS8_A, 5D3L_A, 4O7C_A, 5E0R_A
4XLI_B	43	4OTF_A, 4Y93_A, 4Y95_A, 5JRS_A, 3GEN_A, 5P9H_A, 4NWM_A, 4ZLY_A, 3PIY_A, 4RFY_A
4AUA_A	75	2R3Q_A, 3EZR_A, 3TN8_A, 5FGK_A, 2R3M_A, 5IDN_A, 2YIY_A, 5HBJ_A, 4CFU_A, 3DOG_A
5IOB_A	100	5OPB_A, 5OPR_A, 2YM4_A, 2C3K_A, 3F69_B, 4FTT_A, 4FTO_A, 3TKH_A, 2BRG_A, 2YDK_A
4FIM_A	6	5LPV_A, 5LPW_A, 5LPY_A, 4OH4_B, 4Q5J_A, 5LPB_A
2CY0_A	10	3PHJ_A, 4FQ8_B, 3PHH_A, 3PGJ_C, 3DON_A, 4FOS_A, 3PHG_A, 3DOO_A, 4FR5_B, 4FPX_A

The CrossDock dataset, first processed by Luo et al. (2021), is a commonly used benchmark in the SBDD field. Luo et al. (2021) claimed to have employed the MMseqs2 method to filter the test set with a 30% sequence similarity threshold. However, upon rigorous examination, we found that the test set still contains proteins with sequence similarity higher than 30%.

Table 6 summarizes the calculated statistics for these similar proteins, where we only show the randomly sampled 10 PDB

IDs when there are more than 10 sequences that exceed  $> 30\%$  sequence identity. Fig. 8 illustrates structural alignments between some random proteins in the test and training sets that exhibit high sequence similarity. The near-identical overlap suggests that the dataset split may not be challenging enough. To address these potential problems with CrossDock evaluation, we propose using PoseBusters (Buttenschoen et al., 2024) as a held-out test set.

### C.2. Curation of the held-out PoseBusters test set

The PoseBusters Benchmark (Buttenschoen et al., 2024) set exclusively contains complexes released since 2021. In contrast, the CrossDock dataset comprises data collected from the PDB before 2020. This makes PoseBusters a challenging time-split dataset, ideal for evaluating the generalizability of models to real-world, unseen scenarios. We further apply MMseqs clustering (Steinegger & Söding, 2017), and filter any test protein with any chain that has  $> 30\%$  sequence identity threshold to CrossDock training sequences. This process leaves us with 180 data points, making most of the baselines directly available for the curated held-out test. Additionally, the protein-ligand complex structures in CrossDock are mainly generated using docking software, which inevitably introduces noisy poses (Francoeur et al., 2020). In contrast, PoseBuster contains real-world crystal structures resolved from wet-lab experiments, making it suitable for evaluating the SBDD model’s ability to capture genuine molecular interactions accurately.

### C.3. Experimental Setup

**Baselines** We provide a brief overview of all SBDD baselines as follows:

- **Autoregressive methods:** AR (Luo et al., 2021) utilizes MCMC sampling to reconstruct molecules atom-by-atom based on voxel-wise density predictions. Pocket2Mol (Peng et al., 2022) generates molecules atom-by-atom with bonds using an E(3)-equivariant network, predicting frontier atoms to improve sampling efficiency. FLAG (ZHANG et al., 2023) is a fragment-based model that assembles molecular fragments by predicting their positions and torsion angles.
- **Diffusion:** DiffSBDD (Schneuing et al., 2022) constructs an E(3)-equivariant continuous diffusion model for full-atom generation, applying noise to both atom types and coordinates, while TargetDiff (Guan et al., 2022) adopts a hybrid diffusion process to separately handle continuous coordinates and discrete atom types. DecompDiff (Guan et al., 2023) incorporates chemical priors by decomposing molecules into scaffolds and contact arms.
- **BFN:** MolCRAFT (Qu et al., 2024) uses Bayesian Flow Network (BFN) with advanced variance reduction sampling technique, demonstrating notable improvements over diffusion-based models in molecular design.

