FLOW DIVERSE AND EFFICIENT: LEARNING MOMENTUM FLOW MATCHING VIA STOCHASTIC VELOCITY FIELD SAMPLING

Anonymous authors

000

001

002

004

006

008 009 010

011

013

014

015

016

018

019

021

024

025

026

027

028

029

031

033

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Recently, the rectified flow (RF) has emerged as the new state-of-the-art among flow-based diffusion models due to its high efficiency advantage in straight path sampling, especially with the amazing images generated by a series of RF models such as Flux 1.0 and SD 3.0. Although a straight-line connection between the noisy and natural data distributions is intuitive, fast, and easy to optimize, it still inevitably leads to: 1) Diversity concerns, which arise since straight-line paths only cover a fairly restricted sampling space. 2) Multi-scale noise modeling concerns, since the straight line flow only needs to optimize the constant velocity field vbetween the two distributions π_0 and π_1 . In this work, we present Discretized-RF, a new family of rectified flow (also called momentum flow matching models since they refer to the previous velocity component and the random velocity component in each diffusion step), which discretizes the straight path into a series of variable velocity field sub-paths (namely "momentum fields") to expand the search space, especially when close to the distribution p_{noise} . Different from the previous case where noise is directly superimposed on x, we introduce noise on the velocity v of the sub-path to change its direction in order to improve the diversity and multi-scale noise modeling abilities. Experimental results on several representative datasets demonstrate that learning momentum flow matching by sampling random velocity fields will produce trajectories that are both diverse and efficient, and can consistently generate high-quality and diverse results.

1 Introduction

Flow-based diffusion models (Lipman et al.; Bartosh et al., 2024; Luo et al., 2024; Liu et al., 2023b) have recently attracted widespread attention, which generate a wide variety of realistic natural images from pure noise by modeling trajectories from noise distributions to data distributions. As a milestone work, the most popular flow models currently are the rectified flow (RF) models (Liu et al., 2023a) built

upon straight-line trajectories, which significantly improves their sampling efficiency by establishing the shortest straight-line connection between the noise distribution π_0 and the data distribution π_1 . Furthermore, the models can be easily optimized by directly calibrating this straight-line trajectory $dx_t/dt = v_\theta$ at a constant rate $x_1 - x_0$. Due to its high sampling efficiency (even enabling one-step diffusion generation), RF is also considered to be one of the fastest flow-based optimal transport models.

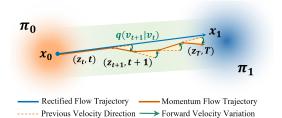


Figure 1: Graphical momentum flow trajectories.

Momentum Flow (orange) vs. Rectified Flow (blue).

Though remarkable success has been witnessed, RF still suffers from limitations in *diversity* and *multi-scale noise modeling*. Specifically, *1) Diversity concerns*, which arise since the straight-line path can only cover a fairly restricted sampling space. *2) Multi-scale noise modeling concerns*, since the straight-line flow only needs to directly optimize the constant velocity field $v_{\theta} \rightarrow (x_1 - x_0)$ between the two distributions π_0 and π_1 , instead of considering multi-scale progressive denoising. At the other extreme, the diffusion probability models (*e.g.*, DDPM) based on fluctuation trajectories

have extremely strong diversity and multi-scale noise modeling capabilities but face the challenge of training- and sampling-efficiency because they require a large number of time steps to sample, and each step should be optimized iteratively to achieve high-fidelity modeling of the reverse trajectory. In this work, to strike a balance and take into account both efficiency and diversity (especially the potential diversity when close to noise distribution π_1), we propose a Discretized-RF model, also known as the momentum flow matching (MFM) model. For clarity, we first give a unified definition of the flow transport problem and then introduce our momentum flow transport.

Flow Transport Problem Definition. Given empirical observations of two distributions $x_0 \sim \pi_0$ (real data distribution) and $x_1 \sim \pi_1$ (noise distribution) on \mathbb{R}^d , define a forward transport trajectory $T_o: \mathbb{R}^d \to \mathbb{R}^d$ that satisfies $x_1 := T_o(x_0) \sim \pi_1$ when $x_0 \sim \pi_0$. At the same time, the forward flow transport should own the estimable property of the reverse solution trajectory, that is, $x_0 := \tilde{T}_\theta(x_1) \sim \pi_0$ when $x_1 \sim \pi_1$, which requires the trajectory to be continuous and tractable.

Momentum Flow Transport (Discritized-RF). Given the shortest optimal transport $dx_t/dt = v$ (straight-line trajectory) at a constant rate $v = x_1 - x_0$ and a series of discretized anchor points $\{z_1, \cdots, z_{T-1}\}$, define a segmented straight-line trajectory $T_{x_0\mapsto x_1}=\{x_0, z_1, \cdots, z_{T-1}, x_1\}$ that satisfies $dz_t/dt = v_t$ and $v_t = \sqrt{\gamma}v_{t-1} + \sqrt{(1-\gamma)}\epsilon_t, \epsilon_t \sim \mathcal{N}(0, \mathbf{I})$. Meanwhile, the endpoint transport of this momentum flow trajectory are respectively defined as: $T_{x_0\mapsto z_1}: dz_t/dt = v_0$ (initialized by $x_1 - x_0$) and $T_{z_{T-1}\mapsto x_1}: dz_t/dt = \epsilon, \epsilon \sim \mathcal{N}(0, \mathbf{I})$. The momentum flow ensures that the velocity is Gaussian divergent when approaching π_1 , while the velocity is more deterministic and faster when approaching π_0 . Note that the acceleration ϵ follows the same Gaussian distribution $\mathcal{N}(0,\mathbf{I})$ and can therefore be easily estimated by the neural model ϵ_θ to obtain a learnable and tractable inverse transport trajectory $\tilde{T}_{x_1\mapsto x_0;\theta}=\{x_1,z_{T-1;\theta},\cdots,z_{1;\theta},x_{0;\theta}\}$.

Beyond image generation, the challenge of balancing efficiency and diversity is even more pronounced when generating 3D geometric structures, where the data lies on inherently non-Euclidean manifolds instead of the common Euclidean manifold \mathbb{R}^d . Consequently, to further demonstrate the scope and applicability of our momentum flow, we extend it to the **Special Euclidean group** SE(3) for protein backbone generation. In this context, each amino acid residue is represented by a *frame* (i.e., 3D rigid body) in SE(3), parameterizing its spatial orientation and position. This extension is profoundly advantageous: by leveraging the Lie algebra $\mathfrak{se}(3)$, which is the tangent space of SE(3) and linearly isomorphic to \mathbb{R}^6 , we transform the complex nonlinear manifold of protein structures into a vector space where rotations and translations are seamlessly unified within a single stochastic momentum field, achieving diverse and efficient frame sampling without expensive SE(3) geodesic calculations.

The goal of this work is to extend the constant velocity field model to the acceleration field model by learning the momentum flow matching via stochastic velocity field sampling, so as to finally derive a compromise transport path with both speed and diversity. Main contributions are summarized below:

- A momentum-driven flow model for reasonable diversity-efficiency trade-off: Is the straighter the flow, the better? Unlike rectified flows that are modeled on a straight-line trajectory or diffusion models that adopt completely stochastic paths, our momentum flow discretizes the straight path into a series of variable velocity field sub-paths. This makes the trajectory more deterministic (efficient) near the data distribution and more random (diverse) near the noise distribution, thus achieving a proper balance between diversity and efficiency without sacrificing straight-path advantages.
- Segmented straight-line sampling for multi-scale noise modeling: Our proposed Discretized-RF solution trajectory (i.e., segmented flow trajectory where each small segment is a straight line) is a better approximation of multi-scale noise-adding and denoising. It is easier to optimize than the stochastic differential equation (i.e., fluctuation flow trajectory) and can better model multi-scale noise than the constant velocity field differential equation (i.e., rectified flow trajectory).
- Superior performance on multiple image datasets: Extensive experiments show that our Momentum Flow achieves competitive FID and recall scores with substantially fewer denoising steps. Specifically, on multiple image datasets, including CIFAR-10 (Krizhevsky, 2009), CelebA-HQ (Karras et al., 2018) and ImageNet (Deng et al., 2009), it consistently matches or even outperforms the performance of Rectified Flow while requiring only half the number of sampling steps.
- High adaptability to SE(3) for protein generation: By unifying rotations and translations in a single stochastic momentum field via $\mathfrak{se}(3)$ - \mathbb{R}^6 isomorphism, our MFM enables efficient and diverse *frame* sampling, demonstrating its better adaptability to non-Euclidean manifolds than RF.

2 RELATED WORK

Diffusion Models: High Diversity at the Cost of Efficiency. Diffusion models (Song & Ermon, 2019; Ho et al., 2020; Song et al., 2020b; Nichol & Dhariwal, 2021; Kawar et al., 2022; Ma et al., 2024d; Ma et al.; 2025; 2024b) have emerged as a powerful class of generative models, known for their impressive sample diversity. However, their stochastic diffusion trajectories typically require hundreds or thousands of sampling steps, leading to significant computational costs. To overcome this inefficiency, researchers have proposed some optimization methods along two primary directions: sampling acceleration strategies (Liu et al., 2022a; Salimans & Ho, 2022; Gonzalez et al., 2023; Meng et al., 2023; Song et al., 2023; Sauer et al., 2024; Xu et al., 2024) and model architecture improvements (Li et al., 2023; Zhao et al., 2024; Xu et al., 2024; Li et al., 2024a; Ma et al., 2024c). For instance, DDIM (Song et al., 2020a) introduces a non-Markovian reverse process that decouples temporal dependencies, substantially reducing the number of sampling steps. DeepCache (Ma et al., 2024a) accelerates inference by caching and retrieving features across adjacent denoising stages to avoid redundant computations. On the architectural side, some works enhance model efficiency by employing custom multi-decoder U-Net designs that combine time-specific decoders with a shared encoder (Zhang et al., 2024), or by enabling parallel decoder execution to speed up the denoising process (Li et al., 2024b). Moreover, in the field of protein backbone generation, FrameDiff (Yim et al., 2023b) develops a SE(3)-invariant diffusion model on SE(3) N for protein modelling, thereby generating designable, novel and diverse monomers beyond the Protein Data Bank (PDB) (Berman et al., 2000) without relying on a pretrained protein structure prediction network. Despite these advances, diffusion-based models still depend on curved stochastic paths, which remain inherently more expensive to compute than deterministic or straight-path methods. As a result, the fundamental trade-off remains: high sample diversity comes at the expense of computational efficiency.

Rectified Flows: Faster Sampling Meets Less Diversity. Rectified flow models (Liu, 2022; Liu et al., 2023a;b; Wang et al., 2024a; Zhu et al., 2024b; Gat et al., 2024) significantly improve sampling efficiency over diffusion models by optimizing straight-line trajectories in probability space. However, their deterministic and straight sampling paths fundamentally limit their diversity. To solve this issue, various techniques have been proposed to enhance sample diversity while maintaining efficiency. Some methods focus on optimizing noise sampling techniques (Yan et al., 2024; Wang et al., 2024c; Liu et al., 2024), such as training on perceptually relevant noise scales (Esser et al., 2024) or sampling from multi-modal flow directions (Guo & Schwing, 2025). Other efforts aim to improve generation quality (Lee et al., 2024; Li et al., 2024c; Dalva et al., 2024) include applying flow matching in the latent space of pretrained autoencoders (Dao et al., 2023), mitigating numerical errors in the ODE-solving process (Wang et al., 2024b) or introducing posterior-mean-based optimal estimators (Ohayon et al., 2025). Moreover, some protein-related methods (Yim et al., 2023a; Campbell et al., 2024; Lin et al., 2024) utilize RF for protein modelling, achieving speedup during frame sampling. However, the trade-off between sampling speed and diversity persists, motivating the development of adaptive flow-based methods that preserve computational efficiency while enhancing sampling diversity.

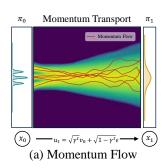
Unlike previous methods, our momentum flow matching model introduces a momentum field into the forward process, where multi-scale noise dynamically adjusts the trajectory directions to promote sampling diversity. To improve computational efficiency, the reverse trajectory is discretized into multiple sub-paths, each optimized via rectified flow. As a result, our model retains the fast sampling speed of rectified flow while recovering much of the sample diversity achieved by diffusion models.

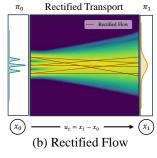
3 Method

In this section, we propose Momentum Flow Transport, a novel flow-based diffusion model family that aims to achieve an effective balance between diversity and efficiency via a brand-new momentum flow matching technique in Sec. 3.1. Momentum Flow is a dynamically compromise approximation of multi-scale noise-adding (or de-noising) between a straight line and a fluctuating line by combinating: 1) fluctuation flow trajectory (close to x_1) for diversity and 2) rectified flow trajectory (close to x_0) for efficiency. We then further introduce the momentum-guided forward process in Sec. 3.2, the acceleration fields-driven reverse process in Sec. 3.3 and its extension to SE(3) in Sec. 3.4.

