# Optimal Stepsize for Diffusion Sampling

Jianning Pei[1,2], Han Hu[2], Shuyang Gu[2†]

[1]University Chinese Academic of Science   [2]Tencent Hunyuan Research

jianningpei22@mails.ucas.ac.cn, cientgu@tencent.com

Figure 1. Flux sampling results using different stepsize schedules. Left: Original sampling result using 100 steps. Middle: Optimal stepsize sampling result within 10 steps. Right: Naively reducing sampling steps to 10.

## Abstract

*Diffusion models achieve remarkable generation quality but suffer from computational intensive sampling due to suboptimal step discretization. While existing works focus on optimizing denoising directions, we address the principled design of stepsize schedules. This paper proposes Optimal Stepsize Distillation, a dynamic programming framework that extracts theoretically optimal schedules by distilling knowledge from reference trajectories. By reformulating stepsize optimization as recursive error minimization, our method guarantees global discretization bounds through optimal substructure exploitation. Crucially, the distilled schedules demonstrate strong robustness across architectures, ODE solvers, and noise schedules. Experiments show 10× accelerated text-to-image generation while preserving **99.4%** performance on GenEval. Our code is available at* [https://github.com/bebebe666/OptimalSteps](https://github.com/bebebe666/OptimalSteps).

## 1. Introduction

Score-based generative models [28, 29] have emerged as a groundbreaking paradigm in generative modeling, achieving state-of-the-art results across diverse domains. Unlike traditional approaches that directly model complex data distributions, these methods decompose the target distribution into a sequence of conditional distributions structured as a Markov chain. While this decomposition enables tractable training through iterative noise prediction, it inherently necessitates computational expensive sampling procedures.

From a theoretical perspective, diffusion models approximate the data generation process as a continuous transformation from noise to structured data, parameterized by an infinite series of conditional distributions mapping. But in practice, computational constraints require discretizing this continuous trajectory into a finite sequence of steps, introducing a critical trade-off between sampling efficiency and approximation fidelity.

The discretization includes two key factors: the update direction (determined by the score function) and the stepsize (controlling the progression towards the target distribution), as shown in Figure 2. While significant attention has been devoted to optimizing update directions, including higher-order solvers [18, 19] and corrected network predictions [38], the principled design of stepsize remains underexplored. Existing methods predominantly rely on heuristic schedules, such as the uniform timesteps in DDIM [27] or the empirically designed schedule in EDM [10], lacking theoretical guarantees of optimality under constrained step budgets.
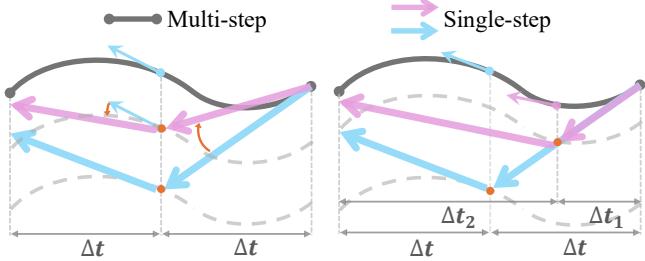
---

[†]Corresponding Author.

Figure 2. Two key factors in diffusion sampling: direction strategy(left) and stepsize strategy(right).

In this work, we propose a principled framework for searching the **O**ptimal **S**tepsize in diffusion **S**ampling (**OSS**), compatible with arbitrary direction strategies. Specifically, we regard the stepsize searching problem as a knowledge distillation problem, where a "student" sampling process with limited steps aims to approximate the results of a "teacher" sampling process with abundant steps. Crucially, we identify that this distillation objective exhibits recursive substructure—the optimal step schedule for $N$ steps inherently contain optimal sub-schedules for fewer steps. Exploiting this insight, we develop a dynamic programming algorithm that systematically derives theoretically optimal step sequences, ensuring maximal proximity to the teacher model's output under any pre-defined student sampling steps. It guarantees that the distilled schedule achieves the closest approximation to the teacher's continuous trajectory under computational constraints.

Importantly, our experimental analysis demonstrates the universal applicability and robustness of the proposed stepsize distillation framework. It exhibits consistent performance across various model architectures (e.g., U-Net [10] vs. Transformer-based [21]), ODE solver orders (1st to 3th-order [19]), noise schedules (linear [27]/Flow Matching [15, 17]/EDM [10]), and diverse denoising direction strategies [18, 19, 38]. Furthermore, when applying the optimal stepsize to text2image and text2video diffusion pipelines, our method achieves **10×** speedups while preserving **99.4%** of the teacher model's performance metrics on GenEval [5] benchmark. Such universal acceleration capability suggests broad practical potential for deploying latency-efficient diffusion models without sacrificing output quality.

Above all, our key contributions are as follows:

1. Theoretically Optimal Stepsize Distillation: We formulate stepsize optimization as a dynamic programming problem, proving that our solution can achieve global error minimization.

2. Architecture-Agnostic Robustness: We demonstrate that the distilled schedules generalize across datasets, noise schedules, and ODE solvers, outperforming heuristic baselines in diverse settings.

3. Efficient Adaptation: Our method enables lightweight schedule calibration across tasks, achieving 5-10× speedup with negligible performance degradation.

## 2. Related Work

### 2.1. Diffusion sampling direction strategy

Recent advances focused on correcting high-order direction for diffusion sampling acceleration. The earliest work [27] adopted first-order Euler for discretization diffusion ODEs. EDM [10] pioneers adopting Heun's[1] method for second-order gradient estimation between steps, reducing the Euler approximation error. PNDM [16] figured out the way of solving the denoising differential equation on manifold. AMED [40] utilized an additional network to predict the middle point of each step in the second order Heun solver. DPM-Solver [18] leveraged the semi-linear structure of diffusion ODEs, decomposing backward processes into linear/nonlinear components through Taylor expansion. Building on this, DPM-Solver++ [19] introduced a multistep gradient reuse and shifts schedule, further saving the computational efforts. UniPC [38] generalized these via a predictor-corrector framework using Lagrange polynomial interpolation, enabling arbitrary-order accuracy with minimal computational overhead. These innovations collectively established theoretical foundations for balancing discretization precision and computational overhead in diffusion sampling.

