
Scaling Flow Matching Models at Inference-Time by Search and Path Exploration

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Inference-time compute scaling has become a powerful means of improving discrete
2 generative models, yet its counterpart for continuous-time generative models
3 remains under-explored. We close this gap for flow-matching (FM) models, whose
4 deterministic ordinary-differential (ODE) and stochastic differential (SDE) sam-
5 plers underpin state-of-the-art continuous generative models in domains such as
6 image-generation and protein-folding. We introduce an inference-time strategy that
7 injects analytically constructed divergence-free perturbations into the learned veloc-
8 ity field during ODE sampling. The perturbations preserve the method’s defining
9 linear interpolation path and exactly maintain the continuity equation, ensuring that
10 probability mass is conserved while enabling the sampler to explore high quality
11 trajectories, without modifying trained parameters. We additionally include SDE
12 sampling as an alternative approach. Experiments on ImageNet and FoldFlow
13 demonstrate our methods ability to trade-off computation time with sampling qual-
14 ity over multiple key metrics for each domain. Our work positions inference-time
15 scaling as a principled, training-free lever for enhancing flow-matching models and
16 invites future exploration across diverse continuous generative tasks.

17 1 Introduction

18 Inference-time compute scaling has emerged as a transformative paradigm in generative model-
19 ing, enabling performance improvements through increased computational resources at test time
20 without requiring model retraining. While this approach has shown remarkable success in discrete
21 domains—from language models like OpenAI’s O1 ? to diffusion models for image generation ?—its
22 application to continuous-time generative models remains largely unexplored.

23 Flow Matching (FM) ? has established itself as a leading framework for continuous generative
24 modeling, offering deterministic sampling through ordinary differential equations (ODEs) and
25 achieving state-of-the-art results across diverse domains including image generation ? and protein
26 design ?. However, the deterministic nature of FM sampling presents a fundamental challenge for
27 inference-time scaling: unlike diffusion models that naturally provide stochasticity through noise
28 injection, FM models follow a single, predetermined trajectory from noise to data, leaving limited
29 opportunities for exploring alternative generation paths.

30 This work introduces the first principled approach to inference-time scaling for flow matching
31 models. Our key insight is that we can inject carefully constructed divergence-free perturbations
32 into the learned velocity field during sampling, creating a family of alternative ODE trajectories
33 while preserving the fundamental mathematical properties of the original model. Specifically, our
34 perturbations maintain the continuity equation at the trajectory level, ensuring exact conservation of
35 probability mass throughout the generation process.

36 We make several key contributions:

37 **Theoretical Framework:** We prove that adding divergence-free perturbations to the velocity field
 38 preserves the continuity equation, providing mathematical guarantees that our modifications do not
 39 violate the underlying probability flow. Our approach maintains the linear interpolation paths and
 40 ODE structure that define flow matching, unlike concurrent work that converts to SDE sampling.

41 **Practical Algorithms:** We develop two complementary strategies: (1) divergence-free path explo-
 42 ration that generates multiple trajectories from a single initial condition, and (2) a two-stage approach
 43 combining random search over initial conditions with divergence-free exploration around promising
 44 candidates. This mirrors the exploration-exploitation tradeoff in reinforcement learning.

45 **Empirical Validation:** We demonstrate the effectiveness of our approach across two distinct domains.
 46 On ImageNet 256×256 using SiT-XL/2, our methods achieve significant improvements in FID,
 47 Inception Score, and DINO accuracy as compute budget increases. On protein design using FoldFlow,
 48 we show substantial improvements in structural quality (TM-score) from 0.743 to 0.868, establishing
 49 the cross-domain applicability of our approach.

50 **Methodological Insights:** Our experiments reveal that different inference-time scaling strategies
 51 excel under different conditions: random search provides broad exploration particularly effective for
 52 discrete sequence-structure relationships, while divergence-free perturbations enable focused local
 53 search around high-quality regions. The combination of both approaches consistently outperforms
 54 either method in isolation.

55 Our work establishes inference-time scaling as a viable and principled approach for enhancing
 56 flow matching models, opening new directions for improving continuous generative models across
 57 scientific and creative applications without the computational cost of retraining.

58 2 Preliminaries

59 **Notation.** Random variables are uppercase (e.g., X), and realizations are lowercase (e.g., x).
 60 Distributions at time $t \in [0, 1]$ are denoted π_t with density p_t ; endpoints are $\pi_0 = \pi_{\text{ref}}$ and
 61 $\pi_1 = \pi_{\text{data}}$. Vector fields $v :^d \times [0, 1] \rightarrow ^d$ are time-dependent and (unless stated otherwise) Lipschitz
 62 in x and measurable in t . The continuity equation is

$$\partial_t p_t(x) + \nabla \cdot (p_t(x) v(x, t)) = 0, \quad p_0, p_1 \text{ given.} \quad (1)$$

63 2.1 Flow Matching

64 Flow Matching (FM) ? defines a continuous bridge between a reference distribution π_{ref} , typically a
 65 standard Gaussian, and a data distribution π_{data} , by modeling trajectories x_t that interpolate between
 66 samples $x_0 \sim \pi_{\text{ref}}$ and $x_1 \sim \pi_{\text{data}}$. A common choice is the linear path $x_t = (1 - t)x_0 + tx_1$,
 67 though other conditional paths are possible and lead to generalized or conditional flow matching ???.
 68 This is done via a learned velocity field $v_\theta(x, t)$ that satisfies the probability-flow ODE:

$$\frac{dx_t}{dt} = v_\theta(x_t, t), \quad t \in [0, 1].$$

69 To train v_θ , a supervised loss is used where the ground truth velocity is known analytically:

$$v^*(x_t, t) = x_1 - x_0.$$

70 This target arises because $\frac{dx_t}{dt} = x_1 - x_0$ under linear interpolation. The training loss is then

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{x_0 \sim \pi_{\text{ref}}, x_1 \sim \pi_{\text{data}}, t \sim \mathcal{U}[0, 1]} \left[\|v_\theta(x_t, t) - (x_1 - x_0)\|^2 \right].$$

71 Unlike diffusion models that rely on stochastic sampling from SDEs or discrete Markov chains, FM
 72 offers a deterministic, fast, and interpretable sampling process. Because the interpolant is linear and
 73 the learned dynamics are smooth, FM enables fewer sampling steps while maintaining sample quality.

74 2.2 Minibatch Optimal Transport Flow Matching (OT-FM)

75 While FM defines its objective using i.i.d. sample pairs (x_0, x_1) , such pairs are typically poorly
 76 coupled in high-dimensional space. Minibatch OT-FM ? addresses this by using an optimal transport

(OT) plan computed within each minibatch to generate more meaningful pairs. That is, given batches $\{x_{0,i}\}_{i=1}^B$ from π_{ref} and $\{x_{1,j}\}_{j=1}^B$ from π_{data} , an optimal permutation matrix $\pi^* \in \Pi_B$ is computed to minimize a transport cost:

$$\pi^* = \arg \min_{\pi \in \Pi_B} \sum_{i,j} \pi_{ij} \cdot c(x_{0,i}, x_{1,j}), \quad c(\cdot, \cdot) = \|x_{0,i} - x_{1,j}\|^2$$

Computing exact OT scales cubically in batch size and quadratically in memory, so minibatch approximations are necessary for large datasets. Although such plans only approximate the global coupling, prior work has shown minibatch OT to be effective in practice for generative modeling and domain adaptation. These better-aligned pairs shorten transport paths, stabilize training, and improve generalization.