3.1 Momentum Flow Matching

Optimal Transport (OT). The optimization problem from noise distribution π_1 to data distribution π_0 can be regarded as an optimal transport (OT) problem. Since it is extremely difficult to directly





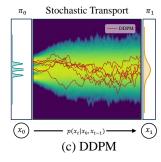


Figure 2: Overview of Momentum Flow. Compared with Rectified Flow (Liu et al., 2023a) (*Efficiency*-OT) and DDPM (Ho et al., 2020) (*Diversity*-OT), the momentum flow tends to explore diversity when close to noise distribution π_1 , and tends to focus on efficiency when close to data distribution π_0 .

solve the trajectory from π_1 to π_0 , recent flow-based methods (Lipman et al.; Liu et al., 2023a) usually first give a tractable forward trajectory T_o to transport any $x_0 \sim \pi_0$ to $x_1 \sim \mathcal{N}(0, \mathbf{I})$ (approximation of π_1), and then solve the posterior $p(\pi_0|\pi_1) = \tilde{T}_{\theta}(\pi_1)$ via a flow-matching trajectory estimator \tilde{T}_{θ} ,

$$\boldsymbol{\pi}_{1} = T_{o}(\boldsymbol{\pi}_{0}) = \int d\boldsymbol{z}_{t} T_{o}\left(\boldsymbol{\pi}_{1} \mid \boldsymbol{z}_{t}\right) \boldsymbol{\pi}\left(\boldsymbol{z}_{t}\right), \ \boldsymbol{\pi}_{0} = \tilde{T}_{\theta}(\boldsymbol{\pi}_{1}) = \int d\boldsymbol{z}_{(0:T)} \boldsymbol{\pi}\left(\boldsymbol{z}_{T}\right) \prod_{t=1}^{T} p\left(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_{t}\right). \tag{1}$$

Stochastic Transport (*Diversity*-OT) and Rectified Transport (*Efficiency*-OT). Stochastic Transport (Ho et al., 2020; Song et al., 2020a) and Rectified Transport (Liu et al., 2023a;b) are two common optimal transport methods, which are respectively known for their high sampling quality (*diversity*) and fast sampling speed (*efficiency*). However, they both struggle with the balance between efficiency and diversity, either relying on overly divergent sampling steps (trajectory variance $\beta_T \to \infty$) or predefined straight trajectories ($\beta_T = 0$). Our work aims to find a balanced trajectory T_o in terms of optimal efficiency and optimal diversity so that the trajectory variance tends to 0 when close to data distribution π_0 (for *efficiency*) and tends to ∞ when close to noise distribution π_1 (for *diversity*).

Momentum Field (Acceleration Field). In order to find a balanced trajectory, we introduce the momentum field. That is a variable velocity field referring to the previous velocity component and the random velocity component in each diffusion step. Let $\boldsymbol{\nu} = \{\boldsymbol{v}_t\}_0^{T-1}$ represent the momentum field (for guiding \boldsymbol{x}_0 to \boldsymbol{x}_1), \boldsymbol{v}_t denote the velocity vector from time t to time t+1, we have:

$$\frac{d\mathbf{z}_t}{dt} = \mathbf{v}_t, \quad \mathbf{v}_t = \begin{cases} \beta(\boldsymbol{\epsilon}_0 - \mathbf{x}_0) & \text{if } t = 0\\ \sqrt{\gamma_t} \mathbf{v}_{t-1} + \sqrt{1 - \gamma_t} \beta \boldsymbol{\epsilon}_t & \text{if } 0 < t < T\\ \beta \boldsymbol{\epsilon}_T & \text{if } t = T. \end{cases}$$
(2)

Here $z_t \sim \pi(z_t)$ is the middle noise-perturbed distribution during the forward diffusion process and $\{\gamma_t\}_1^{T-1}$ is the momentum decay coefficient, which can be chosen as a constant γ (γ < 1) or a positive decreasing series. We choose the former in our work. For convenience, β denotes the normalization coefficient $\beta:=(\sqrt{\gamma}-1)/(\sqrt{\gamma^T}-1)$ and $\{\epsilon_t\}_0^T \sim \mathcal{N}(0,\mathbf{I})$ denotes the standard Gaussian noises. Under the influence of this momentum field, for $\forall \, x_0 \sim \pi_0$ and $x_1 \sim \pi_1$, data x_0 will gradually transform into noise x_1 via the trajectory $T_{x_0\mapsto x_1}=\{x_0,z_1,\cdots,z_{T-1},x_1\}$. Note this momentum field $\{v_t\}_0^{T-1}$ maintains the dynamics of the rectified flow (Liu et al., 2023a) during the initial noise-adding stage with the fastest initial vector $v_0=\beta(\epsilon_0-x_0)$. As the velocity v_0 is gradually noise-perturbed until approaching the noise $\beta\epsilon_T$, we complete the progressively diverse modeling of an OT trajectory $T_{x_0\mapsto x_1}$. Similar to DDPM (Ho et al., 2020), we can directly obtain the momentum v_t at any timestep t via the one-step update formula as (see App. B for details),

$$v_t = \sqrt{\bar{\gamma}_t} v_0 + \sqrt{1 - \bar{\gamma}_t} \beta \epsilon_t, \quad v_0 = \beta (\epsilon_0 - x_0),$$
 (3)

where $\bar{\gamma}_t := \prod_{i=1}^t \gamma_i$. As derived from eq. (3), the proportion of v_0 in v_t decays exponentially with increasing t. This indicates that during the forward process, the momentum v_t gradually deviates from the linear direction defined by v_0 , thereby progressively expanding the exploration diversity.

Momentum Flow Matching Objective. Building upon the flow matching framework for velocity field regression, we optimize the optimal transport (OT) problem by minimizing the MSE between

Algorithm 1: Momentum Flow Transport: Forward Process

- 1: Procedure: $T_o = \text{MomentumField}((\boldsymbol{z}_0, \boldsymbol{z}_T))$: 2: Input: $\boldsymbol{z}_0 \sim \boldsymbol{\pi}_0, \boldsymbol{z}_T = \boldsymbol{\epsilon}_0 \sim \boldsymbol{\pi}_1, T, \{\boldsymbol{\gamma}_t\}_1^{T-1}, \beta, \boldsymbol{v}_0 = \beta(\boldsymbol{\epsilon}_0 \boldsymbol{z}_0)$. 3: For $t \leftarrow 1$ to T do repeat noise disturbance:

- $\begin{aligned} \bullet & \ \ \boldsymbol{z}_t = \boldsymbol{z}_{t-1} + \boldsymbol{v}_{t-1}. \\ \bullet & \ \ \boldsymbol{v}_t = \sqrt{\gamma_t} \boldsymbol{v}_{t-1} + \sqrt{1 \gamma_t} \beta \boldsymbol{\epsilon}_t. \end{aligned}$
- 4: **Return:** Trajectory $T_o = \{z_0, z_1, \dots, z_{T-1}, z_T\}$

predicted momentum and ground-truth. The momentum flow matching objective is formulated as:

$$\mathcal{L}_{\text{MFM}}(\theta) = \mathbb{E}_{t \sim U[0,1]} \| \boldsymbol{u}_{\theta}(\boldsymbol{z}_t, t) - \boldsymbol{v}_t \|^2, \tag{4}$$

where θ denotes learnable parameters for neural network $u_{\theta}(\cdot,t)$, and $t \sim \mathcal{U}[0,1]$. In the inference phase, once the momentum estimate $v_{\theta:t} = u_{\theta}(\cdot, t)$ is obtained, the reverse OT trajectory T_{θ} can be derived, as detailed in the subsequent forward and reverse processes.

3.2 FORWARD PROCESS OF THE MOMENTUM FLOW

Let $z_0=x_0\sim\pi_0$ and $z_T=x_1\sim\pi_1$ respectively denote the data and noise distributions on \mathbb{R}^d . When applying our momentum flow to the forward diffusion process, we can obtain the intermediate noisy distribution $\{z_t \sim \pi(z_t)\}_1^{T-1}$ at discretized anchor points z_t . In general, the forward diffusion process is defined as a Markov chain that progressively injects Gaussian noise ϵ into x_0 over Ttimesteps according to the forward coefficient a_t and b_t , which can be formally expressed as

$$q(\boldsymbol{z}_t|\boldsymbol{z}_{t-1}) = \mathcal{N}(\boldsymbol{z}_t; a_t \boldsymbol{z}_{t-1}, b_t^2 \mathbf{I}), \ q(\boldsymbol{z}_{(1:T)}|\boldsymbol{z}_0) = \prod_{t=1}^T q(\boldsymbol{z}_t|\boldsymbol{z}_{t-1}).$$
 (5)

We subsequently introduce our momentum field into the classical diffusion process to adjust the balance between optimal diversity and optimal efficiency for the exploration of the forward compromise trajectory T_o , and the detailed process is presented in Algorithm 1.

Forward Momentum Flow. From eq. (3), we can observe that the recursive formulation of our momentum field shares the similar form as that in DDPM (Ho et al., 2020), allowing us to directly obtain the prior probability distribution $q(\mathbf{v}_t|\mathbf{v}_{t-1})$ of the momentum flow:

$$q(\mathbf{v}_t|\mathbf{v}_{t-1}) = \mathcal{N}(\mathbf{v}_t; \sqrt{\alpha_t}\mathbf{v}_{t-1}, (1-\alpha_t)\beta^2 \mathbf{I}), \ q(\mathbf{v}_t|\mathbf{v}_0) = \mathcal{N}(\mathbf{v}_t; \sqrt{\bar{\alpha_t}}\mathbf{v}_0, (1-\bar{\alpha_t})\beta^2 \mathbf{I}), \quad (6)$$

where $\alpha_t := \gamma$ and $\bar{\alpha_t} := \prod_{i=1}^t \alpha_i = \gamma^t$ ($\gamma < 1$ is a fixed constant). Notably, the formal alignment between the momentum flow and the forward process in DDPM (Ho et al., 2020) also allows a straightforward derivation of the posterior distribution $p_{\theta}(v_{t-1}|v_t)$ of momentum flow, see Sec. 3.3.

Forward Data Flow. Based on the above momentum flow, we can further build the forward trajectory (i.e., data flow) $T_o = \{z_0, z_1, \dots, z_{T-1}, z_T\}$, which is represented in the form of a conditional probability distribution $q(z_t|z_{t-1})$. According to eq. (2) and eq. (3), the one-step forward data distribution can be obtained (see App. C for a detailed derivation) as

$$q(\boldsymbol{z}_t|\boldsymbol{z}_0) = \mathcal{N}\left(\boldsymbol{z}_t; (1 - (\frac{\sqrt{\gamma^t} - 1}{\sqrt{\gamma} - 1})\beta)\boldsymbol{z}_0, ((\frac{\sqrt{\gamma^t} - 1}{\sqrt{\gamma} - 1})^2 - \frac{\gamma^t - 1}{\gamma - 1} + t)\beta^2 \mathbf{I}\right).$$
(7)

Due to $\beta := (\sqrt{\gamma} - 1)/(\sqrt{\gamma^T} - 1)$, when computing the noise distribution $\pi(z_T)$, we can eliminate the complex coefficient in front of z_0 in eq. (7) to derive a zero-mean Gaussian distribution (independent of the data distribution π_0), which can be formally expressed as

$$q(\boldsymbol{z}_T|\boldsymbol{z}_0) = \mathcal{N}\left(\boldsymbol{z}_T; 0, \left(\left(\frac{\sqrt{\gamma^T} - 1}{\sqrt{\gamma} - 1}\right)^2 - \frac{\gamma^T - 1}{\gamma - 1} + T\right)\beta^2 \mathbf{I}\right),\tag{8}$$

This simplified formula facilitates the subsequent reverse momentum transport process and significantly reduces computational complexity during training and inference.

Algorithm 2: Momentum Flow Matching: Reverse Process

- 1: **Procedure**: $\tilde{T}_{\theta} = \text{MomentumFlow}((\boldsymbol{z}_0, \boldsymbol{z}_T))$:
 2: **Input:** Momentum model $\boldsymbol{u}_{\theta}(\cdot, t) : \mathbb{R}^d \times [0, 1] \to \mathbb{R}^d$ with parameters θ .
 3: **Training:** $\hat{\theta} = \arg\min_{\theta} \sum_{t=1}^T \mathbb{E}\left[\|\boldsymbol{u}_{\theta}(m\boldsymbol{z}_t + (1-m)\boldsymbol{z}_{t-1}, m) (\boldsymbol{z}_t \boldsymbol{z}_{t-1})\|^2\right]$, with $m \sim \mathcal{U}[0, 1]$.
- 4: For $t \leftarrow T$ to 1 do repeat sampling:

270

271

272 273 274

275

276

277 278 279

281

283

284

287

289

290

291

292

293

299 300 301

302

303

304

305

306

307

308

310 311

312

313 314

315

316 317

318

319

320

321

322

323

- Draw $(\boldsymbol{z}_{t-1}, \boldsymbol{z}_t)$ from $\boldsymbol{\pi}(\boldsymbol{z}_{t-1}) \times \boldsymbol{\pi}(\boldsymbol{z}_t)$, with $\boldsymbol{z}_{t-1} \sim \boldsymbol{\pi}(\boldsymbol{z}_{t-1})$ and $\boldsymbol{z}_t \sim \boldsymbol{\pi}(\boldsymbol{z}_t)$.
 Solve ODE: $\frac{d\boldsymbol{z}_t}{dt} = \boldsymbol{u}_{\theta}(\boldsymbol{z}_t^m, m)$, with $\boldsymbol{z}_0 \sim \boldsymbol{\pi}_0$.
 Return: Sub-trajectory $\boldsymbol{z}_t = \{\boldsymbol{z}_t^m : m \in [0, 1]\}$.