### 2.2. Diffusion sampling stepsize strategy

Traditional diffusion models typically employ predefined fixed-step strategies. For instance, DiT [21] and SDXL [22] adopt uniform timestep allocations, while EDM [10] and CM [30] rely on manually crafted schedules. These strategies often yield suboptimal results in few-step sampling scenarios, prompting recent studies to explore stepsize optimization.

Existing methods fall into two categories: training-free optimization and learning-based strategies. The former approaches, including DM [36], which minimizes inference error bounds via constrained trust-region methods, have a poor generalization ability since it shares the stepsize strategy across different diffusion networks and different solvers. Similarly, AYS [25] optimizes stepsize through KL divergence upper bounds (KLUB) but remains confined to SDE frameworks. [8] gets a better stepsize according to the similarity with the results of the high-order solver, which mixes the direction and stepsize of denoising. [35] selects the sampling steps based on dynamic programming and minimizes the sum of ELBO but incurs significant performance degradation. GITS [3] also employs dynamic programming to optimize stepsize but exhibits deviations from minimizing global error due to fixed cost matrices. The latter category, exemplified by LD3 [32], leverages a monotonic network for stepsize prediction, which is supervised by global errors. However, it requires additional high computational training. GGDM [34] adds trainable parameters

| Criterion | AYS [25] | GITS [3] | DM [36] | LD3 [32] | OSS |
|---|---|---|---|---|---|
| Training-free | ✓ | ✓ | ✓ | ✗ | ✓ |
| Stability | ✗ | ✓ | ✓ | ✓ | ✓ |
| Solver-aware | ✓ | ✗ | ✗ | ✓ | ✓ |
| Global error | ✗ | ✗ | ✗ | ✓ | ✓ |

Table 1. Related work comparison. Training-free denotes whether eliminating backpropagation during optimization. Stability denotes whether requires specialized stabilization techniques. Solver-aware means the algorithm integrates solver dynamics and network states. Global error means optimizing the final calibration error.

to the schedule and achieves high-quality samples through well-designed training. We summarize some recent methods in Table 1 from different perspectives.

## 3. Method

### 3.1. Diffusion sampling preliminaries

The continuous diffusion model progressively adds Gaussian noise to a data point $\boldsymbol{x_0}$ in the forward process:

$$q(\boldsymbol{x_t}|\boldsymbol{x_0}) = \mathcal{N}(\boldsymbol{x_t}|\alpha_t\boldsymbol{x_0}, \sigma_t^2\boldsymbol{I}) \tag{1}$$

where $\lambda_t = \frac{\alpha_t^2}{\sigma_t^2}$ denotes the Signal-to-Noise-ratio (SNR). The denoising process generates images from either Stochastic Differential Equation (SDE) solvers [7, 20] or Ordinary Differential Equation(ODE) solvers [17, 27]. The key difference between them is whether there is randomness in the generation process. Since ODE eliminates randomness, it often achieves significantly better results with fewer sampling steps. Therefore, this paper mainly focuses on the ODE solvers that follows:

$$d\boldsymbol{x_t} = [f_t\boldsymbol{x_t} - \frac{1}{2}g_t^2\nabla_{\boldsymbol{x}}logp_t(\boldsymbol{x_t})]dt \tag{2}$$

where $f_t = \frac{dlog\alpha_t}{dt}$ and $g_t^2 = \frac{d\sigma_t^2}{dt} - 2\frac{dlog\alpha_t}{dt}\sigma_t^2$. The generation process essentially involves solving the ODE and integrating the following expression from timestep $T$ to 0:

$$\boldsymbol{x_0} = \boldsymbol{x_T} + \int_T^0 f_t\boldsymbol{x_t} - \frac{1}{2}g_t^2\nabla_{\boldsymbol{x}}logp_t(\boldsymbol{x_t})dt \tag{3}$$

In practical applications, we usually discretize it into a finite number of $N$ steps to speed up the inference process:

$$\boldsymbol{x_0} = \boldsymbol{x_T} + \sum_{i=N}^{1} \boldsymbol{v_{\theta,i}}(t_{i-1} - t_i) \tag{4}$$

where $\boldsymbol{x_T} \sim N(0,1)$ is the initial noise, $\boldsymbol{\theta}$ denotes the parameter of denoising model, $\boldsymbol{v_{\theta,i}}$ denotes the optimization direction on step $i$. For simplicity, we use $\boldsymbol{v_\theta}$ to represent the direction strategy. $t_{i-1} - t_i$ denotes the stepsize on

each optimization step. Many previous works have explored the selection of optimization direction $\boldsymbol{v_i}$. For example, in Heun [1, 10], $\boldsymbol{v_i} = \frac{1}{2}(\boldsymbol{v_{t_i}} + \boldsymbol{v_{t_{i+1}}})$, in DPM-Solver [18], $\boldsymbol{v_i} = \boldsymbol{v}_{t_\lambda(\frac{\lambda_{t_i}+\lambda_{t_{i+1}}}{2})}$, $t_\lambda$ is the inverse function mapping from $\lambda$ to $t$.

However, few works have discussed the selection of stepsize $t_i$. Most of the previous works can be divided into two categories. The first class selects a uniform $t_i$:

$$t_i = \frac{i}{N} \quad where \; i = 0 \dots N \tag{5}$$

such practices include LDM [23], DiT [21], FLUX [11], SDXL [22], etc. The second class includes EDM [10], CM [30], which adopts a handcraft stepsize schedule:

$$t_i = (t_1^{1/\rho} + \frac{i-1}{N-1}(t_N^{1/\rho} - t_1^{1/\rho}))^\rho \tag{6}$$

However, these ad-hoc schedules always lead to a suboptimal result. In this paper, we emphasize that the choice of stepsize is particularly crucial, especially in scenarios with limited steps. Besides, the stepsize operates orthogonally to the directional strategy.