**Metrics** We evaluate the generated molecules using the following commonly adopted metrics:

- **Affinity Metrics** are calculated using AutoDock Vina (Eberhardt et al., 2021), these include *Vina Score* as the raw binding energy of a molecular pose in the pocket, *Vina Min* as the binding energy after local energy minimization of the molecular pose, and *Vina Dock* as the lowest binding energy obtained after an extended search for the optimal pose.
- **Molecular Properties** including **QED** for drug-likeness and **SA** (synthetic accessibility score) are calculated using RDKit. They are desired to fall within reasonable ranges.
- **Connected Ratio** is the percentage of fully connected molecules.
- **Diversity** assesses the variety of generated molecules for each binding site by averaging Tanimoto similarity over Morgan fingerprints across all test proteins, following Luo et al. (2021). It is worth noting that this is not necessarily the higher the better, as the desirable bioactive compounds against a protein target often cluster in the molecular space.
- **Key Interactions** such as the formation of critical non-covalent interactions with protein binding sites, calculated using ProLIF (Bouysset & Fiorucci, 2021).
- **Strain Energy** reflects the internal energy of generated poses, indicating pose quality by Harris et al. (2023).
- **RMSD** reports the percentage of molecules where the RMSD between generated poses and ground-truth or Vina redocked poses is within 2 Å, indicating consistent binding modes. For better differentiation, we refer to the latter as self-consistency RMSD (scRMSD).

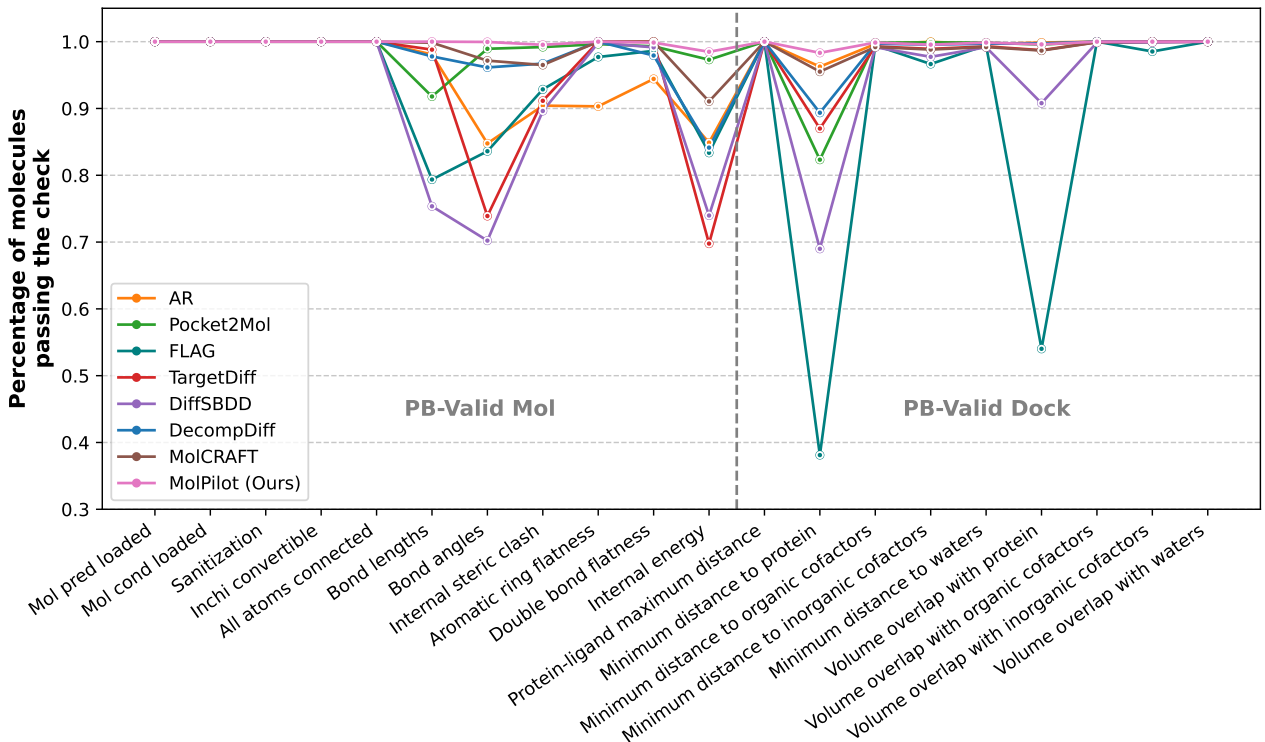


Figure 9. Percentage of generated molecules that have passed the PoseBusters validity checks on ID CrossDock test set. *PB-Valid Mol*: intramolecular validity. *PB-Valid Dock*: intermolecular validity. Reported *PB-Valid*: *PB-Valid Mol* & *PB-Valid Dock*.