- 5: **Return:** Trajectory $T_{\theta} = \{z_t : t \in [0, 1]\}.$

REVERSE PROCESS OF THE MOMENTUM FLOW

The reverse process of the momentum flow aims to restore noise distribution π_1 to data distribution π_0 via an inverse trajectory $T_{\theta} = \{z_T, z_{T-1;\theta}, \cdots, z_{1;\theta}, z_{0;\theta}\}$, which is estimated by a neural network for approximating $z_{t;\theta} \sim \pi(z_{t;\theta})$. To achieve this, we can approximate the momentum field $\{v_t\}_0^{T-1}$ and then utilize the relationship $z_{t-1;\theta} = z_{t;\theta} - v_{t-1;\theta}$ to estimate $z_{t-1;\theta}$ from $z_{t;\theta}$ and $v_{t-1:\theta}$, as illustrated in Algorithm 2. We denote the estimated values of (z_t, v_t) as $(z_{t:\theta}, v_{t:\theta})$. Based on this framework, we discuss two ways to approximate the momentum field. The first way is to approximate v_t by estimating $p_{\theta}(v_{t-1}|v_t)$. Benefiting from the formal similarity between the forward momentum flow and the DDPM formulation (Ho et al., 2020), we can derive directly the corresponding posterior distribution and estimate $v_{t-1:\theta}$ from $v_{t:\theta}$ by training a noise predictor ϵ_{θ} :

$$p_{\theta}(\boldsymbol{v}_{t-1}|\boldsymbol{v}_{t}) = \mathcal{N}\left(\frac{\sqrt{\gamma}\left(1-\gamma^{t-1}\right)+\sqrt{\gamma^{t-1}}(1-\gamma)}{1-\gamma^{t}}\boldsymbol{v}_{t} - \frac{(1-\gamma)\beta\epsilon_{\theta}}{\sqrt{\gamma(1-\gamma^{t})}}, \frac{(1-\gamma)\left(1-\gamma^{t-1}\right)}{1-\gamma^{t}}\beta^{2}\mathbf{I}\right),$$
(9)

$$\boldsymbol{v}_{t-1;\theta} = \frac{\sqrt{\gamma} \left(1 - \gamma^{t-1}\right) + \sqrt{\gamma^{t-1}} (1 - \gamma)}{1 - \gamma^t} \boldsymbol{v}_{t;\theta} - \frac{(1 - \gamma)\beta \boldsymbol{\epsilon}_{\theta}}{\sqrt{\gamma(1 - \gamma^t)}} + \sqrt{\frac{(1 - \gamma)(1 - \gamma^{t-1})}{1 - \gamma^t}} \beta \boldsymbol{\epsilon}. \quad (10)$$

The second method directly approximates v_t by employing rectified flow on each sub-path. Specifically, between each adjacent intermediate noise-perturbed distribution pair $(\pi(z_t), \pi(z_{t-1}))$ at the discretized anchor point pair (z_t, z_{t-1}) , we insert M intermediate points $z_{t-1}^{(m)}$ via linear interpolation:

$$\mathbf{z}_{t-1}^{(m)} = m\mathbf{z}_t + (1-m)\mathbf{z}_{t-1},$$
 (11)

where $m \sim \mathcal{U}[0,1]$. To enhance sampling efficiency, we apply rectified flow to formulate a straight path for each sub-path $\{\pi(z_t) \to \pi(z_{t-1})\}_1^T$, with the network $u_{\theta}(\cdot,t)$ trained to match the corresponding velocity $v_{t-1} = z_t - z_{t-1}$. Therefore, the original objective, i.e., eq. (4), is reformulated into the following optimization objective:

$$\mathcal{L}_{MFM}(\theta) = \sum_{t=1}^{T} \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\| \boldsymbol{u}_{\theta}(m\boldsymbol{z}_{t} + (1-m)\boldsymbol{z}_{t-1}, m) - (\boldsymbol{z}_{t} - \boldsymbol{z}_{t-1}) \|^{2} \right].$$
 (12)

The second method achieves much higher computational efficiency than DDPM by using rectified flow to optimize the trajectory. Therefore, we follow the second method in all experiments.

3.4 MOMENTUM FLOW MATCHING ON SE(3)

We now describe the extension of our MFM to protein backbone generation. The backbone atom positions of each residue in a protein backbone are parameterized by a rigid transformation $T \in SE(3)$. Each frame T=(r,x) consists of a rotation matrix $r \in SO(3)$ and a translation vector $x \in \mathbb{R}^3$. A protein backbone consists of N residues meaning it can be parameterized as $\mathbf{T} = [T^{(1)}, \dots, T^{(N)}]$ with $T \in SE(3)^N$. For notational simplicity, our extension focuses on a single frame but applies to all frames in a backbone since $SE(3)^N$ is a product space and we use an additive metric over frames.

Different from previous methods where noise is directly superimposed on $T \in SE(3)$, we introduce noise on its tangent space $\mathfrak{se}(3)$ to characterize momentum. Specifically, the Lie algebra $\mathfrak{se}(3)$ consists of all infinitesimal generators of rigid body motions and can be formally represented as:

$$\mathfrak{se}(3) = \left\{ \begin{pmatrix} [\boldsymbol{\omega}]_{\times} & \boldsymbol{v} \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{4\times4} \mid \boldsymbol{\omega}, \boldsymbol{v} \in \mathbb{R}^{3} \right\},$$
 (13)

where $[\omega]_{\times}$ denotes the rotation generator corresponding to the angular velocity ω and v denotes the translation generator, i.e., linear velocity. Thus, each element in $\mathfrak{se}(3)$ is uniquely determined by 6 parameters $(\omega, v) \in \mathbb{R}^6$. Moreover, benefiting from the linear isomorphism $\mathfrak{se}(3) \cong \mathbb{R}^6$, calculations can be simplified smoothly from the complex nonlinear manifold SE(3) to the vector space \mathbb{R}^6 .

4 EXPERIMENTS

In this section, we conduct experiments using the rectified flow framework implemented in PyTorch to evaluate the image generation diversity and efficiency of the proposed momentum flow model. The primary objectives are to compare the generating performance between momentum flow and rectified flow, and to analyze the impact of the momentum field on the diversity and speed of the generative process. The results show that momentum flow retains the fast sampling capability of straight velocity fields. In addition, by injecting multi-scale noise through the momentum fields, the diversity and the quality of the generated images are significantly enhanced.

4.1 Unconditioned Image Generation

Experiment Settings. We build upon the official open-source implementation as the foundation of our model framework, and all experiments are conducted as illustrated in Table 4. To maximize performance within our computational budget, we conduct a grid search over learning rates and weight decay parameters. For evaluation, we generate 50,000 samples from each model and evaluate generating quality and diversity using the Fréchet Inception Distance (FID) (Seitzer, 2020) and the recall value (Sajjadi et al., 2018). As 'recall' is defined as the coverage rate of generated samples over the real data distribution, we evaluate the diversity of generated samples by calculating the Recall value. Additional training and implementation details are provided in App. E.

Comparison on CIFAR-10. We report unconditional image generation results on the CIFAR-10 dataset (Krizhevsky, 2009). We train all models for 20,000 steps with a batch size of 1,024. In Table 1, the FID-50K scores are obtained using 50-NFE sampling. All entries employ the same U-Net architecture, applied directly in the pixel space. On this dataset, our method demonstrates a clear advantage over prior approaches.

Comparison on CelebA-HQ and ImageNet. As shown in Table 2, the momentum flow model consistently achieves superior performance compared to the rectified flow model, as evidenced by significantly lower FID and higher recall values across various settings. The improvements on CelebA-HQ dataset are significant, achieving an average improvement of over 11 FID points and 0.06 recall values. In addition, when reducing the number of function evaluations (NFE) from 100 to 10, the performance degradation is minimal, and the model remains competitive with rectified flow under the same sampling budget, highlighting the efficiency of our method

Table 1: Quantitative results on the CIFAR-10 dataset, while $\gamma=0.99$ in our Momentum Flow.

| method | NFE | $\text{FID}{\downarrow}$ |
|-----------------------------------|-----|--------------------------|
| RectifiedFlow (Liu et al., 2022b) | 50 | 50.26 |
| NanoFlow (Zhu et al., 2024a) | 50 | 47.40 |
| MomentumFlow (ours) | 50 | 45.66 |

Table 2: Quantitative results on CelebA-HQ and ImageNet-64 datasets. Here \widehat{Step} is the number of denoising steps in each sub-path ($\gamma = 0.98$).

| N | \widehat{Step} | CelebA-HQ | | | ImageNet-64 | | |
|--------------------|------------------|-----------|-------|----------|-------------|-------|----------|
| | | FID ↓ | NFE ↓ | Recall ↑ | FID ↓ | NFE ↓ | Recall ↑ |
| 1 (Rectified Flow) | 10 | 98.98 | 10 | 0.268 | 61.48 | 10 | 0.366 |
| | 50 | 65.38 | 50 | 0.384 | 42.42 | 50 | 0.457 |
| | 100 | 58.90 | 100 | 0.445 | 41.83 | 100 | 0.451 |
| 2 (Ours) | 5 | 84.61 | 10 | 0.345 | 60.34 | 10 | 0.368 |
| | 25 | 54.07 | 50 | 0.457 | 41.99 | 50 | 0.454 |
| | 50 | 50.57 | 100 | 0.488 | 41.77 | 100 | 0.459 |
| 5 (Ours) | 2 | 110.72 | 10 | 0.177 | 104.06 | 10 | 0.258 |
| | 10 | 99.57 | 50 | 0.249 | 96.70 | 50 | 0.294 |
| | 20 | 94.61 | 100 | 0.261 | 94.87 | 100 | 0.294 |

grounded in balanced transport (OT). These results indicate that the momentum flow model preserves the fast sampling efficiency of rectified flow while generating higher-quality images.

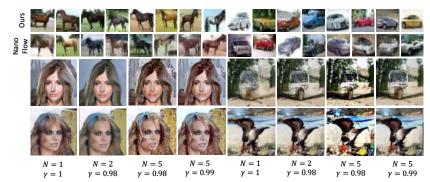


Figure 3: Samples of different datasets. Top: Samples on the CIFAR-10 dataset. Bottom: Momentum Flow samples with varying N and γ , shown for CelebA-HQ (left) and ImageNet-64 (right).

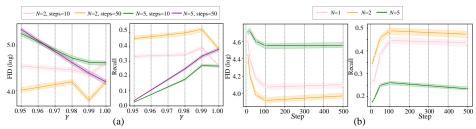


Figure 4: (a) shows the relationship between model performance (FID, Recall) and the γ setting. (b) illustrates the relationship between denoising steps and the quality of the generated images.

As shown in Figure 3, the deterministic nature of straight-line modeling in rectified flow leads to noticeable distortions in local details (e.g., mouth, eyes, and accessories). In contrast, the momentum flow model employs momentum-guided trajectories to explore a broader space, resulting in significantly improved detail generation. By dynamically injecting controllable velocity deviations via the momentum field, our method enhances both generating diversity and fidelity on high-resolution datasets such as CelebA-HQ, highlighting the effectiveness of multi-scale noise in guiding generation.

Acceleration Process of Momentum Flow. We empirically evaluate the efficiency of momentum flow in image generation. Although additional noise is injected into the velocity field, the linear straight structure ensures a constant velocity within each sub-path. This design preserves the efficiency of the original rectified flow. We compare the visual quality of generated images under different total denoising step settings: 10 and 50 steps while more examples can be found in the appendix. Momentum flow achieves comparable or even superior results to rectified flow with only half the number of sampling steps, as highlighted in the red-marked values in Table 2. These results demonstrate momentum flow inherits the acceleration advantages of rectified flows while further benefiting from enhanced flexibility.

Sampling Efficiency and Generating Diversity. Under a fixed image input, we compare the number of sampling steps under the same network architecture to assess their sampling efficiency. The Momentum Flow model achieves superior performance within the same time budget and requires fewer steps to reach comparable results, demonstrating its strong sampling efficiency. As shown in Figure 3, our method exhibits significant advantages in object color, shape, and background, and by adjusting N and γ , Momentum Flow can produce more diverse colors and finer rendering details. These results confirm that our model effectively balances efficiency and diversity in the generative process, as visually shown in Figure 4.

Momentum Decay Coefficient. In our model framework, the decay coefficient γ controls the level of noise perturbation by modulating the influence of the momentum flow, thereby enabling dynamic refinement of the forward trajectory. We observe that decreasing γ causes the momentum flow to deviate more rapidly from the initial momentum direction v_0 , which expands the exploration region and enhances sample diversity. However, excessive decay in the early stages may weaken the guidance from the initial momentum v_0 , so γ should not be set too low. As shown in Figure 4, when the number of noise-injecting steps is set to N=2, a decay coefficient of $\gamma=0.99$ results in significantly lower FID scores and higher recall values compared to $\gamma=0.999$ and $\gamma=0.98$. In contrast, when N=5, the best performance is achieved with $\gamma=1$, indicating that a larger number

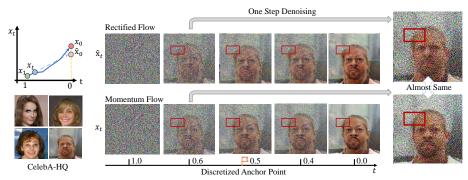


Figure 5: Reverse trajectories of momentum flow and rectified flow at different denoising steps, where the optimized velocity field of momentum flow improves image details.

of forward steps N requires a slower decay (i.e., a γ closer to 1) to maintain effective guidance of v_0 . These experimental results suggest that appropriately selecting the decay coefficient γ and the number of velocity steps N can substantially improve both quality and diversity of generated images.