### 3.2. Definition of optimal denoising stepsize

Although discretizing the ODE trajectory can significantly increase the speed, it will deviate from the original ODE and cause the results to deteriorate to a certain extent. Therefore, the stepsize discretization aims to generate results that are as consistent as possible with the non-discretization. This is similar to traditional distillation, so we formulate our task as stepsize distillation with the difference being that no parameter training is required:

Suppose that we have a teacher inference schedule, including both the direction schedule ($\boldsymbol{v_\theta}$ in Equation 4) and the stepsize schedule ($\{t^N\} = \{t_N, t_{N-1}, ..., t_1, t_0\}$), where the number of steps $N$ is large enough to approximate infinite steps sampling results. We abbreviate the result generated from the initial point $\boldsymbol{x_T}$ through the stepsize strategy $\{t^N\}$ and direction strategy $\boldsymbol{v_\theta}$ according to Equation 4 as $\mathcal{F}(\boldsymbol{x_T}, \boldsymbol{v_\theta}, \{t^N\})$. For one step denoising of state $\boldsymbol{x_i}(t_j < t_i)$ starts from $t_i$ and ends at $t_j$, the results can be format as follows:

$$\mathcal{F}(\boldsymbol{x_i}, \boldsymbol{v_\theta}, \{t_i, t_j\}) = \boldsymbol{x_i} + \boldsymbol{v_\theta}(t_j - t_i) \tag{7}$$

We aim to distill a student stepsize schedule, ($\{s^M\} = \{s_M, s_{M-1}, ..., s_1, s_0\}$), where $\{s^M\}$ is a subset of $\{t^N\}$ and $M$ is smaller than $N$. For a fixed initial noise $\boldsymbol{x_T}$, the optimal student stepsize schedule $\{s^M\}^o$ aims to minimize the global discretization error with the teacher generation result:

$$\{s^M\}^o = \arg\min_{\{s^M\}\subset\{t^N\}} \left\| \mathcal{F}(\boldsymbol{x_T}, \boldsymbol{v_\theta}, \{t^N\}) - \mathcal{F}(\boldsymbol{x_T}, \boldsymbol{v_\theta}, \{s^M\}) \right\|_2^2 \tag{8}$$
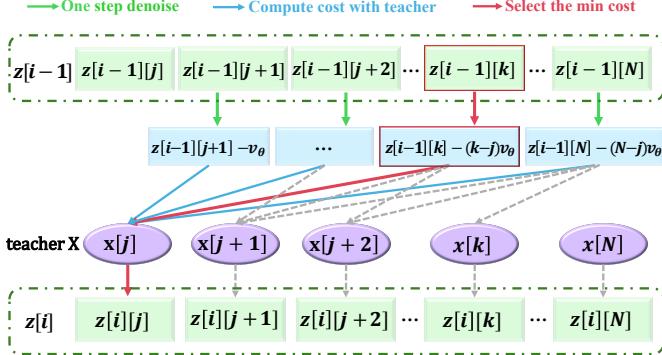
Figure 3. Subtask illustration of the recursive subtasks. The optimal results at timestep $j$ using $i$ step denosing ($z[i][j]$) derives from the $i - 1$ step optimal denosing results($z[i-1]$).

### 3.3. Dynamic programming based calibration

#### 3.3.1. Recursive subtasks

Our goal is to use $M$ student steps to approximate $N$ teacher steps sampling results as closely as possible. We separate it into subtasks: *using $i$ student steps to approximate $j$ teacher steps sampling results*.

Therefore, we make such notations: the sampling trajectory of the teacher model is $\{X^t\} = \{x[t_N], x[t_{N-1}], ..., x[t_1], x[t_0]\}$ based on the stepsize schedule $\{t^N\}$, where $x[t_j] = \mathcal{F}(x_T, v_\theta, \{t_N, ..., t_j\})$ denotes the sampling result at timestep $t_j$. For simplicity, we abbreviate $x[t_j]$ as $x[j]$. For the student model, we denote the optimal denoised result using $i$ step to approximate $x[j]$ as $z[i][j]$. Moreover, We additionally record that $z[i][j]$ comes from timestep $r[j](r[j] > j)$, which means:

$$r[j] = \underset{j+1,...,N}{\arg\min} \|x[j] - \mathcal{F}(z[i-1][r[j]], v_\theta, \{r[j], j\})\|_2^2 \tag{9}$$

We can find that these subtasks demonstrate a Markovian recursive property contingent upon the number of student steps utilized: The student model's $i$-th step approximation for the teacher's $j$-th step is determined by its prior result, where the approximation at timestep $r[j]$ utilized $i - 1$ denosing steps. It can be formally expressed as:

$$z[i][j] = \mathcal{F}(z[i-1][r[j]], v_\theta, \{r[j], j\}) \tag{10}$$

By combining Equation 10 and Equation 7, the recursion relation can be written as:

$$z[i][j] = z[i-1][r[j]] + v_\theta(j - r[j]) \tag{11}$$

The illustration of solving the recursive subtasks is shown in Figure 3.

#### 3.3.2. Dynamic programming

According to the recursion formula, we use dynamic programming to solve the stepsize selection problem for the

---

**Algorithm 1** Optimal Stepsize Distillation.

> $D_\theta$ : Denoising model which outputs the velocity.
> $x[N]$, $z$ samples from $\mathcal{N}(0, 1)$

1:   # get teacher trajectory.
2:   **for** $i = N$ to $1$ **do**
3:       $v = D_\theta(x[i], i)$
4:       $x[i - 1] = x[i]$ - $v$
5:   # use the student i steps result to approach the teacher's j steps result.
6:   **for** $i = 1$ to $M$ **do**
7:       $v = D_\theta(z[i-1, :], range(N, 1))$
8:       $z_{tmp} = \text{ones}(N + 1, N) \times \text{Inf}$
9:       **for** $j$ in $N$ to $1$ **do**
10:          $z_{tmp}[j, :] = z[i-1, j] - v[j] \times range(0, j - 1)$
11:      **for** $j$ in $N - 1$ to $0$ **do**
12:          $cost = MSE(z_{tmp}[:, j], x[j])$
13:          $r[i, j] = argmin(cost)$
14:          $z[i, j] = z_{tmp}[r[i, j], j]$
15:  # retrieve the optimal trajectory.
16:  $R = [0]$
17:  **for** $i$ in $M$ to $1$ **do**
18:      $j = R[-1]$
19:      $R.\text{append}(r[i][j])$
20:  **return** $R[1 :]$

---

student model. This process can be regarded as searching $z[i][j]$ and recording each $r[j]$ as $i$ increases. Then we can obtain the final result of using $M$ student steps approximating $N$ teacher steps $z[M][N]$. Finally, we backtrack the timestep $r[j]$ as the optimal trajectory. In this way, we distilled the optimal stepsize of the student model from the teacher steps. The pseudo code is provided in Algorithm 1.