2.3 Stochastic Interpolants: A Unifying Framework

The stochastic interpolants framework shows that both FM and diffusion models can be described using the same general family of interpolations. A stochastic interpolant is defined by

$$x_t = a(t)x_0 + b(t)x_1 + \sigma(t)\epsilon, \quad \epsilon \sim \mathcal{N}(0, I),$$

with boundary conditions $a(0) = 1$, $b(0) = 0$, $\sigma(0) = 0$ and $a(1) = 0$, $b(1) = 1$, $\sigma(1) = 0$, and with a, b, σ smooth in t . FM is the deterministic special case $a(t) = 1 - t$, $b(t) = t$, $\sigma(t) \equiv 0$.

Crucially, the stochastic interpolants framework allows inference-time reinterpretation of trained models by simply modifying the schedule (a, b, σ) . For instance, a model trained with linear FM ($a(t) = 1 - t$, $b(t) = t$, $\sigma(t) \equiv 0$) can, at test time, be sampled using a different valid schedule. One example is the rotational path $a(t) = \cos(\frac{\pi t}{2})$, $b(t) = \sin(\frac{\pi t}{2})$, $\sigma(t) \equiv 0$, which still interpolates smoothly between endpoints. Another is the variance-preserving (VP) diffusion schedule, usually written in the two-coefficient form $x_t = \alpha_t x_{\text{data}} + \sigma_t \epsilon$ with $\alpha_t^2 + \sigma_t^2 = 1$. This is a special case of the three-coefficient interpolant by reinterpreting $(x_0, x_1) = (\epsilon, x_{\text{data}})$ and setting $a(t) = \sigma_t$, $b(t) = \alpha_t$, $\sigma(t) = 0$. In both cases, the learned parameters can be reused without retraining, revealing a continuum of models and motivating inference-time strategies that use the same velocity or score network.

2.4 Inference-Time Compute Scaling in Generative AI

Inference-time scaling refers to performance improvements achieved by increasing computational resources *after training*, without modifying model parameters. While generative models have traditionally scaled performance by increasing data, model size, or training compute, recent research has explored whether additional computation at inference can yield better outputs.

In discrete domains such as language modeling, inference scaling has been explored through methods that allocate more compute per query. These include using longer reasoning chains, recursive planning, or verifier-guided editing. Notably, OpenAI’s O1 and O3 models and DeepSeek R1 do not use branching across multiple sampled responses. Instead, they increase inference-time compute by allowing more function evaluations per output, for example via deeper chains-of-thought or tree-based planning structures.

In contrast, recent work in diffusion models has implemented inference-time scaling via *search over generation trajectories*. Diffusion models naturally provide stochasticity through their noise-injection process. The inference-time scaling diffusion framework identifies that *some initial noises are better than others*, and proposes to search over this noise space using verifiers to guide sample selection. This turns the generation problem into a two-axis search: one axis governs the verifier that provides feedback (e.g., Inception Score, CLIP, DINO), and the other defines the search algorithm (e.g., random search, zero-order optimization, path-space exploration). This approach allows performance to continue improving beyond what is achievable by simply increasing denoising steps.

These developments motivate the need for inference-time scaling in *continuous-time models* such as Flow Matching (FM), which lack natural sources of stochasticity or branching. The key challenge is that FM models sample deterministically via a single integration of the learned velocity field along a fixed interpolant. Our method addresses this by introducing divergence-free perturbations to the velocity field, thereby defining a family of ODEs that maintain the continuity equation but

diverge in geometry. This enables principled sampling diversity and verifier-guided branching in FM, filling a crucial gap in inference-time scaling for continuous-time generative models.

3 Related Work

3.1 Inference-Time Scaling for Diffusion Models

Inference-time scaling in diffusion models has recently been advanced through optimization in the latent noise space ?. The core idea is to treat the initial noise vector $z \sim \mathcal{N}(0, I)$ as a controllable input, and to iteratively refine it using a verifier score $r(\cdot)$ that evaluates the final generated sample. The method begins with a population of candidate noise vectors $\{z_i\}$, denoises each to obtain $x_0^{(i)}$, scores them, and retains high-scoring candidates. New proposals are then generated by applying small perturbations to the best z_i , repeating the procedure over multiple rounds.

In addition to this search in noise space, the method proposes a "search over paths" strategy, in which denoised samples are partially re-noised and then denoised again. That is, after reaching an intermediate timestep t , the top samples are re-noised forward to $t' > t$, and denoising resumes from t' to $t = 0$. This is intended to explore local perturbations in trajectory space around high-scoring samples.

However, the actual implementation uses hyperparameters such that the search begins at $t = 0.11$, and each re-noising step applies a forward process to $t' = 0.89$. This implies that 89% of the diffusion process is re-applied after a brief denoising. Because most of the signal is lost at this noise level, the re-noised samples are nearly indistinguishable from fresh random samples from the prior. Consequently, while this method technically performs a limited search over paths, in practice it behaves similarly to a best-of- N strategy (often termed random search) conducted over initial noise vectors.

These strategies are highly effective in the diffusion setting, but they rely fundamentally on starting generation from a known distribution, typically $\mathcal{N}(0, I)$. In contrast, our approach applies to flow matching models trained with arbitrary or learned base distributions π_{ref} , which may not admit tractable sampling via random search.

3.2 Particle Sampling Approaches

In addition to search in noise and path space, diffusion models greatly benefit from particle sampling strategies. The simplest is Best-of- N (BoN) ??, which generates multiple samples and retains the highest scoring one. More sophisticated methods integrate the selection into the denoising process: SVDD ? selects the top particles at every step, CoDe ? performs selection at fixed intervals, and Sequential Monte Carlo (SMC) samplers ?? resample particles proportionally to their reward-weighted importance scores. These approaches have proven highly effective in using additional inference-time compute but have not been directly applied to flow matching models due to their deterministic ODE sampling process: given the same initial condition, all particles collapse to a single trajectory.

3.3 Inference-Time Scaling for Flow Matching Models

A concurrent line of work proposes inference-time scaling for flow models by introducing stochasticity and path diversity into the otherwise deterministic sampling process ?. This is achieved through a two-step transformation of the learned model: (1) converting the velocity field from the probability flow ODE to a score-based SDE sampler, and (2) replacing the standard linear interpolant with a variance-preserving (VP) interpolant path to enhance exploration.