## D. More Evaluation Results

### D.1. Dynamics of Better-balanced Modalities

We visualize the modality-specific validation loss curves for the model trained under generalized loss (Eq. 9) in Fig. 11, which shows that the denoising model effectively learns to mutually benefit from cleaner information either modality.

### D.2. In-Distribution De novo Design

We report the results on CrossDock (Francoeur et al., 2020) in an in-distribution (ID) setting. We sample 100 molecules for each of the 100 test proteins, and the results in Table 1 show that our model consistently performs exceptionally across most metrics—leading in PB-Valid, overall Vina affinities, and achieving comparable drug-like properties w.r.t. reference molecules.

For conformation quality, our model achieves the best performance with the highest PoseBusters passing rate, closely matching the reference value of 95.0%. Details of each validity check are shown in Fig. 9, and it can be seen that the leading factors affecting the overall performance are bond lengths and angles (AR, TargetDiff), internal energy (MolCRAFT, DecompDiff) for intramolecular validity, and minimum distance to protein (Pocket2Mol) for intermolecular validity.

For molecular geometries, we visualize the distributions of bond length, bond angle, and torsion angle, between generated molecules and reference molecules in the test set in Fig. 14, 15, 16 for the top-5 frequent types, and summarize the overall Jensen-Shannon Divergence (JSD) in Table 7, averaged over all types with a frequency > 100. While the previous strong-performing method MolCRAFT captures the bond length distributions relatively well, it cannot fit all the bond angle or torsion angle distributions, displaying for example non-standard C-C-C bond angle, smoothed C-C-N bond angle, and diverges more severely in C-C-C-C and C-C-N-C torsion angles compared with reference. Our method excels at modeling all these distributions, underscoring its ability to capture the molecular geometries.

For binding affinities, our model demonstrates the leading performance with an average Vina score of -6.88 kcal/mol,



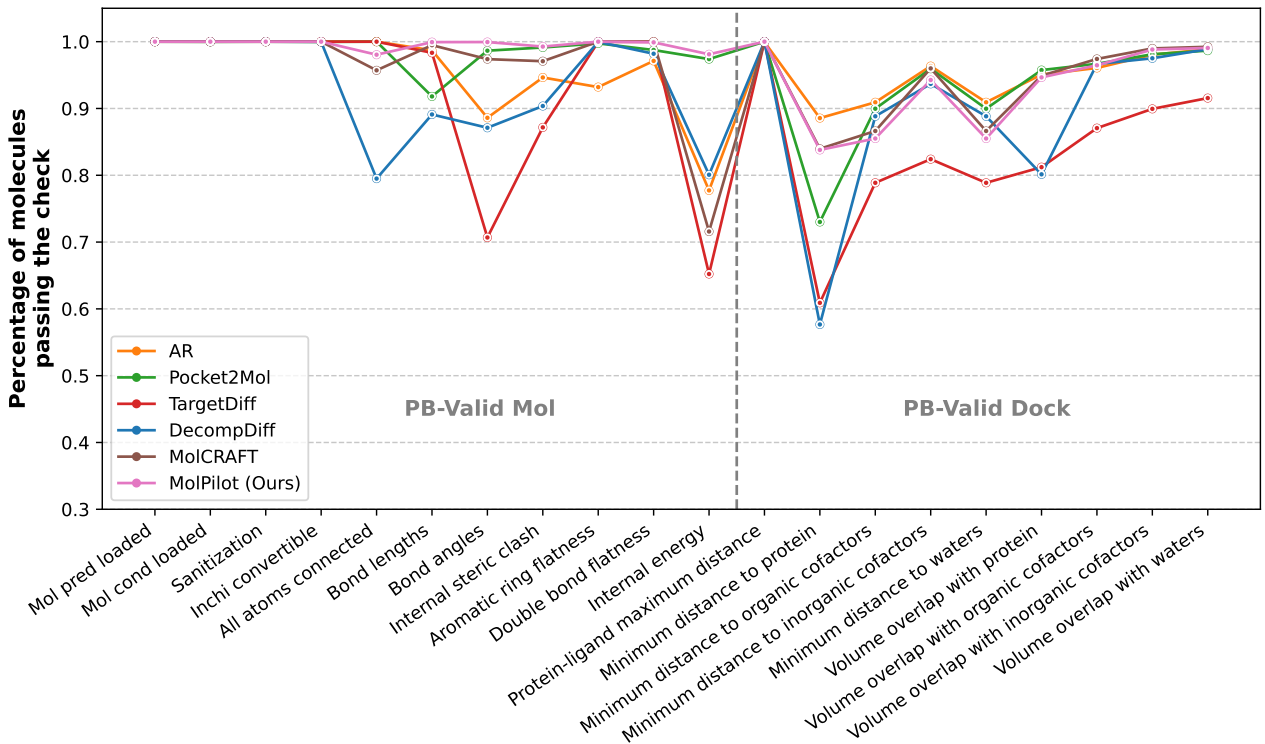


Figure 10. Percentage of generated molecules that have passed the PoseBusters validity checks on OOD PoseBusters test set. *PB-Valid Mol*: intramolecular validity. *PB-Valid Dock*: intermolecular validity. Reported *PB-Valid*: *PB-Valid Mol* & *PB-Valid Dock*.