Broader Analysis. We compare the sampling processes of rectified flow and momentum flow to assess the advantages of our method in the denoising trajectory. In the early stages of generation—specifically the first few sampling steps—the results of both models appear similar. To illustrate this, we select z_t at t=0.4 and perform a single denoising step using the predicted velocity field to obtain a reference image, as shown on the right side of Figure 5. While the early-stage images generated by both methods show no notable differences, particularly in regions such as the hair, momentum flow exhibits a clear advantage in detail fidelity after passing the anchor point (t=0.5).

Ignoring the refinement methods like distillation (Lee et al., 2024; Zhao et al., 2024), the initial velocity field often struggles to effectively bridge the gap between the noise and data distributions due to its reliance on fixed straight-line trajectories. The strength of our method lies in its ability to encourage broader exploration of the data space. By introducing momentum-guided velocity deviations, the model is not constrained to a fixed straight-line trajectory. Instead, it gains the flexibility to adjust its path dynamically.

4.2 RESULTS OF PROTEIN BACKBONE GENERATION

To verify the effectiveness of Momentum Flow on the protein monomer generation task, following GENIE (Lin & AlQuraishi, 2023) and FrameFlow (Yim et al., 2023a), we train it on the SCOPe dataset (Chandonia et al., 2022) with proteins below length 128 for a total of 3,938

Table 3: Protein backbone generation results.

| Model | Sampling | Accuracy N | Confidence Metrics | | |
|-----------|----------|----------------------------|--------------------|--------|-------|
| Wiodei | | scTM (> 0.5) \uparrow | scRMSD↓ | pLDDT↑ | pAE↓ |
| GENIE | SDE | 0.09 (0.0) | 27.97 | 55.03 | 19.65 |
| FrameFlow | ODE | 0.39 (0.15) | 9.92 | 59.09 | 12.64 |
| Ours | MFM | 0.47 (0.45) | 8.05 | 70.09 | 9.50 |

examples. During evaluation, we sample 10 backbones for every length between 60 and 300¹ then use ProteinMPNN (Dauparas et al., 2022) to design 8 sequences for each backbone. We then evaluate the quality of generated proteins based on four metrics: scTM, scRMSD, pLDDT, and pAE. The quantitative results are reported in Table 3. See App. H for complete settings and detailed analysis.

5 Conclusions

Discretized-RF proposes a compromise transport method to balance the trade-off between diversity and efficiency. By injecting multi-scale noise perturbations based on momentum flow and formulating discretized straight-line trajectories, our approach effectively optimizes two key limitations of previous models: restricted generating diversity and high computational cost. Extensive experiments show that the momentum flow model achieves both high-quality image generation and fast sampling speed. Looking ahead, we believe the Discretized-RF framework offers a promising direction for designing more flexible flow trajectories and further exploring the diversity-efficiency optimal method.

¹The upper limit of 300 here differs from the upper limit of 128 during training. We increase the upper limit during evaluation to demonstrate the generalization of our model in generating long sequence proteins.

REFERENCES

- Kingma DP Ba J Adam et al. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 1412(6), 2014.
- Grigory Bartosh, Dmitry P Vetrov, and Christian Andersson Naesseth. Neural flow diffusion models:
 Learnable forward process for improved diffusion modelling. *Advances in Neural Information Processing Systems*, 37:73952–73985, 2024.
 - Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic acids research*, 28(1): 235–242, 2000.
 - Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. arXiv preprint arXiv:2402.04997, 2024.
 - John-Marc Chandonia, Lindsey Guan, Shiangyi Lin, Changhua Yu, Naomi K Fox, and Steven E Brenner. Scope: improvements to the structural classification of proteins—extended database to facilitate variant interpretation and machine learning. *Nucleic acids research*, 50(D1):D553–D559, 2022.
 - Yusuf Dalva, Kavana Venkatesh, and Pinar Yanardag. Fluxspace: Disentangled semantic editing in rectified flow transformers, 2024. URL https://arxiv.org/abs/2412.09611.
 - Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space, 2023. URL https://arxiv.org/abs/2307.08698.
 - Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning–based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
 - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
 - RA Engh and R Huber. Structure quality and target parameters. *International Tables for Crystallog-raphy*, 2012.
 - Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In Forty-first International Conference on Machine Learning, 2024. URL https://openreview.net/forum?id=FPnUhsQJ5B.
 - Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky T. Q. Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 133345–133385. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/f0d629a734b56a642701bba7bc8bb3ed-Paper-Conference.pdf.
 - Martin Gonzalez, Nelson Fernandez Pinto, Thuy Tran, elies Gherbi, Hatem Hajri, and Nader Masmoudi. Seeds: Exponential sde solvers for fast high-quality sampling from diffusion models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 68061–68120. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/d6f764aae383d9ff28a0f89f71defbd9-Paper-Conference.pdf.
 - Pengsheng Guo and Alexander G. Schwing. Variational rectified flow matching, 2025. URL https://arxiv.org/abs/2502.09616.
 - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
 - Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. URL https://arxiv.org/abs/1710.10196.
 - Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *arXiv preprint arXiv:2201.11793*, 2022.
 - Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical Report TR-2009, University of Toronto, 2009.
 - Adam Leach, Sebastian M Schmon, Matteo T Degiacomi, and Chris G Willcocks. Denoising diffusion probabilistic models on so (3) for rotational alignment. In *ICLR 2022 workshop on geometrical and topological representation learning*, 2022.
 - Sangyun Lee, Zinan Lin, and Giulia Fanti. Improving the training of rectified flows. *Advances in Neural Information Processing Systems*, 37:63082–63109, 2024.
 - Muyang Li, Tianle Cai, Jiaxin Cao, Qinsheng Zhang, Han Cai, Junjie Bai, Yangqing Jia, Kai Li, and Song Han. Distrifusion: Distributed parallel inference for high-resolution diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7183–7193, 2024a.
 - Senmao Li, Taihang Hu, Joost van de Weijer, Fahad Shahbaz Khan, Tao Liu, Linxuan Li, Shiqi Yang, Yaxing Wang, Ming-Ming Cheng, and Jian Yang. Faster diffusion: Rethinking the role of the encoder for diffusion model inference. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 85203–85240. Curran Associates, Inc., 2024b. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/9ad996b5c45130de2bc00b60d8607904-Paper-Conference.pdf.
 - Shufan Li, Konstantinos Kallidromitis, Akash Gokul, Zichun Liao, Yusuke Kato, Kazuki Kozuka, and Aditya Grover. Omniflow: Any-to-any generation with multi-modal rectified flows. *arXiv* preprint arXiv:2412.01169, 2024c.
 - Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snapfusion: Text-to-image diffusion model on mobile devices within two seconds. *Advances in Neural Information Processing Systems*, 36:20662–20678, 2023.
 - Haitao Lin, Odin Zhang, Huifeng Zhao, Dejun Jiang, Lirong Wu, Zicheng Liu, Yufei Huang, and Stan Z Li. Ppflow: Target-aware peptide design with torsional flow matching. *arXiv* preprint *arXiv*:2405.06642, 2024.
 - Yeqing Lin and Mohammed AlQuraishi. Generating novel, designable, and diverse protein structures by equivariantly diffusing oriented residue clouds. *arXiv* preprint arXiv:2301.12485, 2023.
 - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*.
 - Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds, 2022a. URL https://arxiv.org/abs/2202.09778.
 - Peng Liu, Dongyang Dai, and Zhiyong Wu. Rfwave: Multi-band rectified flow for audio waveform reconstruction, 2024. URL https://arxiv.org/abs/2403.05010.
 - Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport, 2022. URL https://arxiv.org/abs/2209.14577.
 - Xingchao Liu, Chengyue Gong, et al. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2022b.

- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023a.
 - Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. Instaflow: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023b.
 - Ao Luo, Xin Li, Fan Yang, Jiangyu Liu, Haoqiang Fan, and Shuaicheng Liu. Flowdiffuser: Advancing optical flow estimation with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19167–19176, 2024.
 - Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15762–15772, June 2024a.
 - Zhiyuan Ma, Guoli Jia, Biqing Qi, and Bowen Zhou. Safe-sd: Safe and traceable stable diffusion with text prompt trigger for invisible generative watermarking. In *ACM Multimedia* 2024.
 - Zhiyuan Ma, Guoli Jia, and Bowen Zhou. Adapedit: Spatio-temporal guided adaptive editing algorithm for text-based continuity-sensitive image editing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4154–4161, 2024b.
 - Zhiyuan Ma, Zhihuan Yu, Jianjun Li, and Bowen Zhou. Lmd: faster image reconstruction with latent masking diffusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4145–4153, 2024c.
 - Zhiyuan Ma, Liangliang Zhao, Biqing Qi, and Bowen Zhou. Neural residual diffusion models for deep scalable vision generation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024d.
 - Zhiyuan Ma, Yuzhu Zhang, Guoli Jia, Liangliang Zhao, Yichao Ma, Mingjie Ma, Gaofeng Liu, Kaiyan Zhang, Ning Ding, Jianjun Li, et al. Efficient diffusion models: A comprehensive survey from principles to practices. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
 - Siegfried Matthies, J Muller, and GW Vinel. On the normal distribution in the orientation space. *Texture, Stress, and Microstructure*, 10(1):77–96, 1988.
 - Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14297–14306, June 2023.
 - Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
 - Guy Ohayon, Tomer Michaeli, and Michael Elad. Posterior-mean rectified flow: Towards minimum mse photo-realistic image restoration, 2025. URL https://arxiv.org/abs/2410.00418.
 - Frank C Park and Roger W Brockett. Kinematic dexterity of robotic mechanisms. *The International Journal of Robotics Research*, 13(1):1–15, 1994.
 - Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. In *Advances in Neural Information Processing Systems* (*NeurIPS*), 2018.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models, 2022. URL https://arxiv.org/abs/2202.00512.
 - Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pp. 87–103. Springer, 2024.

- Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. https://github.com/mseitzer/pytorch-fid, August 2020. Version 0.3.0.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv* preprint arXiv:2010.02502, 2020a.
 - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
 - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint *arXiv*:2011.13456, 2020b.
 - Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
 - Brian L Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and Tommi Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. *arXiv* preprint arXiv:2206.04119, 2022.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
 - Fu-Yun Wang, Ling Yang, Zhaoyang Huang, Mengdi Wang, and Hongsheng Li. Rectified diffusion: Straightness is not your need in rectified flow, 2024a. URL https://arxiv.org/abs/2410.07303.
 - Jiangshan Wang, Junfu Pu, Zhongang Qi, Jiayi Guo, Yue Ma, Nisha Huang, Yuxin Chen, Xiu Li, and Ying Shan. Taming rectified flow for inversion and editing, 2024b. URL https://arxiv.org/abs/2411.04746.
 - Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. Frieren: Efficient video-to-audio generation with rectified flow matching. *arXiv e-prints*, pp. arXiv–2406, 2024c.
 - Yanwu Xu, Yang Zhao, Zhisheng Xiao, and Tingbo Hou. Ufogen: You forward once large scale text-to-image generation via diffusion gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8196–8206, 2024.
 - Hanshu Yan, Xingchao Liu, Jiachun Pan, Jun Hao Liew, Qiang Liu, and Jiashi Feng. Perflow: Piecewise rectified flow as universal plug-and-play accelerator. *arXiv preprint arXiv:2405.07510*, 2024.
 - Jason Yim, Andrew Campbell, Andrew YK Foong, Michael Gastegger, José Jiménez-Luna, Sarah Lewis, Victor Garcia Satorras, Bastiaan S Veeling, Regina Barzilay, Tommi Jaakkola, et al. Fast protein backbone generation with se (3) flow matching. *arXiv preprint arXiv:2310.05297*, 2023a.
 - Jason Yim, Brian L Trippe, Valentin De Bortoli, Emile Mathieu, Arnaud Doucet, Regina Barzilay, and Tommi Jaakkola. Se (3) diffusion model with application to protein backbone generation. *arXiv preprint arXiv:2302.02277*, 2023b.
 - Huijie Zhang, Yifu Lu, Ismail Alkhouri, Saiprasad Ravishankar, Dogyoon Song, and Qing Qu. Improving training efficiency of diffusion models via multi-stage framework and tailored multi-decoder architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7372–7381, June 2024.
 - Yang Zhao, Yanwu Xu, Zhisheng Xiao, Haolin Jia, and Tingbo Hou. Mobilediffusion: Instant text-to-image generation on mobile devices. In *European Conference on Computer Vision*, pp. 225–242. Springer, 2024.
 - Kan Zhu, Yilong Zhao, Liangyu Zhao, Gefei Zuo, Yile Gu, Dedong Xie, Yufei Gao, Qinyu Xu, Tian Tang, Zihao Ye, Keisuke Kamahori, Chien-Yu Lin, Stephanie Wang, Arvind Krishnamurthy, and Baris Kasikci. Nanoflow: Towards optimal large language model serving throughput, 2024a. URL https://arxiv.org/abs/2408.12757.