Specifically, the learned velocity $u_t(x)$ is used to define the drift term of a reverse-time SDE:

$$dx_t = f_t(x_t) dt + g_t dW_t, \quad \text{where } f_t(x_t) = u_t(x_t) - \frac{g_t^2}{2} \nabla \log p_t(x_t).$$

The score function $\nabla \log p_t(x_t)$ is estimated analytically from u_t using the stochastic interpolants framework ?:

$$\nabla \log p_t(x_t) = \frac{1}{\sigma_t} \cdot \frac{\alpha_t u_t(x_t) - \dot{\alpha}_t x_t}{\dot{\alpha}_t \sigma_t - \alpha_t \dot{\sigma}_t}.$$

170 In parallel, the interpolation path is converted from a linear interpolant $x_t = (1 - t)x_0 + tx_1$ to a VP
 171 interpolant, such as $x_t = \alpha_t x_0 + \sigma_t x_1$ where $\alpha_t^2 + \sigma_t^2 = 1$. This conversion requires transforming
 172 the original velocity field into one compatible with the new interpolant ?:

$$\bar{u}_s(\bar{x}_s) = \frac{\dot{c}_s}{c_s} \bar{x}_s + c_s \dot{t}_s u_{t_s}(\bar{x}_s/c_s),$$

173 where $c_s = \bar{\sigma}_s/\sigma_{t_s}$ and $t_s = \rho^{-1}(\bar{\rho}(s))$ is defined via signal-to-noise ratio schedules $\rho(t) = \alpha_t/\sigma_t$,
 174 $\bar{\rho}(s) = \bar{\alpha}_s/\bar{\sigma}_s$.

175 This approach enables the use of particle sampling strategies originally developed for diffusion models,
 176 which benefit from diverse sample paths and stochastic exploration. However, by transforming both
 177 the dynamics and the interpolant, the method loses the key benefits of flow matching: fast sampling
 178 via few deterministic steps and linear interpolants.

179 While this work is concurrent to our own, we differentiate ourselves as **our method preserves the**
 180 **original FM structure. We maintain the ODE form of the sampler and the linear interpolant,**
 181 **and introduce diversity solely by injecting divergence-free velocity perturbations during sam-**
 182 **pling.** This guarantees that probability mass is conserved through the continuity equation, while
 183 allowing geometrically distinct trajectories for verifier-guided search.

184 3.4 Noise Injection while Preserving the Continuity Equation

185 A distinct thread of work aims to inject stochasticity during inference without violating the underlying
 186 continuity equation that defines the model’s evolution. This objective is shared by our divergence-free
 187 perturbation strategy, but is also addressed through Langevin-style diffusion sampling schemes.

188 The most widely adopted approach in this category is the so-called "churn" strategy ?, used ex-
 189 tensively in diffusion models. The idea is to introduce noise in a way that preserves the marginal
 190 densities $p_t(x)$ *in expectation*, rather than at the level of individual sample trajectories. Specifically,
 191 the sampler follows an SDE of the form:

$$dx_t = u_t(x_t) - \beta(t) \sigma^2(t) \nabla \log p_t(x_t) dt + \sqrt{2\beta(t)} \sigma(t) dW_t,$$

192 where $\beta(t)$ is a user-defined schedule that governs the amount of stochasticity injected at each
 193 timestep. The drift term and noise are precisely scaled so that they cancel out in the associated
 194 Fokker–Planck equation, thereby preserving p_t .

195 Unlike our method, which satisfies the continuity equation path-wise by ensuring $\nabla_x \cdot (p_t w_t) = 0$ at
 196 every point, the EDM method preserves the distribution only in expectation over trajectories. As a
 197 result, individual sample paths do not strictly conserve mass. This difference becomes critical when
 198 injecting higher magnitudes of noise: EDM samplers tend to lose image quality earlier than our
 199 divergence-free ODE approach, as shown in our empirical study in Section ??.

200 We include this method as a strong baseline in our experiments for its conceptual proximity to our
 201 work.

202 4 Inference Time Scaling for Flow Matching while Preserving the ODE

203 Our investigation of inference-time scaling for flow matching models operates along two independent
 204 axes: **noise injection methods** (Sections 4.1–4.3) that introduce stochasticity during sampling, and
 205 **inference algorithms** (Section 4.4) that leverage this stochasticity for reward-guided search. This
 206 separation allows for flexible combinations of noise types with different search strategies, enabling
 207 exploration of the compute-quality tradeoff space.

208 4.1 Divergence-Free Noise

209 To enhance sample diversity without altering the trained density path $p_t(x)$, we inject a divergence-
 210 free perturbation $w_t(x)$ into the learned velocity field $u_t(x)$ at inference time when simulating the
 211 ODE. Specifically, we modify the velocity field to $u_t(x) + \lambda w_t(x)$, where λ is a scalar hyperparameter
 212 controlling the perturbation magnitude. Crucially, from the flow matching formulation, each modified
 213 vector field $u_t(x) + \lambda w_t(x)$ represents an equally valid "randomized" ODE as defined by the continuity

equation, since the divergence-free property ensures that $\nabla_x \cdot (p_t(x)w_t(x)) = 0$. This theoretical foundation allows us to introduce stochasticity for improved sample diversity while preserving the mathematical structure of the original flow.

4.2 Proof: Adding Divergence-Free Noise at Inference Preserves the Continuity Equation and the Probability Flow ODE

In flow-matching models, the learned velocity $u_t(x)$ satisfies the continuity equation:

$$\partial_t p_t(x) + \nabla_x \cdot (p_t(x) u_t(x)) = 0, \quad (\text{CE})$$

ensuring that the evolution of densities $\{p_t\}$ is consistent with an underlying deterministic ODE. We aim to add a perturbation $w_t(x)$ (e.g. divergence-free "swirl") to improve diversity at inference time, and want to confirm that the modified flow still respects the continuity equation. This is essential for ensuring that samples remain consistent with the trained marginal densities $p_t(x)$.

Proposition. Let p_t and u_t satisfy (CE). If a vector field w_t satisfies

$$\nabla_x \cdot (p_t w_t) = 0 \quad \text{for all } t \in [0, 1], \quad (1)$$

then the modified drift $\tilde{u}_t := u_t + \lambda w_t$, where $\lambda \in \mathbb{R}$ is a scalar hyperparameter controlling the amount of added noise, yields the same continuity equation:

$$\partial_t p_t + \nabla_x \cdot (p_t \tilde{u}_t) = 0.$$

Proof. Plug $\tilde{u}_t = u_t + \lambda w_t$ into (CE):

$$\partial_t p_t + \nabla_x \cdot (p_t(u_t + \lambda w_t)) = \underbrace{[\partial_t p_t + \nabla_x \cdot (p_t u_t)]}_{=0 \text{ by (CE)}} + \lambda \nabla_x \cdot (p_t w_t).$$

The second term vanishes by assumption (1), hence the entire expression equals zero and the modified drift preserves the same continuity equation. ■

Remark (local criterion). Using the identity

$$\nabla_x \cdot (p_t w_t) = p_t (\nabla_x \cdot w_t + s_t \cdot w_t), \quad s_t := \nabla_x \log p_t,$$

we see that condition (1) is equivalent to the pointwise or expected constraint:

$$\boxed{\nabla_x \cdot w_t + s_t \cdot w_t = 0.}$$

In practice, this is satisfied (in expectation) by choosing

$$w_t(x) = (I - \hat{s}_t \hat{s}_t^\top) \varepsilon, \quad \hat{s}_t = \frac{s_t}{\|s_t\|}, \quad \varepsilon \sim \mathcal{N}(0, I),$$

which projects Gaussian noise onto the subspace orthogonal to the score direction. This guarantees that $s_t \cdot w_t = 0$ exactly, while $\nabla_x \cdot w_t = 0$ holds in expectation because w_t is a linear transformation of ε with coefficients that are independent of the spatial variable x . Thus, $\nabla_x w_t = 0$ and the continuity equation is preserved.