maintaining the closest gap between Vina Score, Vina Min, and Vina Dock. This shows our model’s superiority in capturing spatial interactions in the generated pose.

For molecular properties like QED and SA, our model displays competitive performance, and ranks at the top among all non-autoregressive models.

### D.3. Out-of-Distribution De novo Design

We select competitive baselines in the ID settings, and evaluate them on the OOD PoseBusters test set. We sample 100 molecules for each of the 180 test proteins, and report the overall results in Table 1.

We provide more detailed statistics of PoseBusters passing rate in Fig. 10. It can be seen that the performance drop in PB-Valid is mainly attributed to two major factors: (1) *Intramolecular validity* drops due to distorted molecular geometry such as non-standard bond lengths, bond angles, ring flatness and the most prominently, internal energy as can be viewed as the overall indicator of molecular geometries, where all the non-autoregressive baselines and AR are around or below 80%. (2) *Intermolecular validity* drops due to clashes w.r.t. protein surface. TargetDiff appears to be mostly affected, and nearly all models cannot effectively avoid clashing by maintaining a feasible distance to protein atoms, waters or other cofactors. For protein atom clashing, it is worth pointing out that the high performance over CrossDock might come from the potential information leakage caused by data splits (Appendix C.1). Thus the passing rates on held-out PoseBusters test set can serve as a more informative indicator of the performances between different methods, showing that voxel-based AR and BFN-based models are performing better than diffusion-based models. However, it should be noted that though with better clashing performance, AR generates significantly fewer molecules that are also with smaller sizes, undermining its practicality. Besides, current SBDD models have not yet taken waters and cofactors into consideration, suggesting that there might be a calling for an advanced and comprehensive formulation of SBDD.

For a finer-grained quantification of the internal energy, we employ PoseCheck test suite (Harris et al., 2023) to measure the Strain Energy of generated molecules, and report the percentiles of distributions in Table 8, associated with the internal

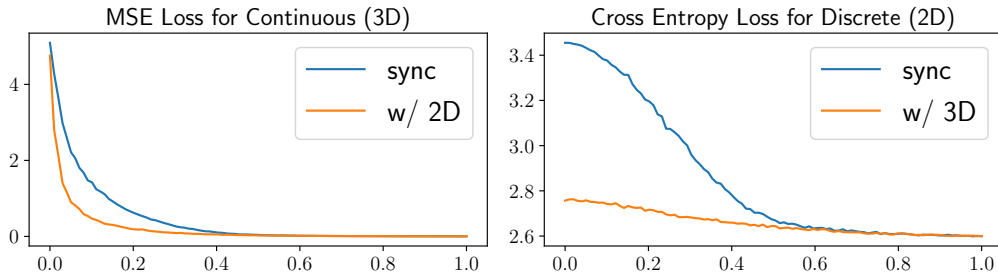


Figure 11. Validation loss curves for the model trained under generalized objective w.r.t. timestep in the generative process. *Sync*: Modalities at the same timestep. *w/ 2D*: Discrete modality always at  $t = 1$ . *w/ 3D*: Continuous modality always at  $t = 1$ .

Table 7. Results of Jensen-Shannon Divergence (JSD) of bond length, bond angle, torsion angle distributions between generated molecules and CrossDock reference molecules, averaged over all types with a frequency  $> 100$ .

Methods	Length ( $\downarrow$ ) Avg. JSD	Angle ( $\downarrow$ ) Avg. JSD	Torsion ( $\downarrow$ ) Avg. JSD
AR	0.544	0.507	0.545
Pocket2Mol	0.472	0.482	0.467
TargetDiff	0.365	0.435	0.411
DecompDiff	0.332	0.410	0.338
MolCRAFT	0.318	0.384	0.322
Ours	<b>0.252</b>	<b>0.351</b>	<b>0.287</b>

Table 8. Strain Energy results calculated by PoseCheck v1.1, and PoseBusters internal energy passing rate for SBDD baselines.