#### 3.3.3. Theoretically optimal stepsize

To rigorously establish the theoretical optimality of the dynamic programming algorithm in attaining the optimal stepsize schedule, we demonstrate that this subtask decomposition simultaneously satisfies both *ovealapping subproblems* and *optimal substructure*.

**Overlapping Subproblems**. This can be easily found since each denoised step is only dependent on the previous one step result, which is independent of the denoised history before. Take the one order solver as an example. Suppose we want to denoise from timestep $s$ to timestep $t$, we have:

$$x_t = x_s + v_\theta(t - s) \tag{12}$$

One special case is the multistep high order solvers such as DPM-Solver [19] and UniPC [38]. Compared with one order solver, the only difference is that the denoising process needs to rely on multi-step results. In this case, the overlapping subproblem changes into denoising from multiple

4

states with different noise intensities. We leave the details about searching with high order solvers in the Appendix.

**Optimal Substructure**. The optimal substructure enables pruning suboptimal denoising trajectories in advance by dynamic programming. Here, we illustrate this substructure in Theorem 3.1 and leave the proof in the Appendix.

**Lemma 3.1.** *The optimal m step denosing results of student solver $z[m]$ always derives from the optimal m-1 step results $z[m-1]$ with additional one step denosing.*
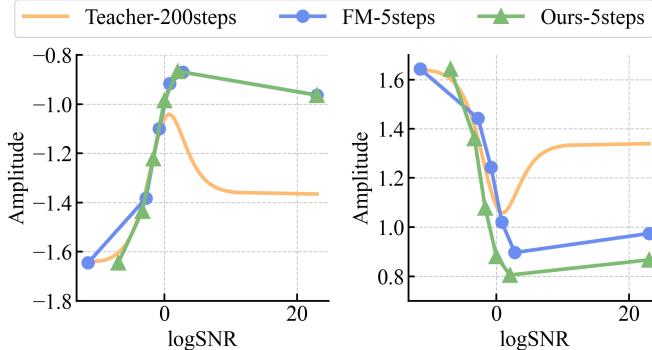


Figure 4. Amplitude of input tensor throughout denoising steps. The left and right plot the quantile of 5%, and 95% respectively.

### 3.4. Amplitude calibration

While our step optimization enables efficient trajectory approximation, we identify a critical challenge when applying few-step inference: systematic amplitude deviation becomes pronounced in low-step regimes. In Figure 4, we plot the the quantile of 5% and 95% percentiles of the input tensor to denoising models across varying noise intensities. We find that the distribution changes significantly between the few-step (5) and multiple-step (200) trajectories, especially in the later stages of the denoising process.

To address this, we propose a per-step affine transformation that aligns student outputs with teacher amplitude characteristics. For each noise level $\lambda_t$, we calibrate the transformation using a linear function:

$$\hat{\boldsymbol{x}}_t = \frac{S_{95} - S_5}{x_{t,95} - x_{t,5}} \boldsymbol{x}_t \tag{13}$$

where $S_{95}$ and $S_5$ denote the quantile ranges of 95%-5% of the teacher on 512 samples, which are calculated in advance. $x_{t,95}$ and $x_{t,5}$ are the quantile ranges of 95% - 5% of the student at step $t$ during denosing. In Section 4.2.1, we show that applying this calibration at the test time can significantly improve details such as color and texture.

## 4. Experiments

**Implementation details.** Our algorithm features plug-and-play compatibility, allowing for seamless integration into

various diffusion frameworks with minimal adaptation effort. To ensure consistent parameterization across different implementations, we adopt the velocity prediction paradigm from Flow Matching(FM)[15, 17] as our unified inference target. This requires wrapping the existing models with two standardization operations: (1) reparameterize the prediction target to velocity space, and (2) align the sampling trajectory based on the Signal-to-Noise ratio (SNR) matching. As theoretically guaranteed by [12], this standardization does not affect the sampling result while enabling unified stepsize optimization across architectures. In the Appendix, we will explain this unification in detail.

**Evaluation Protocol** We prioritize minimizing the divergence of the trajectory from the teacher models as our main optimization objective. For quantitative evaluation, we employ PSNR to measure pixel-level fidelity between student and teacher output. Unless otherwise specified, all teacher models utilize 200 sampling steps to empirically balance computational efficiency with output quality across most test domains.

### 4.1. Robustness analysis

We conduct comprehensive ablation studies to validate our method's robustness across four critical dimensions: noise schedules (Section 4.1.1), number of teacher steps (Section 4.1.2), ODE solver orders (Section 4.1.3), and diverse modeling frameworks (Section 4.1.4).

#### 4.1.1. Robustness to noise schedules

To evaluate schedule invariance, we performed class-conditional image generation experiments on ImageNet-64[24] using an EDM-based [10] diffusion model. We test three established noise schedules as teacher configurations:
- EDM's exponential noise schedule [10].
- DDIM's linear noise schedule [27].
- Flow Matching's schedule [15]

Each teacher generates reference samples using 200 steps. We compare our optimized steps against naive uniform step reduction, measuring similarity through PSNR between student and teacher outputs. As shown in Table 2, our method achieves 2.65 PSNR improvements on average over baseline step reduction across all schedules, demonstrating consistent performance regardless of the teacher's schedule.