4.3 Baseline Noise Injection Methods

While our primary contribution focuses on divergence-free noise injection, we also evaluate alternative approaches for introducing stochasticity into flow matching sampling as baseline comparisons:

Standard SDE Sampling: The most straightforward approach adds Gaussian noise directly to the velocity field: $dx_t = u_t(x_t) dt + \sigma dW_t$, where σ controls the noise magnitude. While simple to implement, this approach does not preserve the continuity equation and can lead to rapid degradation of sample quality as noise levels increase.

Score SDE Sampling: Following the stochastic interpolants framework ?, we can convert the learned velocity field to a score-based SDE formulation with drift $f_t(x_t) = u_t(x_t) - \frac{g_t^2}{2} \nabla \log p_t(x_t)$ that incorporates the analytically computed score function. This method preserves the marginal distributions in expectation but may lose trajectory-level mass conservation.

These baseline methods enable us to demonstrate the advantages of our divergence-free approach in maintaining sample quality while introducing beneficial stochasticity.

4.4 Empirical Demonstration of Diversity Without Reducing Quality

To validate this approach, we evaluate how the injection of noise at inference time affects generation quality across several methods. We compare our divergence-free ODE method to a deterministic baseline, a simple Euler-Maruyama SDE, and a Langevin-style stochastic sampler derived from EDM ?. All experiments are conducted using the pretrained SiT-XL/2 flow-matching model on ImageNet 256×256. For each configuration, we generate 1,000 samples and compute Fréchet Inception Distance (FID), Inception Score (IS), and sample diversity.

Setup. The x -axis in all figures represents the average magnitude of injected noise per step, normalized to the average magnitude of the velocity field $u_t(x)$. The y -axis reports sample diversity, FID, or IS. We compare four sampling strategies:

- **ODE-divfree:** Our proposed method, which injects divergence-free perturbations $w_t(x)$ that preserve the continuity equation at the level of individual trajectories:

$$dx_t = (u_t(x) + \lambda w_t(x)) dt,$$

where λ controls the noise scale.

- **SDE:** Samples are generated from a simple Euler-Maruyama SDE:

$$dx_t = u_t(x) dt + \sigma dW_t,$$

where $dW_t \sim \mathcal{N}(0, dt)$, and σ scales the noise.

- **EDM-SDE:** A stochastic method that adjusts both drift and diffusion to preserve $p_t(x)$ in expectation:

$$dx_t = u_t(x) dt - \beta(t) \sigma^2(t) \nabla \log p_t(x) dt + \sqrt{2\beta(t)} \sigma(t) dW_t,$$

where $\beta(t)$ is a user-defined schedule. See Related Work for details (Section ??).

- **Score-SDE:** Uses the Score SDE formulation with analytically computed score functions from the stochastic interpolants framework.

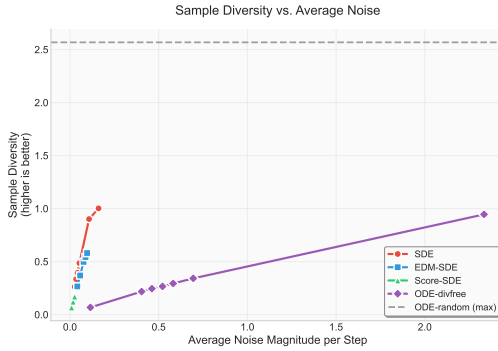


Figure 1: Sample diversity across increasing noise levels. Higher is better.

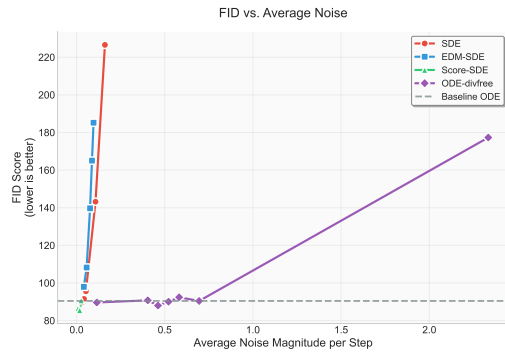


Figure 2: FID across increasing noise levels. Lower is better.

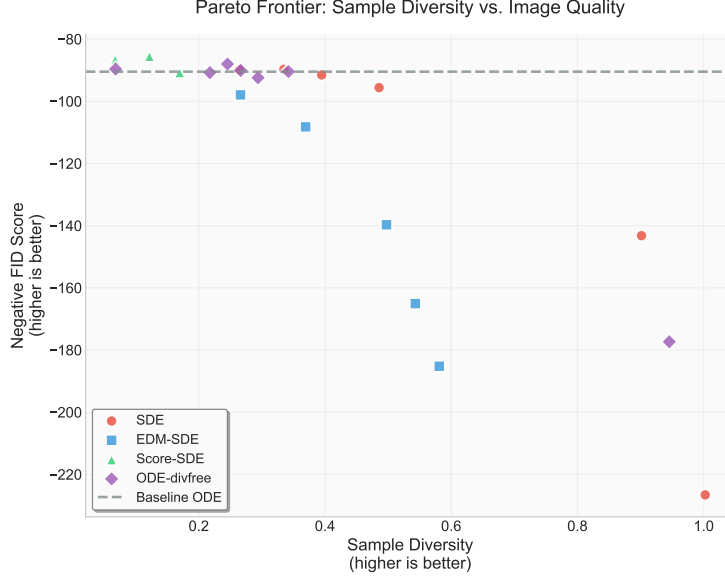


Figure 3: Pareto frontier analysis showing the trade-off between sample diversity and image quality (negative FID, where higher is better for both axes). Our divergence-free method achieves favorable positions on the frontier, maintaining high image quality while enabling substantial diversity gains compared to other stochastic sampling approaches.

Results. As shown in the figures, our divergence-free ODE approach allows substantially higher levels of noise to be injected without degrading sample quality. Sample diversity increases smoothly with noise level (Fig. ??), while FID remains stable across a wide noise range (Fig. ??). The Pareto frontier analysis (Fig. ??) demonstrates that our method achieves superior trade-offs between diversity and quality compared to alternative stochastic approaches. In contrast, standard SDE sampling degrades rapidly as noise increases. EDM sampling performs better than the naïve SDE—consistent with its theoretical guarantee of preserving $p_t(x)$ in expectation—but still deteriorates earlier than our method. Similar analysis for Inception Score shows comparable results (see Appendix). This suggests that satisfying the continuity equation per trajectory, as our approach does, offers stronger robustness to stochasticity during sampling.

4.5 Inference Algorithms

Given any of the above noise injection methods, we employ different search strategies to leverage the introduced stochasticity for improved sample quality. We present multiple inference algorithms that operate independently of the chosen noise injection method, enabling us to flexibly combine noise strategies with these sampling strategies. Each algorithm requires a verifier function $r(\cdot)$, which scores complete samples at $t = 1$. See Appendix A for implementation details for each algorithm

Best-of-N Sampling: The simplest inference scaling approach generates N independent samples from different initial conditions and take the top- K samples. This method provides broad exploration of the initial condition space and scales linearly with computational budget.

Path Exploration: Path exploration generates multiple trajectories from the same initial condition by introducing stochastic branching during the integration process. At predetermined timesteps, we create multiple branches by applying different noise realizations, simulate each branch to completion, score the final samples, and continue with the highest-scoring branches. This enables focused exploration around promising trajectories. **Computational Consideration:** Path exploration requires simulating intermediate samples forward to completion ($t = 1$) for scoring, which adds significant computational overhead beyond the reported compute budgets. This look-ahead scoring ensures accurate quality assessment but substantially increases the actual computational cost compared to other methods.