Yixuan Zhu, Wenliang Zhao, Ao Li, Yansong Tang, Jie Zhou, and Jiwen Lu. Flowie: Efficient image enhancement via rectified flow. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 13-22, June 2024b.

A ORGANIZATION OF THE SUPPLEMENTARY

In this supplementary, we first provide a detailed proof of the one-step momentum update in App. B. In App. C, we derive the one-step forward data distribution. In App. D, we conduct a toy experiment to illustrate the theoretical background of our momentum flow. In App. E, we present additional image experiments and more complete ablation studies. In particular, we recall in App. F some important theoretical preliminaries about SO(3) and SE(3). Using these, we introduce in App. G our protein backbone parameterization, the conversion between coordinates and frames, and the architecture of FramePred. In App. H, we show more detailed analysis of protein experiments.

B PROOF OF THE ONE-STEP MOMENTUM UPDATE

In this section, we provide a detailed proof of the one-step update based on the recursive formulation (refer to eq. (2)) of the momentum field $\{v_t\}_0^{T-1}$. The specific proof process is as follows:

$$v_{t} = \sqrt{\gamma_{t}}v_{t-1} + \sqrt{1 - \gamma_{t}}\beta\epsilon_{t}$$

$$= \sqrt{\gamma_{t}\gamma_{t-1}}v_{t-2} + \sqrt{\gamma_{t} - \gamma_{t}\gamma_{t-1}}\beta\epsilon_{t} + \sqrt{1 - \gamma_{t}}\beta\epsilon_{t-1}$$

$$= \sqrt{\gamma_{t}\gamma_{t-1}}v_{t-2} + \sqrt{1 - \gamma_{t}\gamma_{t-1}}\beta\epsilon_{t}$$

$$= \cdots$$

$$= \sqrt{\gamma_{t}\gamma_{t-1}}\cdots\gamma_{1}v_{0} + \sqrt{1 - \gamma_{t}\gamma_{t-1}}\cdots\gamma_{1}\beta\epsilon_{t}$$

$$= \sqrt{\overline{\gamma_{t}}}v_{0} + \sqrt{1 - \overline{\gamma_{t}}}\beta\epsilon_{t}$$

$$= \sqrt{\overline{\gamma_{t}}}v_{0} + \sqrt{1 - \overline{\gamma_{t}}}\beta\epsilon_{t}$$
(14)

where $\bar{\gamma_t} := \prod_{i=1}^t \gamma_i$.

C PROOF OF THE ONE-STEP FORWARD DATA DISTRIBUTION

In the forward process, the data distribution evolves under the momentum field $\{v_t\}_0^{T-1}$. According to the one-step momentum update formula (refer to eq. (3)), the forward data distribution $\pi(z_t)$ originates from the π_0 , perturbed by a exponentially scaled decaying contribution of the initial momentum v_0 , along with random noise ϵ_t . The specific derivation process is as follows:

$$z_{t} = z_{t-1} + v_{t-1}$$

$$= z_{0} + v_{0} + v_{1} + v_{2} + \dots + v_{t-1}$$

$$= z_{0} + v_{0} \left(1 + \sqrt{\gamma} + \sqrt{\gamma^{2}} + \dots + \sqrt{\gamma^{t-1}}\right) + \sum_{i=1}^{t-1} \sqrt{1 - \gamma^{i}} \beta \epsilon_{i}$$

$$= z_{0} + v_{0} \left(\frac{\sqrt{\gamma^{t}} - 1}{\sqrt{\gamma} - 1}\right) + \sqrt{t - (1 + \gamma + \gamma^{2} + \dots + \gamma^{t-1})} \beta \epsilon_{t}$$

$$= z_{0} + (\epsilon_{0} - z_{0}) \left(\frac{\sqrt{\gamma^{t}} - 1}{\sqrt{\gamma} - 1}\right) \beta + \sqrt{t - \frac{\gamma^{t} - 1}{\gamma - 1}} \beta \epsilon_{t}$$

$$= \left(1 - \left(\frac{\sqrt{\gamma^{t}} - 1}{\sqrt{\gamma} - 1}\right) \beta\right) z_{0} + \sqrt{\left(\frac{\sqrt{\gamma^{t}} - 1}{\sqrt{\gamma} - 1}\right)^{2} - \frac{\gamma^{t} - 1}{\gamma - 1} + t} \beta \epsilon_{t}.$$
(15)

Thus we have

$$q(\boldsymbol{z}_t|\boldsymbol{z}_0) = \mathcal{N}\left(\boldsymbol{z}_t; (1 - (\frac{\sqrt{\gamma^t} - 1}{\sqrt{\gamma} - 1})\beta)\boldsymbol{z}_0, ((\frac{\sqrt{\gamma^t} - 1}{\sqrt{\gamma} - 1})^2 - \frac{\gamma^t - 1}{\gamma - 1} + t)\beta^2 \mathbf{I}\right).$$
(16)

D MFM TOY EXPERIMENT

D.1 TOY EXPERIMENT PARAMETERIZATION

To illustrate the theoretical background of our momentum flow, we provide an example in Figure 6, demonstrating the expected momentum flow in the forward process and the optimal transport path in

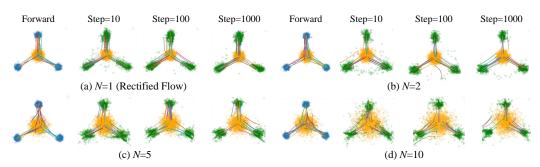


Figure 6: Forward and reverse trajectories of our momentum flow with different numbers N of discretized anchor points. In our settings, blue points are sampled from π_0 , orange points are sampled from π_1 , green points denote generated samples, and lines represent transport trajectories.

the reverse process. Momentum flow is simulated utilizing the Euler method with a constant step size of 1/N, computed at N discrete anchor points, where N denotes the number of such points, corresponding to the value of T in Momentum Flow. The number of sampling points on each segment is defined as \widehat{Step} , and the total number of sampling steps is defined as $\widehat{Step} \times N$. Note that in all experiments in this section, this notation is used by default. Moreover, we define the use of a fully connected neural network with two hidden layers to estimate the momentum field. In practice, the model is trained using full-batch gradient descent and optimized with the Adam optimizer.

D.2 ADDITIONAL ANALYSIS FOR MFM TOY EXPERIMENTS

As shown in Figure 6, increasing the number of discretized anchor points causes significant fluctuations in the velocity field near π_0 during the forward process, highlighting the impact of trajectory complexity on learning. Furthermore, our method encourages exploration of diverse trajectories, as evidenced by the "turning-back" phenomenon observed in the early stages of the reverse process when the number of discretized anchor points (N) increases. This allows for more exploration in the space rather than directly pointing to the π_0 distribution. Despite the unpredictability of the varying velocity field, the residual correlation between the forward and reverse velocity fields, enabled by the momentum field, facilitates velocity field prediction. Additionally, the piecewise linear nature of the trajectory preserves the accelerated denoising capability of Rectified Flow, enabling the generation of high-quality samples with a small number of steps while maintaining high denoising efficiency.

E ADDITIONAL IMAGE EXPERIMENTS

E.1 EXPERIMENT DETAILS

Dataset Description: We use three datasets for training, including

- 1. CIFAR-10: Images with a resolution of 32×32 from the CIFAR-10 training split.
- 2. CelebA-HQ: Images from the 'img_align_celeba_png.7z' version of the CelebA-HQ dataset, resized to 256×256 resolution.
- 3. ImageNet: Images from ImageNet resized respectively to 32×32 and 64×64 resolutions.

Note that during training, images are normalized to have zero mean and unit variance.

Implementation Details: The experiments are implemented in PyTorch (version 2.6.0) and conducted on an NVIDIA A800-SXM4-80GB GPU. Random seeds are set to 42 for reproducibility.

Performance Details: Given that our training commenced from scratch and employed a relatively simple network architecture (U-net), our baseline performance is not as robust as that of most diffusion models with more intricate designs. However, since all our experiments were conducted within the same architectural framework that we designed, our comparisons remain fair and persuasive.

Additional Notes: Ablation studies are conducted to analyze the impact of different hyperparameters. Hyperparameter tuning is performed using grid search over the learning rate and batch size.

| Dataset | CIFAR-10 | CelebA-HQ | ImageNet | Scope | |
|--|---|--------------------------|-------------------------|---------------------|--|
| resolution params (M) step batch size | 32 35 20k 1024 | 256 120 70k 128 | 64 120 70k 512 | - 45 1k 40 | |
| optimizer learning rate ema decay | $\begin{array}{c} {\rm Adam} \\ 3e-4 \\ 0.9999 \end{array}$ | | | | |

Table 4: Configurations for different datasets.

E.2 EXPERIMENTS ON CELEBA-HQ AND IMAGENET

As shown in Table 2, the momentum flow model consistently achieves significant improvements on CelebA-HQ and ImageNet. By balancing N and γ ,we verify that more anchor points bring greater gains and γ should increase with the number of anchor points to reduce velocity abruptness in momentum field. As shown in 7 and 8xuezh, a larger N brings more significant diversity, but it may compromise image quality, so the noise intensity γ should be correspondingly adjusted. The momentum flow model retains the fast sampling efficiency of rectified flow, while also generating higher-quality images, as shown by these results.

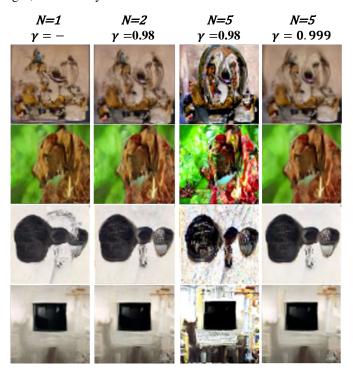


Figure 7: The impact of adjusting N and γ on image details in Momentum Flow (Step=50). Here '–' means that γ do not influence the trajectories while N=1.

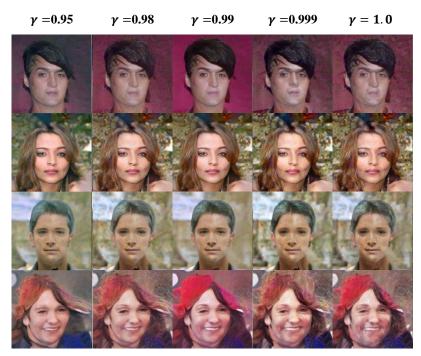


Figure 8: Generated images under different values of gamma, where both excessively large and excessively small gamma values can deteriorate image quality (N = 2, Step = 50).



Figure 9: Face generation by Rectified Flow and Momentum Flow under different sampling steps.

F THEORETICAL PRELIMINARIES ABOUT SO(3) AND SE(3)

This section synthesizes the theoretical foundations of Lie groups SO(3) and SE(3) from two complementary perspectives: geometric structure and representation theory. By integrating these theories, we establish a rigorous mathematical toolkit for rigid-body transformations in computational biology.

F.1 SO(3) LIE GROUP

The Special Orthogonal group in 3 dimensions, SO(3) consists of the 3D rotation matrices:

$$SO(3) = \{ r \in \mathbb{R}^{3 \times 3} : r^{\top} r = r r^{\top} = I, \det r = 1 \}.$$
 (17)

F.1.1 LIE ALGEBRA OF SO(3) AND HAT OPERATION

SO(3) is a matrix Lie group and its Lie algebra $\mathfrak{so}(3)$ consists of all 3×3 skew-symmetric matrices:

$$\mathfrak{so}(3) = \left\{ \mathfrak{r} \in \mathbb{R}^{3 \times 3} : \mathfrak{r}^{\top} = -\mathfrak{r} \right\}. \tag{18}$$

 $\mathfrak{so}(3)$ is 3-dimensional and is isomorphic to the \mathbb{R}^3 vector via the *hat* operation $\widehat{(\cdot)}:\mathbb{R}^3\to\mathfrak{so}(3)$ as

$$\mathbf{r} = \widehat{\boldsymbol{\omega}} = \begin{pmatrix} 0 & -\omega_3 & \omega_2 \\ \omega_3 & 0 & -\omega_1 \\ -\omega_2 & \omega_1 & 0 \end{pmatrix} \in \mathfrak{so}(3), \ \forall \boldsymbol{\omega} = (\omega_1, \omega_2, \omega_3)^{\top} \in \mathbb{R}^3.$$
 (19)

The matrix \mathfrak{r} can be uniquely identified with a vector $\omega \in \mathbb{R}^3$ such that $\forall \mathbf{v} \in \mathbb{R}^3$, $\mathfrak{r}\mathbf{v} = \widehat{\omega}\mathbf{v} = \omega \times \mathbf{v}$, where \times indicates the cross product. The vector ω is known as the *rotation vector*, i.e., *angular velocity*. Moreover, the Lie bracket on $\mathfrak{so}(3)$ corresponds to the cross product in \mathbb{R}^3 :

$$[\widehat{\boldsymbol{\omega}}_1, \widehat{\boldsymbol{\omega}}_2] = \widehat{\boldsymbol{\omega}}_1 \widehat{\boldsymbol{\omega}}_2 - \widehat{\boldsymbol{\omega}}_2 \widehat{\boldsymbol{\omega}}_1 = \widehat{\boldsymbol{\omega}}_1 \times \boldsymbol{\omega}_2. \tag{20}$$

This $\mathfrak{so}(3)$ - \mathbb{R}^3 isomorphism allows the rotation vector $\boldsymbol{\omega} \in \mathbb{R}^3$ to encode both rotation axis (direction of $\boldsymbol{\omega}$) and angle ($||\boldsymbol{\omega}||$) in a unified framework. Specifically, the magnitude of this vector, $\boldsymbol{\theta} = ||\boldsymbol{\omega}||$ denotes the *angle* of rotation and the direction of this vector, $e_{\boldsymbol{\omega}} = \frac{\boldsymbol{\omega}}{||\boldsymbol{\omega}||}$ denotes the *axis* of rotation.