#### 4.1.2. Robostness to number of steps

We analyze the impact of the number of teacher steps based on ImageNet experiments. The teacher solver adopts the DDIM schedule and generates images using various numbers of steps from 100 to 1000 (student fixed at 20 steps). We measure the PSNR against 1000-step teacher results. Table 3 reveals that the performance basically converged in 200 steps. This indicates that the 200-step teacher model already provides an adequate search space for the dynamic

| | steps | EDM | DDIM | FM |
|---|---|---|---|---|
| | 50 | 33.71 | 34.31 | 31.78 |
| Naive step | 20 | 25.76 | 26.12 | 23.97 |
| reduction | 10 | 20.83 | 20.89 | 19.00 |
| | 5 | 16.38 | 15.96 | 15.69 |
| | 20 | 28.21(**+2.45**) | 28.07(**+1.95**) | 28.51(**+4.54**) |
| OSS | 10 | 22.88(**+2.05**) | 22.79(**+1.9**) | 23.02(**+4.02**) |
| | 5 | 18.41(**+2.03**) | 18.05(**+2.09**) | 18.53(**+2.84**) |

Table 2. ImageNet-64 pixel space conditional image generation with different teacher schedules. All the student solvers are learned from the corresponding 200-step teacher schedule. Results are shown in PSNR calculated with 200-step corresponding teacher results.
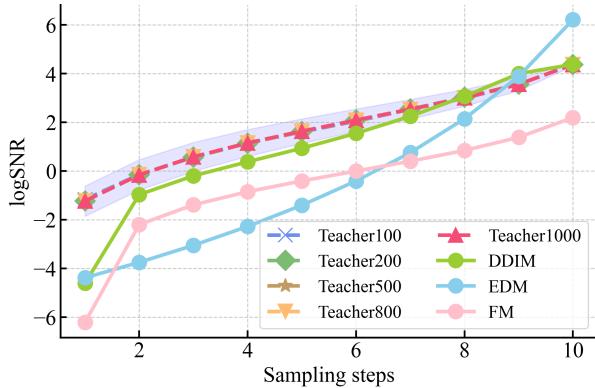


Figure 5. Step schedule for different solvers. Our method achieves nearly identical results from different teacher model steps.

programming algorithm to identify the optimal 20-step policy. Moreover, our method also maintains fidelity well, even when the teacher model has only 100 steps.

We further analyze the variance of step distribution across samples in Figure 5. For teacher solvers with different steps, the optimal step sequences almost overlap. Besides, different samples demonstrate minor differences (light purple area in the figure). This allows us to generalize to unsearched samples by averaging over a selected subset of instances, thereby reducing the inference overhead in a zero-shot manner. We named this setting as *OSS-ave*. Specifically, by employing 512 samples to compute the mean of the optimal stepsize schedule, we achieve results that closely approximate those obtained through instance-specific optimization, as shown in Table 3.

### 4.1.3. Robustness to higher order solver

Our optimal stepsize searching algorithm can be applied to higher-order solvers. We validate it using DiT-XL/2 model[21] on ImageNet 256×256 using classifier-free guidance as 1.5. Figure 6 demonstrates that our step sequences combined with third-order solvers achieves significant improvements over the baseline DDIM schedule.

This verifies that optimal step scheduling and the solver

| Teacher steps | 100 | 200 | 500 | 800 | 1000 |
|---|---|---|---|---|---|
| DDIM teacher | 37.97 | 44.24 | 52.93 | 57.30 | - |
| OSS | 27.04 | 28.07 | 27.11 | 27.12 | 27.13 |
| OSS-ave | 26.41 | 26.47 | 26.47 | 26.48 | 26.48 |

Table 3. ImageNet-64 pixel space conditional image generation results with different teacher steps. All the teacher solver leverages DDIM schedules. Students adopt 20-steps for inference. Results are shown in PSNR calculated with the 1000-step teacher.

order constitute complementary axes for few-step sampling improvement - their synergy enables practical acceleration without quality degradation.
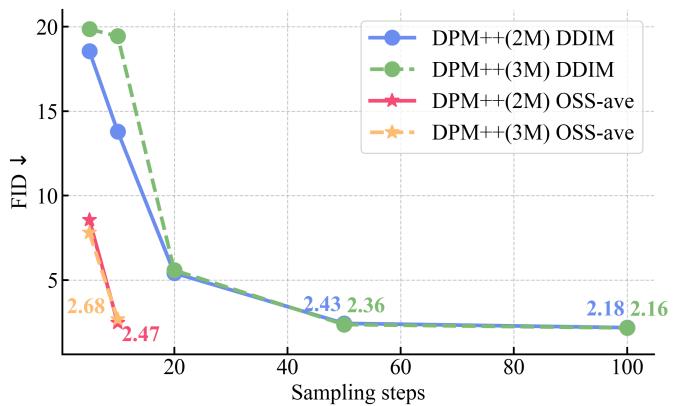


Figure 6. ImageNet-256 latent space conditional image generation results using different order solvers.

### 4.1.4. Generalization Across frameworks

We validate framework-agnostic applicability through two different extensions: **Masked Autoregressive (MAR) Generation:** By applying the optimal stepsize schedule to the MAR [13] ImageNet-256 model, we reduce the sampling steps from 500 to 5, achieving $100\times$ acceleration while maintaining competitive performance (FID=4.67). When decreasing steps to 50, the FID remained nearly unchanged ($2.66 \rightarrow 2.59$), demonstrating robust stability. Note that the model MAR-Base contains 32 tokens, but we only search for the first token and apply the result steps to the other 31 tokens, which further demonstrates the robustness.

**Video Diffusion:** In Open-Sora [39] frameworks, our optimized schedules enable $10\times$ speedups while preserving visual fidelity. Visualization results are shown in the Appendix. We find that simply reducing the sampling step to 20 steps completely changes the generated video. However, by adjusting the sampling steps, the generated result is almost the same as the teacher with $10\times$ speedup.

Figure 8. Through amplitude adjustment, the synthesized outputs (bottom row) exhibit enhanced details.
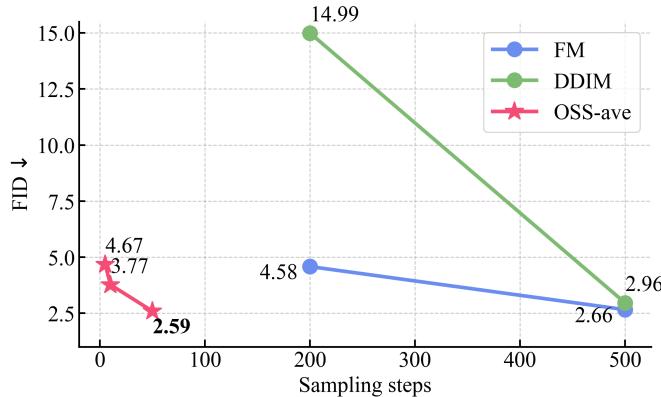


Figure 7. ImageNet-256 latent space conditional image generation results with MAR-Base.