298 **Best-of-N + Path Exploration:** This two-stage approach combines the benefits of initial condition
 299 search with trajectory-level exploration. Stage 1 performs Best-of-N sampling to identify promising
 300 initial conditions. Stage 2 applies path exploration starting from the top- K initial conditions identified
 301 in Stage 1. This mirrors the exploration-exploitation tradeoff in reinforcement learning.

302 **Noise Search:** Noise search generates multiple samples from the same initial condition by applying
 303 different noise realizations throughout the entire sampling process. Unlike path exploration, this
 304 method does not require intermediate scoring and branching, making it computationally efficient
 305 while still enabling local exploration around a given starting point. In practice we extend Noise
 306 search to multiple rounds, where each round starts progressively closer to $t = 1$ using intermediate
 307 samples that correspond to the k top samples from the previous round, allowing us to further explore
 308 high-scoring regions.

309 **Best-of-N + Noise Search:** Similar to the path exploration variant, this two-stage method first
 310 identifies promising initial conditions via Best-of-N sampling, then applies noise search around each
 311 selected starting point. This approach provides both broad initial exploration and focused local search
 312 while maintaining computational efficiency.

313 **Connection to Particle Sampling:** Our inference algorithms share conceptual similarities with
 314 particle sampling methods, where multiple "particles" (sample trajectories) evolve through the state
 315 space and are periodically reweighted or resampled based on fitness criteria. The path exploration
 316 method particularly resembles particle filtering approaches, where particles branch and are selected
 317 based on likelihood estimates. However, our methods operate specifically within the flow matching
 318 framework and leverage verifier functions rather than likelihood-based reweighting.

319 5 Experiments

320 We evaluate the effectiveness of our inference-time scaling framework across two domains: ImageNet
 321 256×256 image generation and protein structure design. Our experiments systematically explore
 322 combinations of noise injection methods with inference algorithms to demonstrate the versatility and
 323 effectiveness of our approach.

324 5.1 Evaluated Method Combinations

325 We evaluate seven distinct combinations of noise injection methods and inference algorithms across
 326 both experimental domains. Each method is assigned a shorthand name for consistent reference
 327 throughout our results:

- 328 • **Best-of-N:** Random search over initial conditions using deterministic ODE sampling
- 329 • **Path-DivFree:** Path exploration with divergence-free noise injection (Algorithm ??)
- 330 • **Path-SDE:** Path exploration with standard SDE noise injection
- 331 • **Path-ScoreSDE:** Path exploration with score SDE noise injection
- 332 • **BestN+Path-DivFree:** Two-stage Best-of-N followed by divergence-free path exploration
 333 (Algorithm ??)
- 334 • **NoiseSearch-3R-DivFree:** Multi-round (3 rounds) noise search with divergence-free noise
 335 injection (Algorithm ??)
- 336 • **BestN+NoiseSearch-3R-DivFree:** Two-stage Best-of-N followed by 3-round divergence-
 337 free noise search (Algorithm ??)

338 All methods are evaluated across compute budgets of 1×, 2×, 4×, and 8×, where the compute factor
 339 corresponds to the number of final samples retained. The 1× baseline uses deterministic ODE
 340 sampling with 1,000 samples for ImageNet and 64 samples for FoldFlow.

341 5.2 ImageNet 256×256: Inference-Time Scaling with SiT-XL/2

342 We demonstrate the effectiveness of our inference-time scaling approach on ImageNet 256×256 using
 343 the pretrained SiT-XL/2 flow-matching model.

5.2.1 Experimental Setup

We evaluate all seven method combinations described in Section 5.1 on ImageNet 256×256 using the pretrained SiT-XL/2 flow-matching model. All methods generate 1,000 samples for robust statistical evaluation.

5.2.2 Evaluation Metrics and Verifiers

We evaluate performance using four metrics: FID (lower is better), Inception Score (higher is better), DINO Top-1 accuracy (higher is better), and DINO Top-5 accuracy (higher is better). Inception Score measures the quality and diversity of generated images based on a pretrained ImageNet classifier ?. DINO scores evaluate semantic quality using self-supervised vision transformer features that capture richer visual representations ?. FID computes the Fréchet distance between generated and real image feature distributions, providing a comprehensive measure of sample quality and diversity ?. Each experiment uses two different verifier functions for sample selection: Inception Score-based scoring and DINO-based scoring.

5.2.3 Results: Inception Score-Guided Scaling

Figure ?? shows the results when using Inception Score as the verifier for sample selection. We report Inception Score and DINO Top-1 accuracy as the primary metrics, as these capture different aspects of sample quality.

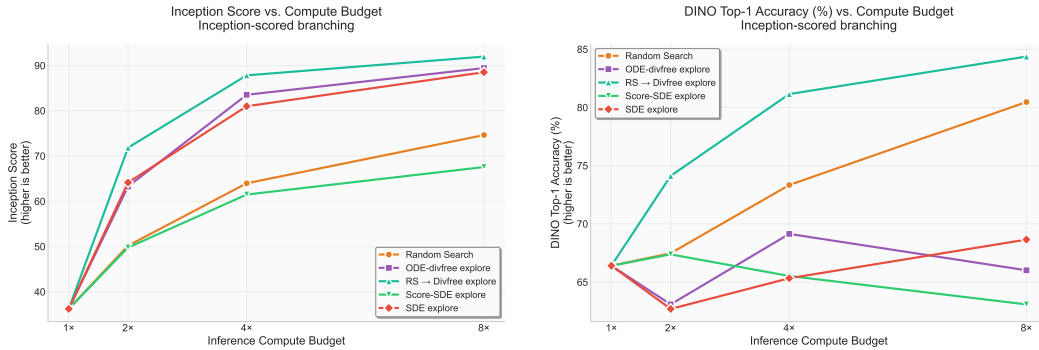


Figure 4: Inference-time scaling results using Inception Score-guided selection. Left: Inception Score vs. compute budget. Right: DINO Top-1 accuracy vs. compute budget.

The RS → Divfree explore method demonstrates the strongest performance, achieving the highest Inception Scores at 4x and 8x compute budgets. This two-stage approach effectively combines the benefits of searching over initial conditions with path-space exploration. The ODE-divfree explore method shows steady improvement but is limited by starting from a single initial condition. Random search provides consistent but modest gains, while the SDE-based methods show more variable performance.

5.2.4 Results: DINO-Guided Scaling

Figure ?? presents results when using DINO-based scoring for sample selection. We focus on FID, Inception Score, and DINO Top-1 accuracy as key performance indicators.

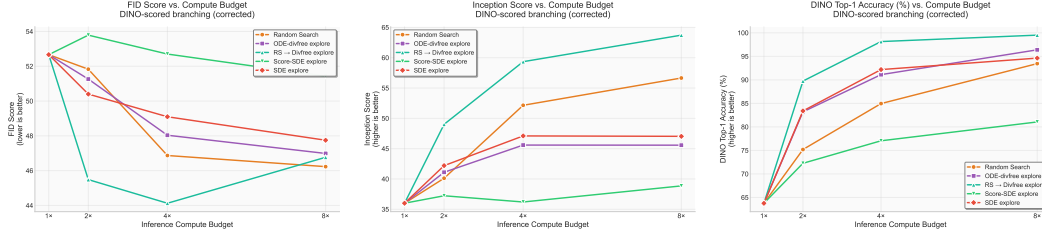


Figure 5: Inference-time scaling results using DINO-guided selection. Left: FID vs. compute budget. Center: Inception Score vs. compute budget. Right: DINO Top-1 accuracy vs. compute budget.