Methods	Strain Energy			
	25%	50%	75%	Passed
Reference	2.6	7.1	18.9	100%
AR	5.6	40.1	238.8	77.7%
Pocket2Mol	<b>1.0</b>	<b>4.4</b>	<b>28.5</b>	<u>97.4%</u>
TargetDiff	41.4	158.2	625.1	65.2%
DecompDiff	6.8	41.5	249.2	80.1%
MolCRAFT	5.3	46.5	27397.8	71.6%
Ours	<u>2.7</u>	<u>10.6</u>	<u>38.1</u>	<b>98.1%</b>

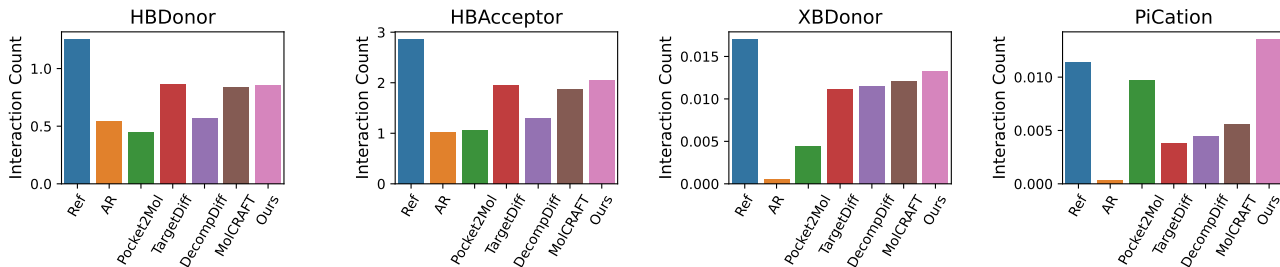


Figure 12. Detailed interaction counts for generated molecules on PoseBusters test. Only PB-Valid molecules are considered.

energy passing rate in PoseBusters checks. It shows that MolCRAFT generates considerably more strained structures as suggested by the tail distribution, and our MolPilot and Pocket2Mol maintains the robust performance.

We provide the detailed interaction counts compared with ground-truth complex structures in PoseBusters. Fig. 12 shows that MolPilot matches the genuine interaction profiles relatively well, especially in Cation- $\pi$  interactions. We additionally report in Table 9 the full statistics of interaction profile similarity calculated from different poses. It can be seen that our MolPilot and MolCRAFT closely match the interaction fingerprint, but MolCRAFT suffers from strained poses, resulting in unrealistic interactions for generated poses with high internal energy. In contrast, our model is able to capture the true interactions reliably.

Table 9. Tanimoto similarity of the interaction profiles between ground-truth structures in PoseBusters and molecules generated by SBDD models. *Gen*: interactions directly calculated from generated poses. *Redock*: interactions calculated from Vina redocked poses. *Gen & PB-Valid*: interactions directly calculated from generated poses that have passed the PoseBusters validity checks.

Methods	Sim. (Gen)	Sim. (Redock)	Sim. (Gen & PB-Valid)
AR	0.394	0.344	0.221
Pocket2Mol	0.447	0.405	0.436
TargetDiff	0.477	0.415	0.458
DecompDiff	0.438	0.332	0.374
MolCRAFT	<b>0.560</b>	<u>0.444</u>	<u>0.498</u>
Ours	<u>0.552</u>	<b>0.477</b>	<b>0.551</b>