F.1.2 PARAMETERIZATIONS OF SO(3)

Here we describe two different possible parameterizations of SO(3) and its Lie algebra $\mathfrak{so}(3)$.

Axis-angle. Let a unit vector $e_{\omega}=(a,b,c)\in\mathbb{S}^2$ represent the rotation axis, where $(a,b,c)\in\mathbb{R}^3$ and $a^2+b^2+c^2=1$, and $\theta\in\mathbb{R}_+$ represent the rotation angle. Hence, any rotation matrix in SO(3) can be formally written via the exponential mapping as $r=\exp(\mathfrak{r})\in\mathrm{SO}(3)$, where $\mathfrak{r}=\theta X\in\mathfrak{so}(3)$ and $X=aX_1+bX_2+cX_3^2$. The parameterization of SO(3) using (e_{ω},θ) is called the *axis-angle* theory. Notably, $X^3=-X$ and the explicit form of r can be given by the *Rodrigues' formula* as

$$r = \exp(\theta X) = I + \sin \theta \cdot X + (1 - \cos \theta)X^{2},\tag{21}$$

which provides a concise way of computing the exponential. In addition, for $\forall (a, b, c), \mathbf{v} \in \mathbb{R}^3$,

$$X\mathbf{v} = (aX_1 + bX_2 + cX_3)\mathbf{v} = e_{\omega} \times \mathbf{v}, \ X^2\mathbf{v} = (aX_1 + bX_2 + cX_3)^2\mathbf{v} = \langle e_{\omega}, \mathbf{v} \rangle e_{\omega} - \mathbf{v}.$$
 (22)

Then substitute these into the expression above to obtain the Rodrigues' rotation formula:

$$r\mathbf{v} = \exp(\theta X)\mathbf{v} = \cos\theta \cdot \mathbf{v} + \sin\theta \cdot (e_{\omega} \times \mathbf{v}) + (1 - \cos\theta)\langle e_{\omega}, \mathbf{v} \rangle e_{\omega}. \tag{23}$$

As this formula shows, $\exp(\theta X)\mathbf{v}$ denotes the rotation of the vector \mathbf{v} of angle θ around the axis $e_{\boldsymbol{\omega}}$. Moreover, an equivalent representation defines the rotation matrix as $r=\exp(\mathfrak{r})=\exp(\widehat{\boldsymbol{\omega}})$, where $\widehat{\boldsymbol{\omega}}\in\mathfrak{so}(3)$ and $\boldsymbol{\omega}=||\boldsymbol{\omega}||_{||\boldsymbol{\omega}||}^{\underline{\boldsymbol{\omega}}}=\theta\,e_{\boldsymbol{\omega}}$ is the rotation vector. So there exists another expression of r:

$$r = \exp\left(\widehat{\omega}\right) = I + \frac{\sin\theta}{\theta}\widehat{\omega} + \frac{1 - \cos\theta}{\theta^2}\widehat{\omega}^2. \tag{24}$$

Notably, it is continuous at $\theta = 0$, yielding the identity matrix I. And the vector rotation formula is:

$$r\mathbf{v} = \exp(\widehat{\boldsymbol{\omega}})\,\mathbf{v} = \cos\theta \cdot \mathbf{v} + \frac{\sin\theta}{\theta}\,(\boldsymbol{\omega} \times \mathbf{v}) + \frac{1 - \cos\theta}{\theta^2}\langle \boldsymbol{\omega}, \mathbf{v} \rangle \boldsymbol{\omega}. \tag{25}$$

Euler angles. Rotation can also be decomposed into sequential elementary rotations about coordinate axes. A common convention is to utilize the x-convention with three angles $(\psi, \phi, \varphi) \in \mathbb{R}^3$. Specifically, the rotation is given by: a rotation about the z-axis by ψ , a second rotation about the former x-axis by ϕ , and a last one about the former z-axis by φ . It can be formally expressed as

$$r = \exp \left[\psi X_3 \right] \exp \left[\phi X_1 \right] \exp \left[\varphi X_3 \right]$$

$$= \left(\begin{array}{ccc} \cos(\psi) & -\sin(\psi) & 0 \\ \sin(\psi) & \cos(\psi) & 0 \\ 0 & 0 & 1 \end{array}\right) \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & \cos(\phi) & -\sin(\phi) \\ 0 & \sin(\phi) & \cos(\phi) \end{array}\right) \left(\begin{array}{ccc} \cos(\varphi) & -\sin(\varphi) & 0 \\ \sin(\varphi) & \cos(\varphi) & 0 \\ 0 & 0 & 1 \end{array}\right). \tag{26}$$

Technically speaking, these three angles ψ , ϕ , φ are called the *Euler angles*: ψ is called the *precession angle*, ϕ is called the *nutation angle*, and φ is called the *angle of proper rotation* (or *spin*).

 $^{^{2}}X$ is a skew-symmetric matrix in $\mathfrak{so}(3)$ and (X_{1}, X_{2}, X_{3}) is the standard basis of $\mathfrak{so}(3)$.

F.1.3 METRIC ON SO(3)

The metric on a Lie group G is a smooth assignment of the inner product to each of its tangent space $\mathcal{T}_g G^3$, where $g \in G$. Thus, a common way to construct a metric on G is to first define the inner product on \mathfrak{g} and then extend it to the entire group via left (or right) translation. A particularly important class is the *bi-invariant* metric, which maintains invariant under both left and right translations.

Let Q be a symmetric positive-definite matrix defining a quadratic form on g, formally expressed as

$$Q = \begin{pmatrix} A & B^{\top} \\ B & C \end{pmatrix}. \tag{27}$$

Depending on the matrix Q, the inner product between two elements $X_0, Y_0 \in \mathfrak{g}$ is given by

$$\langle X_0, Y_0 \rangle_{\mathfrak{g}} = \operatorname{tr}(X_0^\top Q Y_0). \tag{28}$$

Then this inner product on g can be extended to a left-invariant metric on $g \in G$ via

$$\langle u, v \rangle_{\mathcal{T}_q G} = \langle dL_{g^{-1}} u, dL_{g^{-1}} v \rangle_{\mathfrak{g}}, \ \forall u, v \in \mathcal{T}_g G.$$
 (29)

Specifically, to calculate the inner product of two tangent vectors $u,v\in\mathcal{T}_gG$ at $g\in G$, we first use the differential $dL_{g^{-1}}(\cdot)^4$ to "pull back" u,v to the identity element \mathbf{e} , thus becoming $U=dL_{g^{-1}}(u)\in\mathfrak{g}$ and $V=dL_{g^{-1}}(v)\in\mathfrak{g}$, and then use the inner product $\langle\cdot,\cdot\rangle_{\mathfrak{g}}$ defined in advance to calculate the inner product of U and V. This result is further defined as the inner product of u and v at u0 at u1 at u2 at u3 at u4 at u5 are Similarly, this inner product on u5 can also be extended to a right-invariant metric on u5 via

$$\langle m, n \rangle_{\mathcal{T}_h G} = \langle dR_{h^{-1}} m, dR_{h^{-1}} n \rangle_{\mathfrak{g}}, \ \forall m, n \in \mathcal{T}_h G.$$
 (30)

To sum up, the metric is bi-invariant. Moreover, a canonical choice for the metric of SO(3) is obtained by taking Q = 1/2I, resulting in a bi-invariant metric on SO(3). Therefore, the metric is given by

$$\langle \mathfrak{r}_1, \mathfrak{r}_2 \rangle_{\mathfrak{so}(3)} = \operatorname{tr}(\mathfrak{r}_1^\top Q \mathfrak{r}_2) = \frac{1}{2} \operatorname{tr}(\mathfrak{r}_1^\top \mathfrak{r}_2), \ \forall \mathfrak{r}_1, \mathfrak{r}_2 \in \mathfrak{so}(3). \tag{31}$$

Note that the inner product on *Lie groups* essentially acts on elements of the *Lie algebra* and, since the left action is transitive, this inner product is well-defined for all tangent spaces of the group elements. To further verify the bi-invariance of the SO(3) metric, consider the adjoint action $Ad_r(\mathfrak{r}) = r\mathfrak{r}r^{\top}$ for $\forall r \in SO(3)$ and $\mathfrak{r} \in \mathfrak{so}(3)$. Then the specific action process can be formally expressed as

$$\langle \mathrm{Ad}_r \mathfrak{r}_1, \mathrm{Ad}_r \mathfrak{r}_2 \rangle_{\mathfrak{so}(3)} = \frac{1}{2} \operatorname{tr} \left((r \mathfrak{r}_1 r^\top)^\top (r \mathfrak{r}_2 r^\top) \right) = \frac{1}{2} \operatorname{tr} (r \mathfrak{r}_1^\top \mathfrak{r}_2 r^\top) = \frac{1}{2} \operatorname{tr} (\mathfrak{r}_1^\top \mathfrak{r}_2) = \langle \mathfrak{r}_1, \mathfrak{r}_2 \rangle_{\mathfrak{so}(3)}, \tag{32}$$

where we utilize the cyclic property of the trace⁵ and the orthogonality of r, i.e., $r^{\top}r = I$. This result shows the metric is invariant under the adjoint action, which implies bi-invariance on SO(3).

The geodesic distance between two elements $r_1, r_2 \in SO(3)$ induced by this metric is given by

$$d_{SO(3)}(r_1, r_2) = ||\log(r_1^{\top} r_2)||_F, \tag{33}$$

where log is the matrix logarithm mapping and $||\cdot||_F$ is the Frobenius norm.

F.1.4 THE ISOTROPIC GAUSSIAN DISTRIBUTION ON SO(3)

 $\mathcal{IG}_{SO(3)}$ density. The isotropic Gaussian distribution on SO(3), denoted as $\mathcal{IG}_{SO(3)}$, is parameterized by a mean rotation $r \in SO(3)$ and a concentration parameter $\epsilon \in \mathbb{R}$. It can be expressed in the *axis–angle* representation (refer to App. F.1.2), where the rotation axis is sampled uniformly and the rotation angle θ follows a probability density function (abbreviated as pdf) given by

$$f(\theta, \epsilon) = \sum_{l=0}^{\infty} (2l+1)e^{-l(l+1)\epsilon} \frac{\sin\left((l+1/2)\theta\right)}{\sin(\theta/2)}.$$
 (34)

Although this expression involves a complex infinite series, Matthies et al. (1988) has shown that for $\epsilon \leq 1$, it can be accurately approximated by a closed-form expression:

$$f(\theta, \epsilon) = \sqrt{\pi} \epsilon^{-3/2} e^{\frac{\epsilon - \theta^2/\epsilon}{4}} \frac{\left(\theta - e^{-\pi^2/\epsilon} \left((\theta - 2\pi) e^{\pi\theta/\epsilon} + (\theta + 2\pi) e^{-\pi\theta/\epsilon} \right) \right)}{2\sin\left(\frac{\theta}{2}\right)}.$$
 (35)

⁵This property is embodied in: $\frac{1}{2}\operatorname{tr}(r\mathfrak{r}_1^{\mathsf{T}}\mathfrak{r}_2r^{\mathsf{T}}) = \frac{1}{2}\operatorname{tr}(r(\mathfrak{r}_1^{\mathsf{T}}\mathfrak{r}_2)r^{\mathsf{T}}) = \frac{1}{2}\operatorname{tr}((\mathfrak{r}_1^{\mathsf{T}}\mathfrak{r}_2)r^{\mathsf{T}}r) = \frac{1}{2}\operatorname{tr}(\mathfrak{r}_1^{\mathsf{T}}\mathfrak{r}_2)r^{\mathsf{T}}r$

 $^{^3}$ At the identity element $\mathbf{e} \in G$, the tangent space $\mathcal{T}_{\mathbf{e}}G$ coincides with the Lie algebra \mathfrak{g} .

⁴The notation $dL_{g^{-1}}(\cdot)$ is standard in differential geometry, where $L_{g^{-1}}$ denotes left translation by g^{-1} , defined as $L_{g^{-1}}(h) = g^{-1}h$. And the differential $dL_{g^{-1}}$ is a linear mapping that pull back a tangent vector $u \in \mathcal{T}_g G$ at $g \in G$ to the tangent space at the identity element \mathbf{e} , i.e., $dL_{g^{-1}}: \mathcal{T}_g G \to \mathcal{T}_\mathbf{e} G = \mathfrak{g}$.