Under DINO-guided selection, the RS → Divfree explore method again shows superior performance, achieving the best FID scores and DINO Top-1 accuracy at higher compute budgets. Notably, DINO-guided selection produces more substantial improvements in DINO Top-1 accuracy compared to Inception-guided selection, demonstrating the importance of verifier-method alignment. The ODE-divfree explore method shows consistent improvement across all metrics, while random search and SDE methods provide more modest gains.

5.3 FoldFlow: Protein Design

We evaluate our inference-time scaling methods on the task of unconditional protein backbone generation using the pretrained FoldFlow2 model. FoldFlow is a flow matching model that generates protein backbone structures by learning the conditional distribution of protein coordinates given sequence information. This provides an important validation of our approach in a scientific domain with distinct evaluation metrics.

5.3.1 Experimental Setup

We apply the same inference-time scaling framework from our ImageNet experiments to protein generation, evaluating all seven method combinations described in Section 5.1. We generate 64 protein samples of length 100 residues per configuration across compute budgets of 2x, 4x, and 8x. All methods are compared against the standard deterministic ODE sampling baseline (1x compute). During inference, we use self-consistency template modeling score (scTM-score) as the verifier function for sample selection.

5.3.2 Evaluation

Following the FoldFlow evaluation framework, the self-consistency framework measures protein designability by testing structural integrity through a complete folding-refolding cycle: FoldFlow generates a backbone structure, ProteinMPNN predicts 8 amino acid sequences for that structure, ESMFold refolds these sequences into 3D structures, and structural alignment metrics (self-consistency Root Mean Square Deviation, scRMSD, and scTM-score) compare the refolded structures to the original. This tests whether generated structures represent realistic, foldable proteins that follow natural sequence-structure relationships. For both scRMSD and scTM-score, we report both the average scores across the set of sampled proteins per method as well as the designability criteria from the FoldFlow framework.

Designability Thresholds: We measure the percentage of samples with scRMSD below the threshold of 2.0Å, as described in FoldFlow’s evaluation. We additionally use 1.5Å and 1.0Å as thresholds, representing increasingly stringent designability criteria

5.3.3 Results

Figure ?? shows the mean self-consistency TM-scores across different compute budgets for each inference-time scaling method. Additionally, we provide comprehensive evaluation through mean self-consistency RMSD scores and designability threshold analyses.

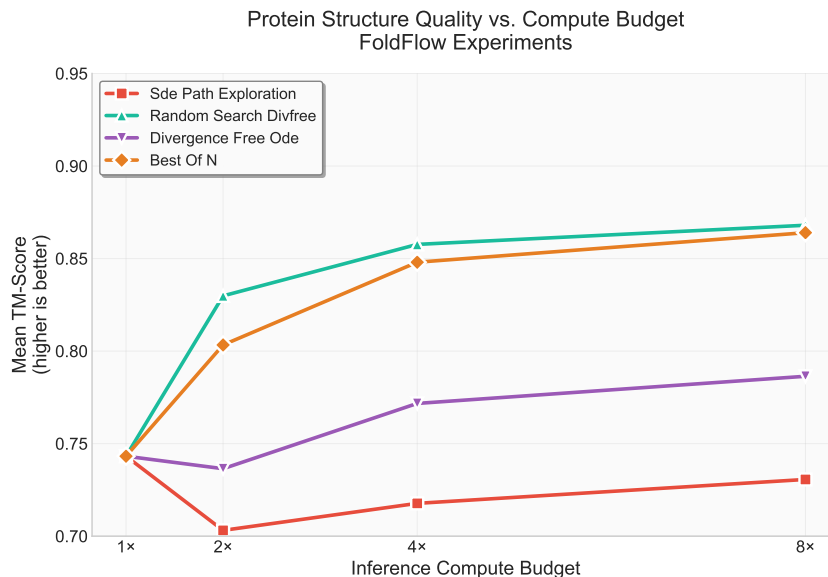


Figure 6: Protein structure quality (self-consistency TM-score) vs. compute budget for FoldFlow experiments. All methods start from the same baseline (standard ODE sampling at 1× compute) and demonstrate how additional inference compute improves protein designability.

The evaluation includes three complementary perspectives on protein generation quality. Mean self-consistency TM-score provides a continuous measure of structural similarity between generated and refolded structures. Mean self-consistency RMSD offers an alternative continuous metric focusing on coordinate-level accuracy. The designability threshold analysis reports the percentage of generated proteins meeting increasingly stringent criteria: $\text{scRMSD} < 2.0\text{\AA}$, $\text{scRMSD} < 1.5\text{\AA}$, and $\text{scRMSD} < 1.0\text{\AA}$. These threshold-based metrics align closely with the FoldFlow evaluation framework and provide insight into the practical utility of generated structures, as proteins meeting tighter thresholds are more likely to represent biologically viable designs.

6 Conclusion

We have introduced the first comprehensive framework for inference-time scaling in flow matching models, addressing a critical gap in continuous-time generative modeling. Our approach preserves the mathematical foundations of flow matching while enabling principled exploration of alternative generation trajectories through divergence-free velocity perturbations.

Key Achievements: Our theoretical analysis proves that divergence-free perturbations maintain the continuity equation at the trajectory level, providing stronger guarantees than methods that preserve distributions only in expectation. This mathematical rigor enables robust performance across diverse noise levels and ensures compatibility with the linear interpolation paths that define flow matching. Our experimental validation demonstrates substantial improvements across two distinct domains: ImageNet image generation and protein structure design, with our two-stage approach consistently achieving the best performance.

Methodological Insights: Our work reveals complementary strengths of different scaling strategies. Random search provides broad exploration particularly valuable for discrete relationships (e.g., protein sequence-structure mappings), while divergence-free path exploration enables focused exploitation around promising regions. The two-stage combination captures benefits of both approaches, suggesting that inference-time scaling benefits from balancing exploration and exploitation.

Practical Impact: The training-free nature of our approach makes it immediately applicable to existing flow matching models across domains. For protein design, improvements from 0.743 to 0.868 TM-score represent meaningful advances in structural quality with implications for drug discovery

434 and protein engineering. For image generation, our method provides a scalable path to higher quality
 435 samples without architectural modifications or retraining.

436 **Limitations and Future Work:** While our approach shows diminishing returns at higher compute
 437 budgets (8×), suggesting practical limits to inference-time scaling, even modest compute increases
 438 (2×-4×) yield substantial quality improvements. Future work could explore adaptive scaling strategies
 439 that optimize compute allocation based on sample complexity, integration with more sophisticated
 440 verifiers for domain-specific applications, and extension to other continuous-time generative models
 441 beyond flow matching.