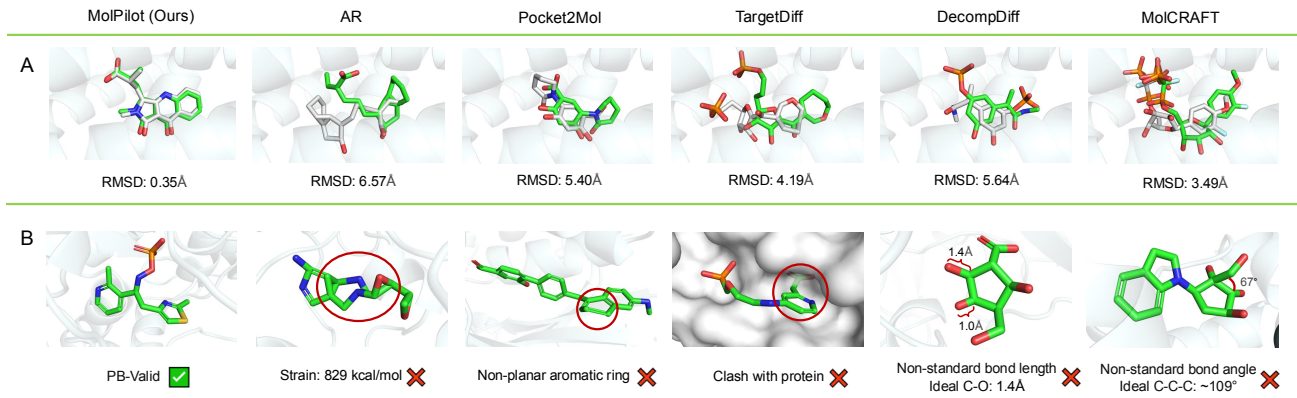


Figure 13. Illustration of molecular geometries for SBDD models. **A**. Binding mode consistency upon redocking, measured by scRMSD between generated poses (Green) and redocked poses (Silver). **B**. Generated molecules with unrealistic geometries, thus failing to pass PoseBusters validity checks.

#### D.4. Molecular Docking

We repurpose the SBDD models from the joint codesign model  $P(\mathbf{x}, \mathbf{h}, \mathbf{A} \mid \mathbf{x}_P)$  to the conditional marginal  $P(\mathbf{x} \mid \mathbf{h}, \mathbf{A}, \mathbf{x}_P)$ . This actually corresponds to another specific line on the loss surface described by Eq. 10, with  $\beta_d(t) \equiv \beta_d(1)$ . We employ the joint schedule  $\beta(t) = (\tilde{\beta}_c(t), \tilde{\beta}_d(1))$  with  $\tilde{\beta}$  being the default noise schedule for continuous modality in Eq. 3.