Sampling from $\mathcal{IG}_{SO(3)}$. Sampling from $\mathcal{IG}_{SO(3)}$ follows the procedure described by Leach et al. (2022). The rotation angle θ is obtained via inverse transform sampling, using the cumulative distribution function (abbreviated as cdf) derived from the pdf above. And this cdf is normalized appropriately, accounting for the uniform base density on SO(3), i.e., $f(\theta) = \frac{1-\cos\theta}{\pi}$. Moreover, the rotation axis is sampled uniformly from the two-sphere \mathbb{S}^2 . Notably, The closed-form approximation of eq. (35) achieves fast computation of the cdf, thus making the sampling process highly efficient.

F.2 SE(3) LIE GROUP

The Special Euclidean group, denoted as SE(3), constitutes the set of all rigid-body transformations (including rotation and translation) in three-dimensional space. It can be formally defined as

$$SE(3) = \left\{ \begin{pmatrix} r & \boldsymbol{v} \\ 0 & 1 \end{pmatrix} : r \in SO(3), \boldsymbol{v} \in (\mathbb{R}^3, +) \right\}, \tag{36}$$

where each element is represented by a 4×4 matrix. And endowed with the group operation of matrix multiplication, SE(3) can also be seen as a subgroup of the general linear group GL(4, \mathbb{R}).

The corresponding Lie algebra of the Lie group SE(3), i.e., $\mathfrak{se}(3)$, is given by

$$\mathfrak{se}(3) = \left\{ \xi = \begin{pmatrix} \mathfrak{r} & \boldsymbol{v} \\ 0 & 0 \end{pmatrix} : \mathfrak{r} \in \mathfrak{so}(3), \boldsymbol{v} \in \mathbb{R}^3 \right\},$$
 (37)

where \mathfrak{r} can also be denoted as $[\omega]_{\times}$, indicating the skew-symmetric matrix form of its axis-angle representation $\omega \in \mathbb{R}^3$. Note that the tangent space of the translation group $(\mathbb{R}^3,+)$ is isomorphic to \mathbb{R}^3 itself so we can directly use the notation v instead of \mathfrak{v} . Hence, each element $\xi \in \mathfrak{se}(3)$ is uniquely determined by 6 parameters $(\omega,v)\in\mathbb{R}^6$ and there further exists an isomorphism between $\mathfrak{se}(3)$ and \mathbb{R}^6 via the mapping: $\xi\mapsto (\omega,v)^6$. Moreover, since the translation group $(\mathbb{R}^3,+)$ is a normal subgroup of SE(3), the full group can be written as a semi-direct product: SE(3) = SO(3) \ltimes ($\mathbb{R}^3,+$).

Metric on SE(3). While numerous metrics can be defined on SE(3), none of them are bi-invariant. Thus, it is common to construct either a left- or right-invariant metric. A straightforward choice for the quadratic form Q from eq. (27) is setting the matrices $A = C = I_3$ and B = 0 (Park & Brockett, 1994). Consequently, the matrix Q after the assignment can be formally expressed as

$$Q = \begin{pmatrix} I_3 & 0 \\ 0 & I_3 \end{pmatrix}. \tag{38}$$

Utilizing this metric we can define an inner product on SE(3) as $\langle \xi_1, \xi_2 \rangle_{\mathfrak{se}(3)} = \operatorname{tr}(\xi_1^\top Q \xi_2)^7$, where tr is the trace operation. For $\xi_1, \xi_2 \in \mathfrak{se}(3)$, the inner product expands explicitly as

$$\operatorname{tr}(\xi_1^{\top} Q \xi_2) = \operatorname{tr}\begin{pmatrix} \mathfrak{r}_1 & \boldsymbol{v}_1 \\ 0 & 0 \end{pmatrix}^{\top} \begin{pmatrix} I_3 & 0 \\ 0 & I_3 \end{pmatrix} \begin{pmatrix} \mathfrak{r}_2 & \boldsymbol{v}_2 \\ 0 & 0 \end{pmatrix}) = \operatorname{tr} \begin{pmatrix} \mathfrak{r}_1^{\top} \mathfrak{r}_2 & \mathfrak{r}_1^{\top} \boldsymbol{v}_2 \\ \boldsymbol{v}_1^{\top} \mathfrak{r}_2 & \boldsymbol{v}_1^{\top} \boldsymbol{v}_2 \end{pmatrix}. \tag{39}$$

After further derivation, we can obtain: $\operatorname{tr}(\xi_1^\top Q \xi_2) = \operatorname{tr}(\mathfrak{r}_1^\top \mathfrak{r}_2) + \operatorname{tr}(\boldsymbol{v}_1^\top \boldsymbol{v}_2) = \operatorname{tr}(\mathfrak{r}_1^\top Q \mathfrak{r}_2) + \boldsymbol{v}_1^\top \boldsymbol{v}_2^8$. Therefore, the metric on SE(3) can be formally decomposed into the metrics on SO(3) and \mathbb{R}^3 :

$$\langle \xi_1, \xi_2 \rangle_{\mathfrak{se}(3)} = \langle \mathfrak{r}_1, \mathfrak{r}_2 \rangle_{\mathfrak{so}(3)} + \langle v_1, v_2 \rangle_{\mathbb{R}^3}. \tag{40}$$

Hence, geodesics on SE(3) can be derived from those on the product manifold SO(3) $\times \mathbb{R}^3$ and the distance between $x_1 = (r_1, v_1) \in SE(3)$ and $x_2 = (r_2, v_2) \in SE(3)$ is given by

$$d_{SE(3)}(x_1, x_2) = \sqrt{d_{SO(3)}(r_1, r_2)^2 + d_{\mathbb{R}^3}(\boldsymbol{v}_1, \boldsymbol{v}_2)^2}.$$
 (41)

where $d_{SO(3)}$ is defined in eq. (33) and $d_{\mathbb{R}^3}$ denotes the standard Euclidean distance.

⁶The isomorphism between $\mathfrak{se}(3)$ and \mathbb{R}^6 identifies each element $\xi \in \mathfrak{se}(3)$ with a twist comprising rotational (i.e., *angular* velocity ω) and translational components (i.e., *linear* velocity v).

⁷Similarly, the inner product on the Lie group SE(3) essentially acts on elements of its Lie algebra $\mathfrak{se}(3)$.

 $^{{}^8\}boldsymbol{v}_1^{\top}$ is a 3-dim row vector and \boldsymbol{v}_2 is a 3-dim column vector, so $\boldsymbol{v}_1^{\top}\boldsymbol{v}_2$ is a 1×1 matrix, and $\operatorname{tr}(\boldsymbol{v}_1^{\top}\boldsymbol{v}_2) = \boldsymbol{v}_1^{\top}\boldsymbol{v}_2$.

G ADDITIONAL DETAILS ABOUT PROTEIN BACKBONE GENERATION

G.1 PROTEIN BACKBONE PARAMETERIZATION

Here our protein backbone parameterization follows the seminal work of AlphaFold2 (Jumper et al., 2021) in that we associate a rigid-body *frame* with each residue in the amino acid sequence. Specifically, an N residue backbone is parameterized by a collection of N orientation preserving rigid transformations, i.e., *frames*, and each frame maps from fixed coordinates of four heavy atoms $N^*, C^*_{\alpha}, C^*, O^* \in \mathbb{R}^3$ centered at $C^*_{\alpha} = (0,0,0)$. Notably, each atom coordinate assumes chemically

idealized bond angles and lengths measured experimentally (Engh & Huber, 2012). Thus, residue $i \in [N]$ is denoted as an action of $T^{(i)}$ on idealized coordinates of the backbone main atoms $[N^{(i)}, C_{\alpha}^{(i)}, C^{(i)}] = T^{(i)} \cdot [N^*, C_{\alpha}^*, C^*],$ where $T^{(i)}$ is a member of the special Euclidean group SE(3), the set of orientation preserving rigid transformations in Euclidean space. Each $T^{(i)}$ can be formally decomposed into two components $T^{(i)} = (r^{(i)}, x^{(i)})$ where $r^{(i)} \in SO(3)$ is a 3×3 rotation matrix and $x^{(i)} \in \mathbb{R}^3$ represents a translation vector. And we collectively denote all N frames as $\mathbf{T} = [T^{(1)}, \dots, T^{(N)}] \in$ $SE(3)^N$. Moreover, to construct the backbone oxygen atom O, we rotate O* around the bond between C_{α} and C with an additional torsion angle $\psi \in SO(2)$. Figure 10 visually shows our backbone parameterization with frames.

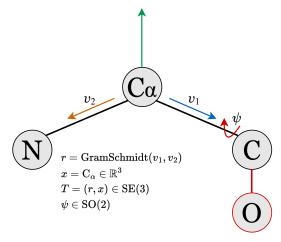


Figure 10: Protein parameterization with frames.

G.2 CONVERSION BETWEEN COORDINATES AND FRAMES

As discussed above, N^* , C^* , O^* are idealized atom coordinates that assumes chemically idealized bond angles and lengths measured experimentally (Engh & Huber, 2012). However, these coordinates differ slightly per amino acid type. Here we uniformly take the idealized values of Alanine which are

$$\begin{split} \mathbf{N}^* &= (-0.525, 1.363, 0.0), \\ \mathbf{C}^*_\alpha &= (0.0, 0.0, 0.0), \\ \mathbf{C}^* &= (1.526, 0.0, 0.0), \\ \mathbf{O}^* &= (0.627, 1.062, 0.0). \end{split}$$

Notably, these idealized values are derived with C_{α}^* as the origin. And using a central frame $T^{(i)}$, we can construct the realistic backbone main atoms of residue i via $[N^{(i)}, C_{\alpha}^{(i)}, C^{(i)}] = T^{(i)} \cdot [N^*, C_{\alpha}^*, C^*]$. And the realistic backbone oxygen requires rotating a idealized oxygen around the $C - C_{\alpha}$ bond:

$$\mathbf{O}^{(i)} = T^{(i)} \cdot T_{\mathrm{psi}}^*(\psi^{(i)}) \cdot \mathbf{O}^*,\tag{42}$$

where $\psi^{(i)}$ is a backbone torsion angle of residue i and $T^*_{psi}(\psi^{(i)}) = (r_x(\psi^{(i)}), x_{psi})$ is a Euclidean transformation that maps the central frame $T^{(i)}$ to a new frame $T^{(i)} \cdot T^*_{psi}(\psi^{(i)})$ centered at C and rotated around the x-axis by $\psi^{(i)}$. Note $\psi^{(i)}$ is a tuple of two values describing a point on the unit circle, $\psi^{(i)} = [\psi_1^{(i)}, \psi_2^{(i)}]$ where $(\psi_1^{(i)})^2 + (\psi_2^{(i)})^2 = 1$. So $r_x(\psi^{(i)})$ and $x_{psi}^{(i)}$ can be expressed as

$$r_x(\psi^{(i)}) = \begin{pmatrix} 1 & 0 & 0\\ 0 & \psi_1^{(i)} & -\psi_1^{(i)}\\ 0 & \psi_2^{(i)} & \psi_1^{(i)} \end{pmatrix}, \ x_{psi} = (1.526, 0.0, 0.0). \tag{43}$$

⁹The translation vector x_{psi} transfers the oxygen atom with the rotation $r_x(\psi^{(i)})$ applied from the coordinate system with the origin C_α to the coordinate system with the origin C_α . Specifically, $T^{(i)}$ maps the ideal coordinate system (origin at C_α^*) to the real coordinate system (origin at $C_\alpha^{(i)}$). In real space, the vector from $C_\alpha^{(i)}$ to $C^{(i)}$ is the same as the vector (1.526,0,0) from C_α^* to C^* in ideal space (because $T^{(i)}$ is a rigid-body transformation, preserving local distances and angles). So the translation vector x_{psi} directly uses the ideal vector (1.526,0,0).

To sum up, the mapping from frames to idealized coordinates, i.e., frame2atom, is implemented by combining $[N^{(i)}, C_{\alpha}^{(i)}, C^{(i)}] = T^{(i)} \cdot [N^*, C_{\alpha}^*, C^*]$ and eq. (42), which can be formally expressed as

$$[N^{(i)}, C^{(i)}_{\alpha}, C^{(i)}, O^{(i)}] = \text{frame2atom}(T^{(i)}, \psi^i).$$
 (44)

We next introduce constructing rigid-body frames from atom coordinates, i.e., atom2frame. Each frame can be obtained as described in the rigidFrom3Points algorithm in AF2 (Jumper et al., 2021):

Algorithm 3: Rigid from 3 points using the Gram-Schmidt process

Data: Coordinates of the three backbone atoms $N^{(i)}, C^{(i)}_{\alpha}, C^{(i)}$ of residue i

Result: The rigid-body frame $T^{(i)}$ corresponding to residue i

The conversion from coordinates to frames can be expressed as $T^{(i)} = \text{atom2frame}(N^{(i)}, C^{(i)}, C^{(i)})$.