## 4.2. Extensions of optimal stepsize

### 4.2.1. Amplitude calibration

As demonstrated in Section 3.4, insufficient sampling steps may induce significant amplitude discrepancies between generated intermediate states and the teacher model's trajectories. To investigate it, we conducted quantitative and qualitative analysis on the ImageNet dataset. Although amplitude adjustment may slightly decrease the PSNR compared to the teacher model's outputs ($18.53 \rightarrow 16.92$), the adjusted samples exhibit enhanced structural detail and demonstrate improved realism, as evidenced in Figure 8.

| | $x_t$-MSE | $x_0$-MSE | $x_0$-IncepV3 | $x_t$-percep |
|---|---|---|---|---|
| PSNR ↑ | 18.54 | 22.45 | 9.91 | 18.27 |

Table 4. Ablation studies on different functions.

### 4.2.2. Metric functions

Conventional continuous diffusion alignment relies on MSE between teacher and student intermediate states $x_t$. However, MSE may induce over-smoothing in pixel-level generation tasks. We thus explore three metrics: (1) MSE on reparameterized clean data $x_0$, this can be regarded as SNR-weighted MSE on $x_t$ [6]; (2) Inception-V3 [31] feature dis-

tance on reparameterized $x_0$; (3) perceptual distance utilizing the denoising network [21] on intermediate feature $x_t$.

We evaluate the three aforementioned metrics on the ImageNet-64 generation with the pretrained EDM model and summarized the results in Table 4. The $x_0$-MSE setting outperforms the default distance metric calculated on $x_t$. However, the Inception-V3 feature metric performs worse, potentially attributed to lacking generalization ability on noisy data. Meanwhile, using the denoising model to calculate the perceptual distance does not demonstrate significant improvement.

| | w/o training | w training |
|---|---|---|
| Uniform | 52.73 | 24.53 |
| OSS | 41.96 | **7.89** |

Table 5. Sampling results on ImageNet-256. Uniform is the schedule adopted in DDIM.

| Method | 5steps | 10steps | 5steps | 10steps |
|---|---|---|---|---|
| | ImageNet-64 [24] | | LSUN-Bed-256 [37] | |
| EDM | 17.72 | 22.33 | 15.60 | 19.61 |
| DDIM | 17.69 | 22.22 | 15.39 | 20.18 |
| FM | 17.80 | 22.47 | 14.43 | 18.23 |
| DM | 17.39 | 20.50 | **18.90** | 20.67 |
| GITS | 16.29 | 19.40 | 11.58 | 13.55 |
| OSS-ave | **17.94** | **22.47** | 16.51 | **21.25** |
| | FFHQ-64 [9] | | MSCOCO-512 [14] | |
| EDM | 18.68 | 23.24 | 13.43 | 15.45 |
| DDIM | 18.37 | 23.57 | 12.31 | 14.58 |
| FM | 17.72 | 21.31 | 11.90 | 13.51 |
| DM | 19.79 | 22.69 | 13.50 | 14.73 |
| GITS | 20.46 | 25.20 | 10.60 | 11.07 |
| OSS-ave | **20.56** | **25.21** | **14.10** | **16.45** |

Table 6. PSNR results of different stepsize schedules.

### 4.2.3. Optimal Stepsize as balanced task partition

Our optimal stepsize schedule not only enables directly accelerating diffusion models sampling without finetuning network, but also inherently reflects the varying difficulty levels of denoising tasks across different noise intensities. This observation suggests that the optimal stepsize schedule can be interpreted as a principled multitask partitioning strategy, which may better balance the subtask relationships compared to conventional uniform timestep division, thereby mitigating distortion on final outputs. To validate this, we conduct experiments on ImageNet generation with the DiT-XL/2 [21] model by splitting the entire denoising process into five consecutive stages, with each stage modeled by independent parameters. Inspired by [2, 26], all submodels are initialized with the pre-trained teacher model to accelerate convergence. The results shown in Table 5 demonstrate that our optimal stepsize can be regarded as

A photo of a yellow dining table and a pink dog



A photo of a cow left of a stop sign

A photo of a toaster below a traffic light

| Teacher | FM | DDIM | GITS | DM | Ours |

Figure 9. Visualization results on Geneval benchmark. Our optimal sampling schedule can produce results that are more similar to those of multi-step teachers, inherited strong instruction following ability.

| Method | steps | Overall | Single object | Two object | Counting | Colors | Position | Color attribution |
|--------|-------|---------|---------------|------------|----------|--------|----------|-------------------|
| Flow matching | 100 | 0.649 | 0.989 | 0.803 | 0.722 | 0.774 | 0.220 | 0.390 |
| Flow matching | 10 | 0.590 | 0.963 | 0.727 | 0.663 | 0.668 | 0.213 | 0.310 |
| DDIM | 10 | 0.623 | 0.971 | 0.742 | 0.694 | 0.758 | **0.22** | 0.353 |
| GITS | 10 | 0.604 | **0.981** | 0.732 | 0.678 | 0.697 | 0.210 | 0.325 |
| DM | 10 | 0.643 | 0.971 | **0.788** | 0.719 | **0.777** | 0.188 | **0.415** |
| OSS-ave | 10 | **0.645** | **0.981** | 0.775 | **0.728** | **0.777** | 0.195 | **0.415** |

Table 7. Results of Geneval based on the Flux.1-dev model with different sampling schedules.

a good task division, better allocating the model capacity. The result can be further improved when combined with other techniques, such as progressive distillation [26].

## 4.3. Compare with other methods

We compare our method with other stepsize strategies including ad-hoc schedules like DDIM, EDM, and flow matching, as well as dynamic adjustment methods such as DM [36] and GITS [3]. Table 6 demonstrates our experimental results across four distinct datasets. Our method achieves results closest to those of the teacher model under 5 and 10 step generation setting. Moreover, we also validate our results on the GenEval [5] benchmark using Flux.1-dev [11] as the foundation model. We list the quantitative results in Table 7, where the OSS-ave setting adopts 32 images for each category. We find that the optimal stepsize schedule achieves $10\times$ acceleration compared to the 100-step teacher model with almost no performance degrada-

tion on the GenEval benchmark. As shown in Figure 9, our method demonstrates the closest proximity to the teacher model comparing with other approaches.