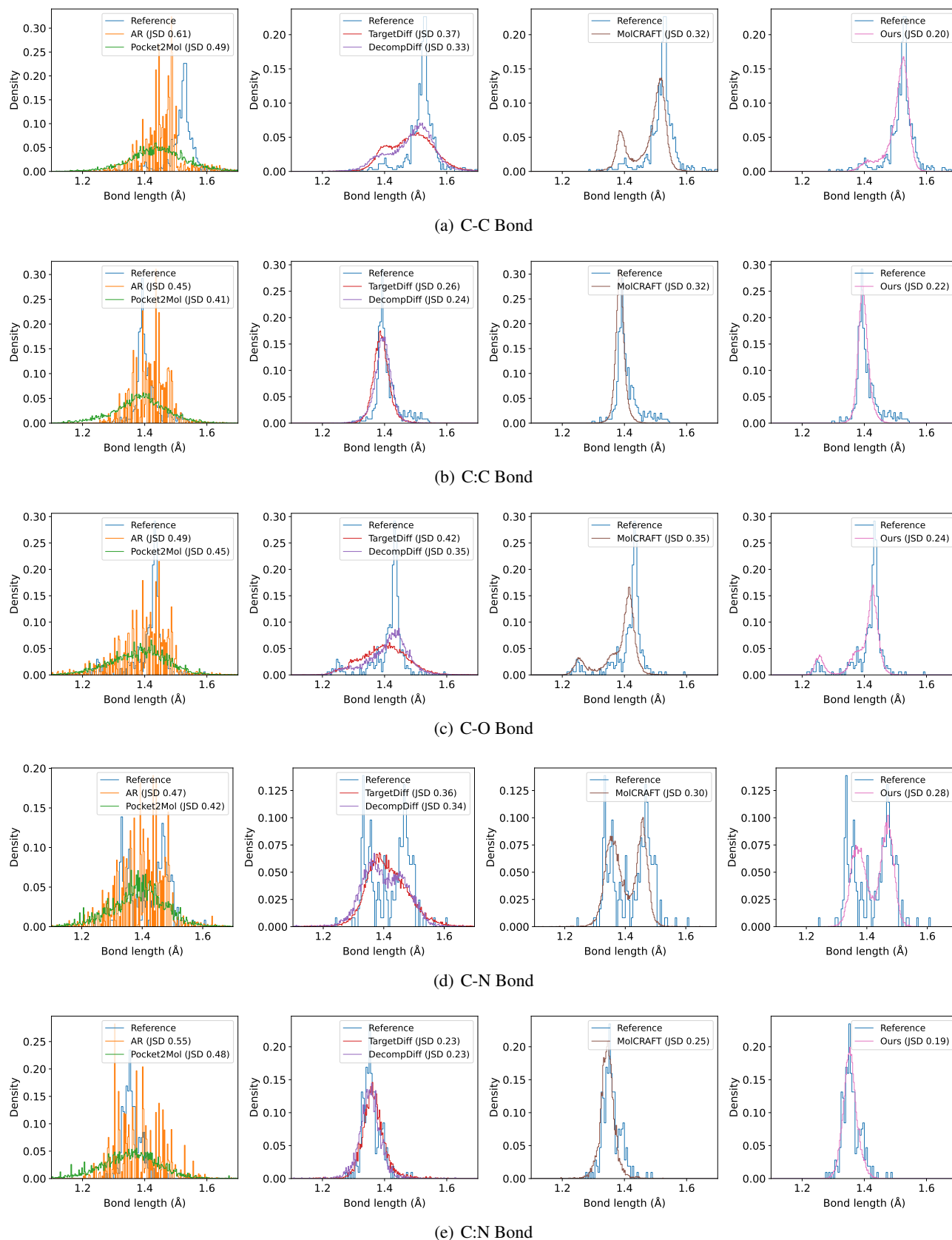


Figure 14. Top-5 frequent bond length distribution of generated molecules compared with CrossDock reference molecules.



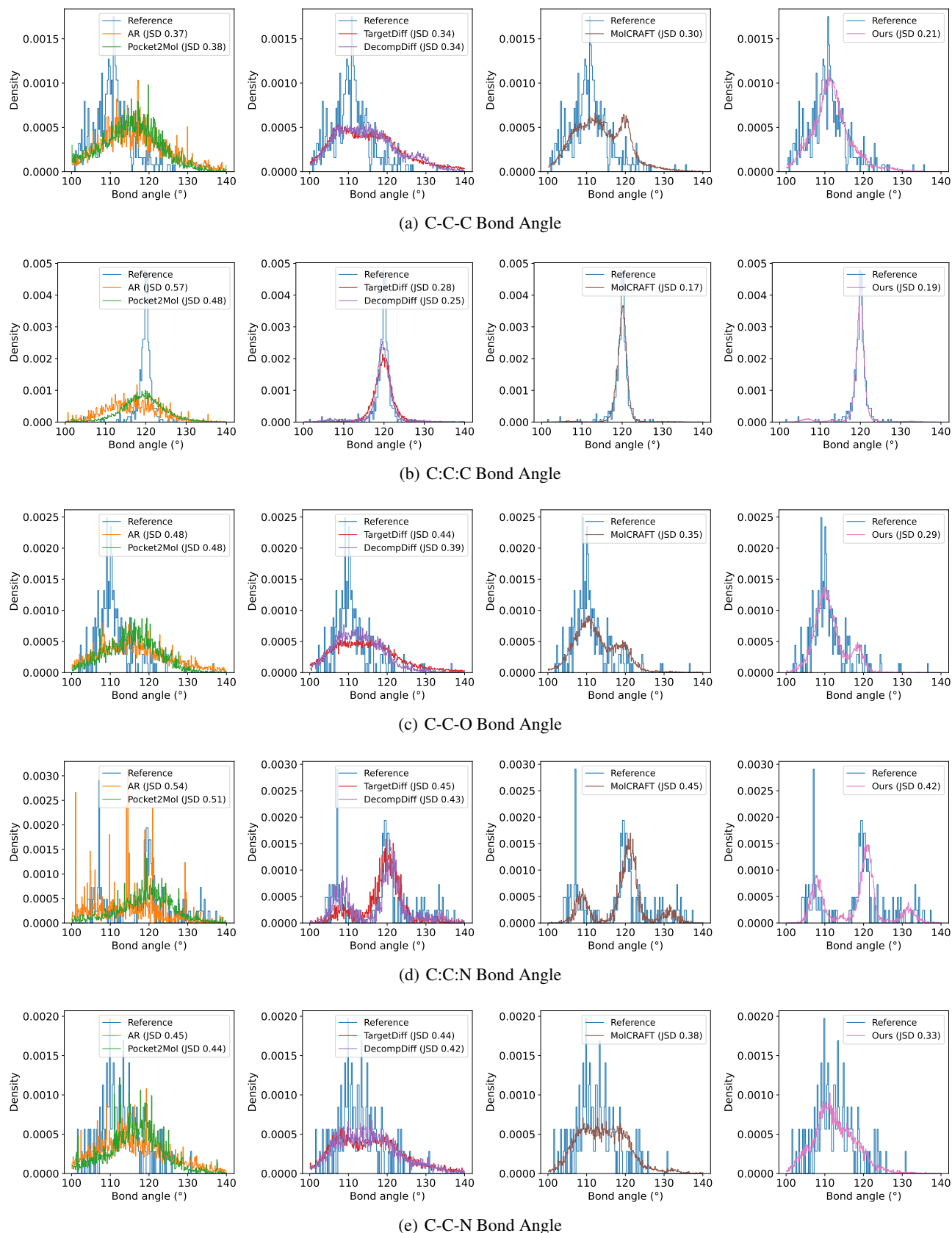


Figure 15. Top-5 frequent bond angle distribution of generated molecules compared with CrossDock reference molecules.