G.3 FRAMEPRED ARCHITECTURE

Overview of FramePred. To predict the rigid-body frame for every protein residue, we utilize the FramePred architecture introduced in FrameDiff (Yim et al., 2023b) which performs iterative updates to the frames over a series of L layers using a combination of *spatial* and *sequence* based attention modules. Specifically, $\mathbf{h}_{\ell} = [h_{\ell}^{(1)}, \cdots, h_{\ell}^{(N)}] \in \mathbb{R}^{N \times D_h}$ are node embeddings of the ℓ -th layer where $h_{\ell}^{(i)}$ is the embedding for residue $i \in [N]$. And $\mathbf{z}_{\ell} \in \mathbb{R}^{N \times N \times D_z}$ are edge embeddings with $z_{\ell}^{(nm)}$ encoding the edge between residues n and m. The frame of each residue at the ℓ -th is denoted as $\mathbf{T}_{\ell} \in \mathrm{SE}(3)^N$. Unless stated otherwise, all instances of Multi-Layer Perceptrons (MLPs) use 3 Linear layers with biases, ReLU activation, and LayerNorm after the final layer. When FramePred is running, node embeddings \mathbf{h}_{ℓ} are first updated using Invariant Point Attention (IPA) (Jumper et al., 2021) with a skip connection. Before Transformer, the initial node embeddings \mathbf{h}_0 and post-IPA embeddings are concatenated. After transformer, we include a skip connection with post-IPA embeddings. The updated node embeddings $\mathbf{h}_{\ell+1}$ are then used to update edge embeddings $\mathbf{z}_{\ell+1}$ as well as predict frame updates $\mathbf{T}_{\ell+1}$. And so on to get the final frames \mathbf{T}_L of all protein residues.

Feature initialization. Following the methodology established by Trippe et al. (2022), node and edge embeddings are initialized using a combination of residue indices and timestep information. Specifically, sinusoidal embeddings $\phi(\cdot)$ (Vaswani et al., 2017) are applied to these features, after which an MLP is used to project them into the desired embedding space. For residue $i \in [N]$, the initial node embedding at layer 0 incorporates the residue index i and the diffusion timestep t, i.e., $h_0^{(i)} = \text{MLP}([\phi(n), \phi(t)]) \in \mathbb{R}^{D_h \cdot 10}$, where D_h denotes the dimension of node embeddings. Moreover, for a residue pair (n, m), the edge embedding $z_0^{(nm)}$ additionally includes the relative sequence distance $\phi(m-n)$ and a binned displacement feature derived from self-conditioned C_α coordinates $\phi(\text{disp}_{sc}^{(nm)})$. The initial edge embeddings can be formally expressed as

$$z_0^{(nm)} = \text{MLP}([\phi(n), \phi(m), \phi(m-n), \phi(t), \phi(\text{disp}_{sc}^{(nm)})]) \in \mathbb{R}^{D_z},$$
(45)

¹⁰Here we stipulate that superscripts without parentheses are used to refer to time step, superscript numbers within parentheses refer to residue indices, and subscripts refer to variable names (layer indices in most cases).

where D_z denotes the dimension of edge embeddings and disp_{sc} is the self-conditioning of predicted C_α displacements. Specifically, let \hat{x}_{sc} be the C_α coordinates (in Å) predicted during self-conditioning. And to prevent over-reliance on self-conditioned outputs, we set $\hat{x}_{sc}=0$ with 50% probability during training. The binned displacement between residues n and m is formally expressed as

$$\operatorname{disp}_{sc}^{(mn)} = \sum_{i=1}^{N_{\text{bins}}} \mathbb{1}\{|\hat{x}_{sc}^{(n)} - \hat{x}_{sc}^{(m)}| < \nu_i\},\tag{46}$$

where $\nu_1, \dots, \nu_{N_{\text{bins}}} = \text{Linspace}(0, 20)^{11}$ are equally spaced intervals between 0 and 20 angstroms.

To construct initial frames, C_{α} coordinates are first zero-centered and all backbone coordinates (N, C_{α}, C, O) are scaled to nanometers as done in AF2 (Jumper et al., 2021) by multiplying coordinates by 1/10. We then construct initial frames for each protein residue i as

$$T^{0,(i)} = (r^{0,(i)}, x^{0,(i)}) = \text{atom2frame}(N^{(i)}, C_{\alpha}^{(i)}, C^{(i)}).$$
(47)

Node update. IPA is first introduced in AF2 (Jumper et al., 2021) and we apply this algorithm without modifications. And Transformer is also used without modification from Vaswani et al. (2017). Node update is formally represented as follows, including specific operations and data dimensions.

$$\begin{split} \mathbf{h}_{ipa} &= \operatorname{LayerNorm}(\operatorname{IPA}(\mathbf{h}_{\ell}, \mathbf{z}_{\ell}, \mathbf{T}_{\ell}) + \mathbf{h}_{\ell}) \in \mathbb{R}^{N,D_h} \\ \mathbf{h}_{skip} &= \operatorname{Linear}(\mathbf{h}_0) \in \mathbb{R}^{N,D_{skip}} \\ \mathbf{h}_{in} &= \operatorname{concat}(\mathbf{h}_{ipa}, \mathbf{h}_{skip}) \in \mathbb{R}^{N,(D_{skip}+D_h)} \\ \mathbf{h}_{trans} &= \operatorname{Transformer}(\mathbf{h}_{in}) \in \mathbb{R}^{N,(D_{skip}+D_h)} \\ \mathbf{h}_{out} &= \operatorname{Linear}(\mathbf{h}_{trans}) + \mathbf{h}_{\ell} \in \mathbb{R}^{N,D_h} \\ \mathbf{h}_{\ell+1} &= \operatorname{MLP}(\mathbf{h}_{out}) \in \mathbb{R}^{N,D_h} \end{split}$$

Edge update. Each directed edge is updated through an MLP of the current edge and the embeddings of the source and target nodes. Edge update is also formally expressed as follows.

$$\mathbf{h}_{\text{down}} = \text{Linear}(\mathbf{h}_{\ell+1}) \in \mathbb{R}^{N,D_h/2}$$

$$z_{\text{in}}^{(nm)} = \text{concat}(h_{\text{down}}^{(n)}, h_{\text{down}}^{(m)}, z_{\ell}^{(nm)}) \in \mathbb{R}^{N,(D_h+D_z)}$$

$$\mathbf{z}_{\ell+1} = \text{LayerNorm}(\text{MLP}(\mathbf{z}_{\text{in}})) \in \mathbb{R}^{N,N,D_z}$$

Notably, in the first line, node embeddings are first projected down to half the dimension.

Backbone update. As for the backbone update, we follow the Backbone Update algorithm proposed in AF2 (Jumper et al., 2021) and present its specific operations in detail as follows.

$$\begin{split} b^i, c^i, d^i, x_{\text{update}}^{(i)} &= \text{Linear}(h_\ell^i) \\ (a^i, b^i, c^i, d^i) &= (1, b^i, c^i, d^i) / \sqrt{1 + b^i + c^i + d^i} \\ r_{\text{update}}^{(i)} &= \begin{pmatrix} (a^i)^2 + (b^i)^2 - (c^i)^2 - (d^i)^2 & 2b^ic^i - 2a^id^i & 2b^id^i + 2a^ic^i \\ 2b^ic^i + 2a^id^i & (a^i)^2 - (b^i)^2 + (c^i)^2 - (d^i)^2 & 2c^id^i - 2a^ib^i \\ 2b^id^i - 2a^ic^i & 2c^id^i + 2a^ib^i & (a^i)^2 - (b^i)^2 - (c^i)^2 + (d^i)^2 \end{pmatrix} \\ T_{\text{update}}^{(i)} &= (r_{\text{update}}^{(i)}, x_{\text{update}}^{(i)}) \\ T_{\ell+1}^{(i)} &= T_\ell^{(i)} \cdot T_{\text{update}}^{(i)} \end{split}$$

where $b^i, c^i, d^i \in \mathbb{R}^{12}$, $r_{\text{update}}^{(i)} \in \text{SO}(3)$, and $x_{\text{update}}^{(i)} \in \mathbb{R}^3$. Here we first constructs a normalized quaternion (2nd line) and then convert it into a valid rotation matrix (3rd line).

Torsion Prediction. We still follow AF2 (Jumper et al., 2021) to predict the torsion angle ψ .

$$\mathbf{h}_{ ext{psi}} = \operatorname{MLP}(\mathbf{h}_L) \in \mathbb{R}^{N,D_h}$$
 $oldsymbol{\psi}_{ ext{unnormalized}} = \operatorname{Linear}(\mathbf{h}_{ ext{psi}} + \mathbf{h}_L) \in \operatorname{SO}(2)^N$
 $\hat{oldsymbol{\psi}} = oldsymbol{\psi}_{ ext{unnormalized}} / ||oldsymbol{\psi}_{ ext{unnormalized}}|| \in \operatorname{SO}(2)^N$

¹¹In our experiments we set $N_{\text{bins}} = 22$.

¹²Due to space limitations, we use superscripts without parentheses instead of superscripts with parentheses.

H COMPLETE PROTEIN EXPERIMENTS

H.1 EXPERIMENTAL SETUP

Training. To verify the effectiveness of Momentum Flow on the protein monomer generation task, following GENIE (Lin & AlQuraishi, 2023) and FrameFlow (Yim et al., 2023a), we train it on the SCOPe dataset (Chandonia et al., 2022) with proteins below length 128 for a total of 3, 938 examples. Our model is trained for 10 hours on two NVIDIA A800-80G GPUs using the batching strategy from FrameDiff (Yim et al., 2023b) of combining proteins with the same length into the same batch to remove extraneous padding and we take the optimal 3 checkpoints for evaluation. We use the Adam (Adam et al., 2014) optimizer with a learning rate of 0.0001, $\beta_1 = 0.9$, and $\beta_2 = 0.999$.

Metrics. During evaluation, we sample 10 backbones for every length between 60 and 300 then use ProteinMPNN (Dauparas et al., 2022) to design 8 sequences for each backbone. Notably, the upper limit of 300 here differs from the upper limit of 128 during training. We increase the upper limit during evaluation to show the generalization of our model in generating long sequence proteins.

The assessment of generated protein structures employs complementary metrics evaluating distinct aspects of protein quality. *Accuracy* is quantified by comparing predictions against native structures in PDB using scRMSD for atomic-level precision and TM-score (referred to as pdbTM) for global topological fidelity, with a TM-score > 0.5 indicating a correct fold. Intrinsic structural plausibility, a proxy for *designability*, is assessed using *confidence* estimates from deep learning models: pLDDT reports per-residue local reliability, while pAE evaluates the self-consistency of long-range interresidue distances. Moreover, scTM and scRMSD also serve as the fundamental distance measures for quantifying *novelty* (against known structures in PDB) and *diversity* (within a generated ensemble).

Baselines. We compare our results to GENIE (Lin & AlQuraishi, 2023) and FrameFlow (Yim et al., 2023a), a diffusion model and a rectified flow model for protein backbone generation, respectively, that do not rely on a pre-trained folding network. We retrain both models according to their default recommended settings. Our baselines are intended to demonstrate tradeoffs in efficiency and diversity.

H.2 HYPERPARAMETERS

Neural network hyperparameters.

 $\begin{array}{lll} \mbox{Global parameters}: & D_h = 256 & D_z = 128 & D_{\rm skip} = 64 & L = 4 \\ \mbox{IPA parameters}: & \mbox{heahs} = 8 & \mbox{query points} = 8 & \mbox{value points} = 12 \end{array}$

Transformer parameters : heads = 4 layers = 2

With these parameters, our neural network has 17, 446, 190 trainable weights.

SDE parameters.

Translations: schedule = linear $\beta_{\min} = 0.1$ $\beta_{\max} = 20$ Rotations: schedule = logarithmic $\beta_{\min} = 0.1$ $\beta_{\max} = 1.5$

H.3 FURTHER EXPERIMENTAL ANALYSIS

We use the Euler-Maruyama integrator for SDE sampling and the Euler integrator for ODE sampling. The number of integration timesteps for all methods is set to 100. Quantitative results are shown in Table 3. We use SDE sampling for GENIE (Lin & AlQuraishi, 2023) and ODE sampling for FrameFlow (Yim et al., 2023a) since these are the methods used in their respective papers.

As illustrated in Table 3, both GENIE (Lin & AlQuraishi, 2023) and FrameFlow (Yim et al., 2023a) exhibit limitations—either in sampling fidelity or structural plausibility—our model consistently outperforms them across all four metrics. Our momentum flow significantly improves structural accuracy, as evidenced by the highest scTM (0.47) and lowest scRMSD (8.05), indicating generated backbones more closely resemble native-like folds. Meanwhile, our momentum flow enhances designability, reflected in the highest pLDDT (70.09) and lowest pAE (9.50), suggesting superior local and global structural self-consistency without relying on ground-truth alignment.

While GENIE has less parameters (4.1M) than FrameFlow/Momentum Flow (17.4M), i.e., the FramePred architecture introduced in App. G.3, it uses expensive triangle updates (Jumper et al., 2021) that requires high memory cost and greater compute for each forward call. Sampling a length 100 protein with 100 timesteps on an NVIDIA A800-80G GPU takes GENIE around 10 seconds while for FrameFlow/Momentum Flow sampling with 100 timesteps takes around 4 seconds.

In conclusion, our Momentum Flow achieves a more favorable efficient-diverse trade-off, where high-quality and various protein backbone samples can be generated with reduced computational overhead—a critical advantage for practical protein design applications.

I LLM USAGE STATEMENT

During the preparation of this manuscript, we employed GPT-5 exclusively for language refinement. The model was instructed to improve grammar, clarity, and readability while preserving the original meaning of the content.