442 **Broader Implications:** Our work positions inference-time scaling as a fundamental capability
 443 for continuous generative models, complementing training-time scaling. As generative models
 444 become increasingly important for scientific discovery and creative applications, the ability to trade
 445 computational resources for sample quality at inference time provides practitioners with valuable
 446 flexibility. This paradigm shift from "one model, one quality level" to "adjustable quality through
 447 inference compute" opens new possibilities for interactive applications, quality-critical deployments,
 448 and resource-constrained scenarios.

449 Looking forward, we envision inference-time scaling becoming a standard tool in the continuous
 450 generative modeling toolkit, enabling researchers and practitioners to extract maximum value from
 451 trained models while adapting to diverse quality requirements and computational constraints.

452 A Inference Algorithm Implementations

453 This section provides detailed algorithmic descriptions for all five inference methods introduced in
 454 Section 4.4.

455 A.1 Best-of-N Sampling

Algorithm 1 Best-of-N Sampling

Require: Flow matching model v_θ , verifier function $r(\cdot)$, number of samples N , number to retain K

Ensure: Top- K generated samples

1: **for** $i = 1$ to N **do**

2: $x_0^{(i)} \sim \pi_{\text{ref}}$

▷ Sample initial condition

3: $x_1^{(i)} \leftarrow \text{ODESolve}(v_\theta, x_0^{(i)})$

▷ Deterministic sampling

4: $s^{(i)} \leftarrow r(x_1^{(i)})$

▷ Compute verifier score

5: **end for**

6: **return** Top- K samples by score $\{s^{(i)}\}$

Algorithm 2 Path Exploration with Noise Injection

Require: Flow matching model v_θ , initial noise x_0 , noise injection method $\mathcal{N}(\cdot)$, verifier function $r(\cdot)$, branching schedule $\{t_b\}$, branches per timestep B **Ensure:** Generated sample x_1

```

1: trajectories  $\leftarrow \{x_0\}$  ▷ Initialize with single trajectory
2: for each branching timestep  $t_b$  do
3:   new_trajectories  $\leftarrow \{\}$ 
4:   for each trajectory  $x_{t_b}$  in trajectories do
5:     for  $j = 1$  to  $B$  do
6:        $x_1^{(j)} \leftarrow \text{NoisySample}(v_\theta, x_{t_b}, t_b, \mathcal{N})$  ▷ Sample to completion
7:        $s^{(j)} \leftarrow r(x_1^{(j)})$  ▷ Score final sample
8:       Add  $(x_{t_b}, s^{(j)})$  to new_trajectories
9:     end for
10:   end for
11:   trajectories  $\leftarrow$  Top scoring trajectories from new_trajectories
12: end for
13: return Best trajectory from final sampling round

```

Algorithm 3 Two-Stage Best-of-N + Path Exploration

Require: Flow matching model v_θ , verifier function $r(\cdot)$, stage 1 samples N , top candidates K , noise injection method $\mathcal{N}(\cdot)$ **Ensure:** Top- K generated samples

```

1: Stage 1:  $\{x_0^*\} \leftarrow \text{BestOfN}(v_\theta, r, N, K)$  ▷ Algorithm ??
2: Stage 2: results  $\leftarrow \{\}$ 
3: for each  $x_0^{(i)} \in \{x_0^*\}$  do
4:    $x_1^{(i)} \leftarrow \text{PathExploration}(v_\theta, x_0^{(i)}, \mathcal{N}, r)$  ▷ Algorithm ??
5:   Add  $x_1^{(i)}$  to results
6: end for
7: return Top- $K$  samples from results by verifier score

```

Algorithm 4 Multi-Round Noise Search

Require: Flow matching model v_θ , initial noise x_0 , noise injection method $\mathcal{N}(\cdot)$, verifier function $r(\cdot)$, rounds R , samples per round N , candidates per round K **Ensure:** Generated sample x_1

```

1: candidates  $\leftarrow \{x_0\}$  ▷ Initialize with single candidate
2: for round  $i = 1$  to  $R$  do
3:    $t_{\text{start}} \leftarrow (i - 1)/R$  ▷ Starting timestep for round
4:   round_samples  $\leftarrow \{\}$ 
5:   for each candidate  $x_{t_{\text{start}}}$  in candidates do
6:     for  $j = 1$  to  $N$  do
7:        $x_1^{(j)} \leftarrow \text{NoisySample}(v_\theta, x_{t_{\text{start}}}, t_{\text{start}}, \mathcal{N})$ 
8:        $s^{(j)} \leftarrow r(x_1^{(j)})$ 
9:       Add  $(x_{t_{\text{start}}}, x_1^{(j)}, s^{(j)})$  to round_samples
10:    end for
11:  end for
12:  if  $i < R$  then
13:    candidates  $\leftarrow$  Top- $K$  starting points by final score
14:  else
15:    return Best final sample by verifier score
16:  end if
17: end for

```

Algorithm 5 Two-Stage Best-of-N + Noise Search

Require: Flow matching model v_θ , verifier function $r(\cdot)$, stage 1 samples N , top candidates K , noise injection method $\mathcal{N}(\cdot)$, search rounds R **Ensure:** Top- K generated samples

```

1: Stage 1:  $\{x_0^*\} \leftarrow \text{BestOfN}(v_\theta, r, N, K)$  ▷ Algorithm ??
2: Stage 2: results  $\leftarrow \{\}$ 
3: for each  $x_0^{(i)} \in \{x_0^*\}$  do
4:    $x_1^{(i)} \leftarrow \text{NoiseSearch}(v_\theta, x_0^{(i)}, \mathcal{N}, r, R)$  ▷ Algorithm ??
5:   Add  $x_1^{(i)}$  to results
6: end for
7: return Top- $K$  samples from results by verifier score

```

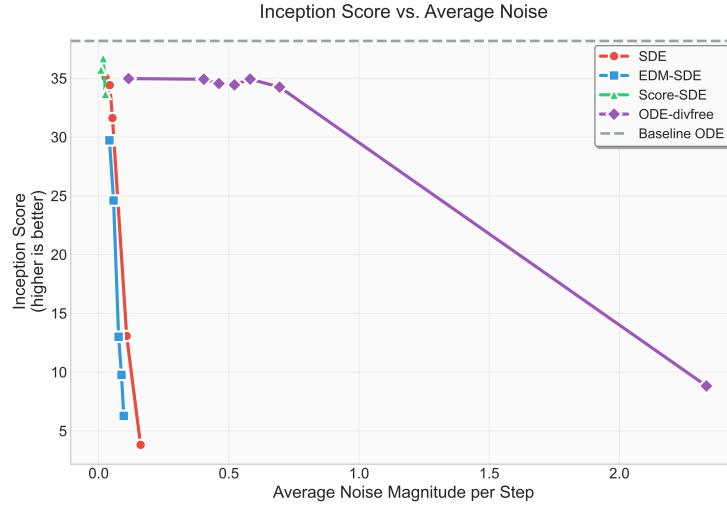


Figure 7: Inception Scores across increasing noise levels. Higher is better. This complements the diversity and FID analysis shown in the main text, demonstrating similar robustness patterns for our divergence-free method.