## 5. Conclusion

This paper proposes a dynamic programming framework for stepsize optimization in diffusion sampling. By reformulating stepsize scheduling as recursive approximation, our method derives theoretically optimal step sequences with low computational overhead. The experiments demonstrate universal robustness across architectures and solvers while maintaining output fidelity. Its seamless integration with existing direction optimization strategies enables practical deployment advantages. This work establishes an alternative pathway for efficient diffusion sampling.

# References

[1] Uri M Ascher and Linda R Petzold. *Computer methods for ordinary differential equations and differential-algebraic equations*. SIAM, 1998. 2, 3

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 7

[3] Defang Chen, Zhenyu Zhou, Can Wang, Chunhua Shen, and Siwei Lyu. On the trajectory regularity of ode-based diffusion sampling. *arXiv preprint arXiv:2405.11326*, 2024. 2, 3, 8

[4] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis, 2024. *URL https://arxiv. org/abs/2403.03206*, 2. 12

[5] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023. 2, 8

[6] Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and Baining Guo. Efficient diffusion training via min-snr weighting strategy. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7441–7451, 2023. 7

[7] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[8] Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021. 2

[9] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 7

[10] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022. 1, 2, 3, 5

[11] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2023. 3, 8, 12

[12] Sangyun Lee, Zinan Lin, and Giulia Fanti. Improving the training of rectified flows. *arXiv preprint arXiv:2405.20320*, 2024. 5, 12

[13] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:56424–56445, 2025. 6

[14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014. 7

[15] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 2, 5

[16] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022. 2

[17] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 2, 3, 5

[18] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 1, 2, 3

[19] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 1, 2, 4

[20] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021. 3

[21] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2, 3, 6, 7, 12

[22] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 3

[23] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3

[24] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 5, 7

[25] Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your steps: Optimizing sampling schedules in diffusion models. *arXiv preprint arXiv:2404.14507*, 2024. 2, 3

[26] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 7, 8

[27] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1, 2, 3, 5, 12

[28] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020. 1

[29] 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*, 2020. 1

[30] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023. 2, 3

[31] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 7

[32] Vinh Tong, Trung-Dung Hoang, Anji Liu, Guy Van den Broeck, and Mathias Niepert. Learning to discretize denoising diffusion odes. *arXiv preprint arXiv:2405.15506*, 2024. 2, 3

[33] WanTeam, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang 1, Wei Wang 2, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. 2025. 12

[34] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers for diffusion models by differentiating through sample quality. In *International Conference on Learning Representations*, 2021. 2

[35] Daniel Watson, Jonathan Ho, Mohammad Norouzi, and William Chan. Learning to efficiently sample from diffusion probabilistic models. *arXiv preprint arXiv:2106.03802*, 2021. 2

[36] Shuchen Xue, Zhaoqiang Liu, Fei Chen, Shifeng Zhang, Tianyang Hu, Enze Xie, and Zhenguo Li. Accelerating diffusion sampling with optimized time steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8292–8301, 2024. 2, 3, 8

[37] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 7

[38] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2, 4

[39] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. *arXiv preprint arXiv:2412.20404*, 2024. 6, 12

[40] Zhenyu Zhou, Defang Chen, Can Wang, and Chun Chen. Fast ode-based sampling for diffusion models in around 5 steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7777–7786, 2024. 2

# Appendix

## A. Proof of the optimal substructure

Here, we aim to proof: *To solve the optimal approximation $z[m][n]$, which indicate the $m$-step student outcome to the $n$-step teacher outcome, it always derives from $z[m-1][k]$, where $k > n$.*

**Assumption 1.** *If $x$ is closer to the teacher denosing result at $x[i]$ than $x'$, as the form $\|x - x[i]\|^2 \leq \|x' - x[i]\|^2$, we have:*

$$\|v_{\theta,k} - v_{\theta^*}\|_2^2 \leq \|v'_{\theta,k} - v_{\theta^*}\|_2^2 \tag{14}$$

We prove our theorem by the method of contradiction. Suppose there exist $x'$ which is suboptimal than $z[m-1][k]$ and satisfies:

$$\|x' - x[k]\|_2^2 \geq \|z[m-1][k] - x[k]\|_2^2 \tag{15}$$

But the optimal result at step $m$ is generated from the suboptimal result at timestep n, which satisfies:

$$\|\mathcal{F}(x', v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{16}$$
$$\leq \|\mathcal{F}(z[m-1][k], v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{17}$$

We use $v'_\theta$ represents the velocity of $x'$ at timestep k, use $v_\theta$ represents the velocity of $z[m-1][k]$ at timestep k:

$$\mathcal{F}(x', v_\theta, \{k, n\}) = x' + v'_\theta(n-k) \tag{18}$$

Also, we assume the teacher solver results are an excellent ground truth approximation. We use $v_{\theta^*}$ represents the ground truth velocity, which follows:

$$x[n] = x[k] + v_{\theta^*}(n-k) \tag{19}$$

*Proof.* Here, we prove this in two cases. When the $x'$ and $x$ are closed to the teacher $x[k]$, we have:

$$\|\mathcal{F}(x', v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{20}$$
$$= \|x' - x[k] + (n-k)(v'_\theta - v_{\theta^*})\|_2^2 \tag{21}$$
$$\approx \|(n-k)(v'_\theta - v_{\theta^*})\|_2^2 \tag{22}$$
$$\geq \|(n-k)(v_\theta - v_{\theta^*})\|_2^2 \tag{23}$$
$$\approx \|\mathcal{F}(z[m-1][k], v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{24}$$
$$\tag{25}$$

This is contradictory to the hypothesis Eq. (16).