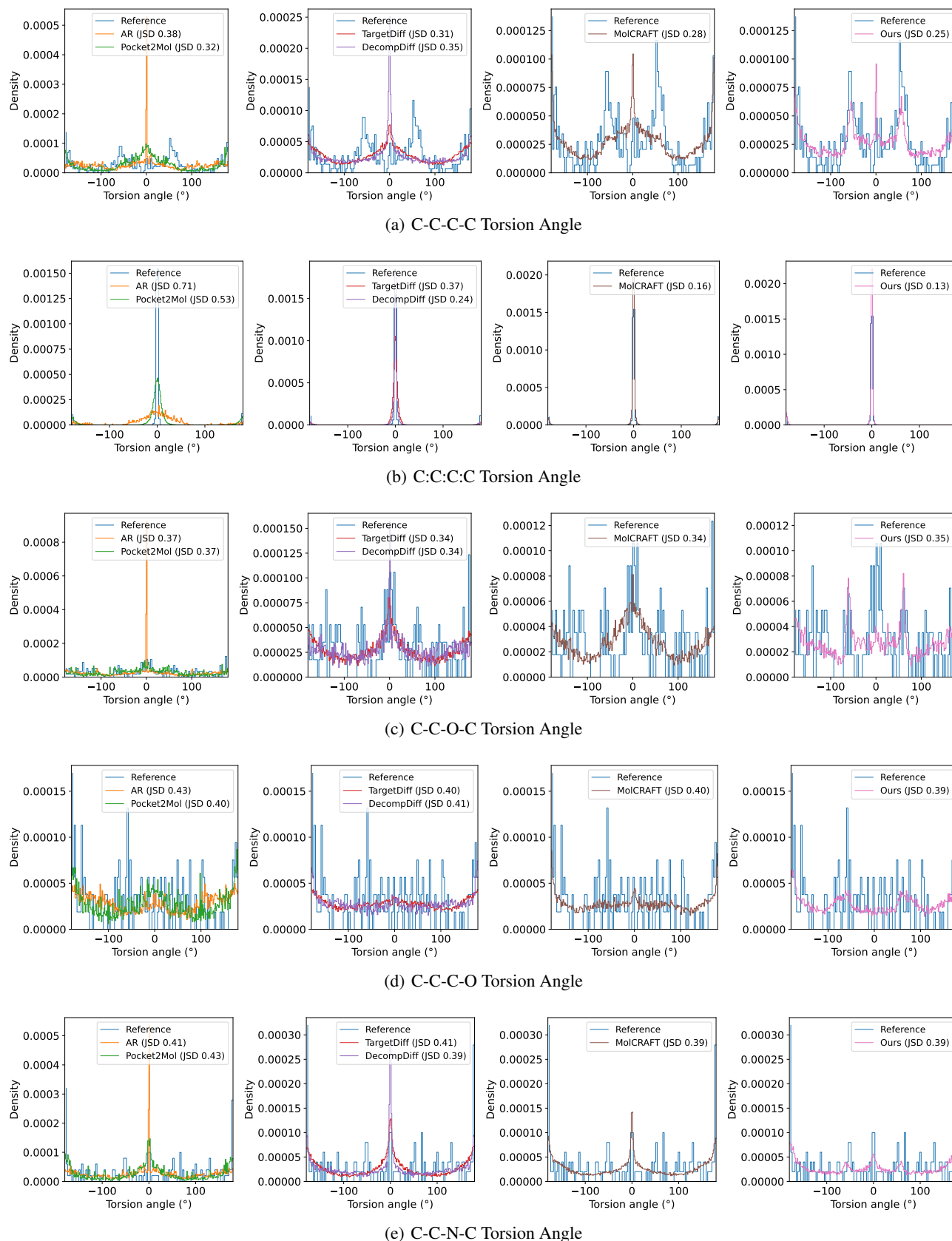


Figure 16. Top-5 frequent torsion angle distribution of generated molecules compared with CrossDock reference molecules.