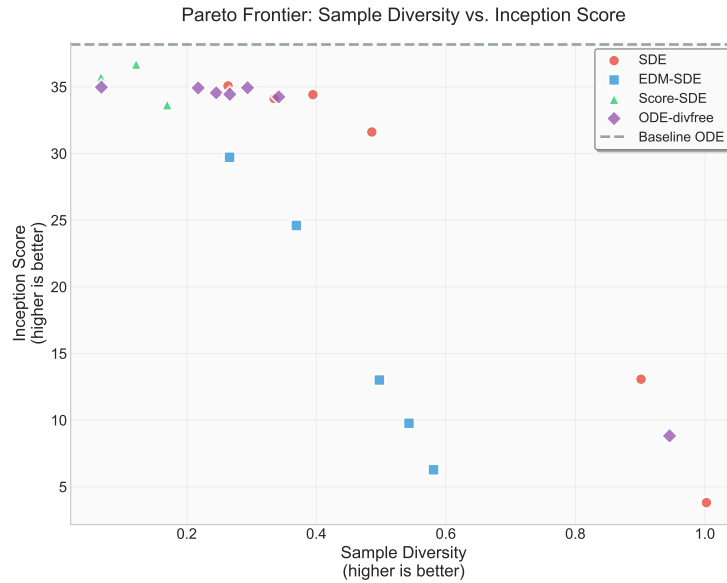


Figure 8: Pareto frontier analysis showing the trade-off between sample diversity and Inception Score (higher is better for both axes). This complements the FID-based Pareto analysis in the main text.

464 B.2 Complete Inference-Time Scaling Results

465 B.2.1 Inception Score-Guided Scaling (All Metrics)

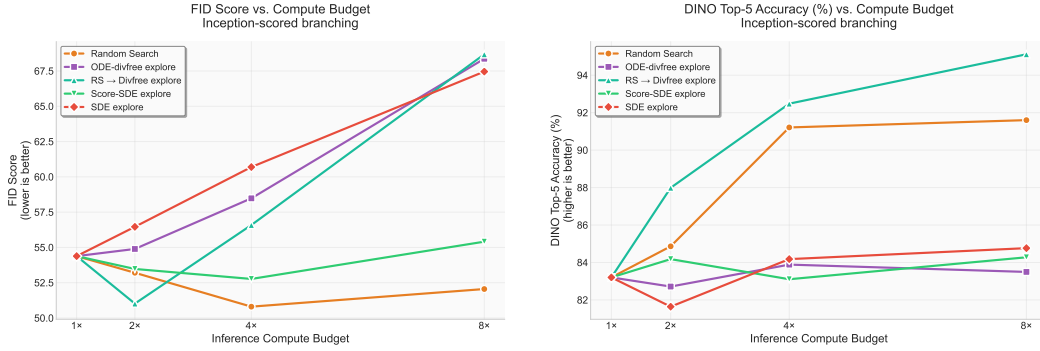


Figure 9: Additional Inception Score-guided scaling results. Left: FID vs. compute budget. Right: DINO Top-5 accuracy vs. compute budget.

466 B.2.2 DINO-Guided Scaling (All Metrics)

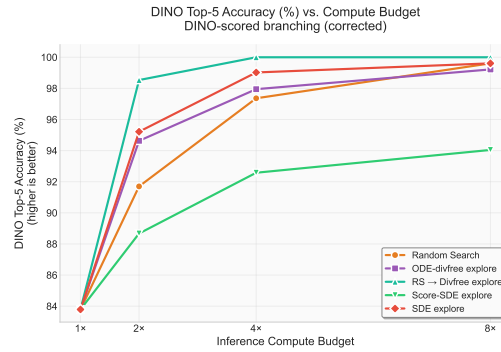


Figure 10: DINO Top-5 accuracy vs. compute budget for DINO-guided scaling experiments.

467 C Implementation Details

468 C.1 Look-Ahead Scoring for Branching Methods

469 Our path exploration methods (ODE-divfree explore, SDE explore, and the two-stage approach)
 470 employ a look-ahead scoring strategy for branch selection. At each branching point, multiple
 471 candidate trajectories are simulated to completion, and each generated sample is evaluated using the
 472 verifier function (Inception Score, DINO score, or TM-score). The scores are then used to select the
 473 most promising branches for continued exploration.

474 Alternative approaches exist for scoring intermediate branches, such as evaluating partially-noised
 475 samples at branching points or using a fixed computational budget for partial trajectory simulation.
 476 However, these methods introduce additional complexity in score interpretation and may not provide
 477 reliable quality estimates for incomplete generations. For simplicity and to ensure fair comparison
 478 across methods, we use complete trajectory simulation for all scoring decisions.

479 It is important to note that we do not include the computational cost of look-ahead scoring when
 480 reporting computational complexity or compute budgets. The reported compute factors (1x, 2x, 4x,
 481 8x) reflect only the final number of samples retained, not the intermediate computational overhead
 482 required for branch evaluation. This choice enables cleaner comparison across methods while
 483 acknowledging that practical deployment would need to account for this additional computational
 484 cost.

C.2 Experimental Hyperparameters

Table ?? summarizes the key hyperparameters used across all experiments. These values were selected based on preliminary experiments to balance sample quality and computational efficiency.

Parameter	ImageNet	FoldFlow
Model	SiT-XL/2	FoldFlow
Number of samples	1,024	64
Integration timesteps	20	50
Protein length	N/A	50 residues
Divergence-free methods		
Noise scale (λ)	0.35	0.6
Branch interval	N/A	0.1 (every 5 steps)
SDE methods		
Noise scale (σ)	0.1	0.3
Branch interval	N/A	0.1 (every 5 steps)
Compute budgets	1×, 2×, 4×, 8×	1×, 2×, 4×, 8×
Verifiers	Inception, DINO	TM-score

Table 1: Experimental hyperparameters for ImageNet and FoldFlow experiments. Branch interval controls the frequency of branching operations during path exploration.

C.3 Domain-Specific Implementation Details

C.3.1 ImageNet Experiments

For ImageNet experiments, we use the pretrained SiT-XL/2 model with 20 integration timesteps, following the standard configuration for efficient sampling. The divergence-free noise scale of $\lambda = 0.35$ was selected to provide substantial diversity without degrading sample quality, as validated in our noise study (Section ??). The SDE noise scale of $\sigma = 0.1$ provides a conservative baseline for stochastic sampling comparison.

All ImageNet experiments generate 1,024 samples per configuration to ensure robust statistical evaluation across the four metrics (FID, Inception Score, DINO Top-1, DINO Top-5). Verifier-based selection uses either Inception Score or DINO features, enabling comparison of verifier-method alignment effects.

C.3.2 FoldFlow Experiments

FoldFlow experiments focus on proteins of length 50 residues, providing a manageable complexity for systematic evaluation while representing realistic protein design scenarios. The model uses 50 integration timesteps by default, but performing branching operations at every timestep would be computationally prohibitive and unnecessary.

The branch interval parameter controls branching frequency: a value of 0.1 indicates branching every 10% of the trajectory. With 50 timesteps, this corresponds to branching every 5 timesteps, resulting in 10 total branching operations per trajectory. This provides sufficient exploration opportunities while maintaining computational tractability.

The higher noise scales for FoldFlow ($\lambda = 0.6$ for divergence-free, $\sigma = 0.3$ for SDE) reflect the different data characteristics and model sensitivities compared to image generation. These values were calibrated to provide meaningful exploration while preserving protein structural validity.

The smaller sample size (64 proteins) reflects the computational cost of protein generation and evaluation, while still providing sufficient data for reliable TM-score estimation and method comparison.

D Additional Theoretical Analysis

TODO: Add extended theoretical analysis and proofs.