When the $x'$ and $x$ are far from the teacher $x[k]$, the input of the denosing network is far away from the training input. Therefore, we consider vectors $x' - x[k]$ and $v'_\theta - v_{\theta^*}$ to be statistically independent and treat them as random

---

**Algorithm 2** Optimal stepsize for high order solver

$D_\theta, x[N], z$ same as one order algorithm.
$od$: Inference Order.
$\mathcal{F}$ : Forward function, $\mathcal{F}(x_{cur}, t_{cur}, v, t_{next}, order)$
1:  # get teacher trajectory $x[i]$
2: **for** $i = N$ to 1 **do**
3:    $v = D_\theta(x[i], i)$
4:     $x[i-1] = \mathcal{F}(x[i], i, v, i-1, od)$
5:  # use the student i steps result to approach the teacher's j steps result.
6: **for** $i = 1$ to $M$ **do**
7:    $v = D_\theta(z[i-1, :], range(N, 1))$
8:    $z_{tmp} = ones(od+1, N+1, N) \times Inf$
9:    **for** $o$ in $od$ to 1 **do**
10:      **for** $j$ in $N$ to 1 **do**
11:        $t_{next} = range(0, j-1)$
12:        $z_{tmp}[o, j, :] = \mathcal{F}(z[i-1, j], j, v[j], t_{next}, o)$
13:    **for** $j$ in $N-1$ to 0 **do**
14:      $cost = MSE(z_{tmp}[:, :, j], x[j])$
15:      $r[i, j] = argmin(cost.flatten())$
16:      $z[k, i] = z_{tmp}[r[i.j]//(N+1), r[i,j]\%(N+1), j]$
17:      $r[i, j] = r[i, j]\%(N+1)$
   # retrieve the optimal trajectory
18: $R = [0]$
19: **for** $i$ in $M$ to 1 **do**
20:    $j = R[-1]$
21:    $R$.append($r[i][j]$)
22: **return** $R[1 :]$

---

vectors in a high-dimensional space, thereby establishing their orthogonality. Then:

$$\|\mathcal{F}(x', v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{26}$$
$$= \|x' - x[k] + (n-k)(v'_\theta - v_{\theta^*})\|_2^2 \tag{27}$$
$$= \|x' - x[k]\|_2^2 + \|(n-k)(v'_\theta - v_{\theta^*})\|_2^2 \tag{28}$$
$$+ 2 * \langle x' - x[k], (n-k)(v'_\theta - v_{\theta^*})\rangle \tag{29}$$
$$\approx \|x' - x[k]\|_2^2 + \|(n-k)(v'_\theta - v_{\theta^*})\|_2^2 \tag{30}$$
$$\geq \|z[m-1][k] - x[k]\|_2^2 + \|(n-k)(v_\theta - v_{\theta^*})\|_2^2 \tag{31}$$
$$\approx \|\mathcal{F}(z[m-1][k], v_\theta, \{k, n\}) - x[n]\|_2^2 \tag{32}$$
$$\tag{33}$$

This is contradictory to the hypothesis Eq. (16). Therefore, the original hypothesis is not valid, the optimal results always come from the previous optimal results, which is the optimal substructure of our searching problems. $\square$

## B. Searching steps with high order solvers

Here, we introduce the algorithm of searching optimal steps for high order solvers. A key distinction lies in the solver's

dynamic order selection capability during the search phase: at each optimization step, the student solver may execute operations corresponding to the specified denoising order or any lower-order approximations. The optimal configuration is selected through minimum-distance alignment with the teacher trajectory across all denoising orders. After searching, the student solver utilize the optimal steps with the order given before. The complete optimization procedure is formalized in Algorithm 2.

## C. Unification of different sampling trajectories

This section elaborates on adapting various pre-trained models for sampling via flow matching. This process comprises three steps: transformation of network prediction targets, sampling trajectory alignment, and model input alignment. The network prediction targets, comprising $\epsilon$, $x_0$, $v$, exhibit mutual convertibility through reparameterization. Our framework unifies their conversion to velocity $v$ through systematic transformation.

Following the methodology in [12], temporal alignment of sampling trajectories requires matching the signal-to-noise ratio (SNR) at each step. The flow matching diffusion process follows:

$$x_t = t\epsilon + (1 - t)x_0 \qquad (34)$$

while we assume the conventional forward trajectory of the pre-trained model adheres to:

$$x_t = \alpha_t x_0 + \beta_t \epsilon \qquad (35)$$

Temporal alignment between flow matching time $t$ and reference trajectory time $t'$ is established through the ratio constraint:

$$\frac{1 - t}{t} = \frac{\alpha_{t'}}{\beta_{t'}} \qquad (36)$$

Furthermore, the network input alignment between the current trajectory position $x_t$ and the pre-trained model's expected input $x_{t'}$ is achieved via linear transformation:

$$x_{t'} = \frac{\alpha_{t'}}{1 - t}x_t \qquad (37)$$

These aligned coordinates ($t'$ and $x_{t'}$) enable direct utilization of pre-trained diffusion models to estimate the current trajectory velocity at timestep $t$, thereby enables sampling through flow matching.

## D. Visualization of optimal sampling schedule

In this section, we present additional visualizations of sampling schedules. We maintain the original teacher schedule used in each open-source implementation: DiT [21]
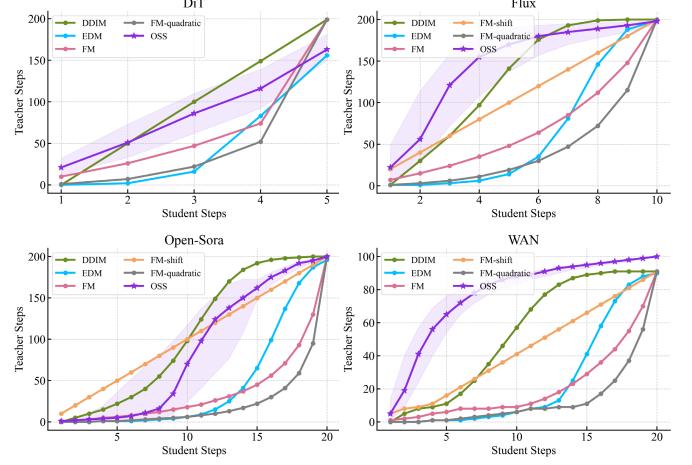


Figure 10. Different sampling schedules of DiT, Flux, Open-Sora, and WAN.

employs DDIM [27] schedule, while Open-Sora [39], Flux [11], and Wan [33] adopt flow matching with timeshift [4] as their default schedules. We compare different few-step sampling strategies and utilize the unified framework in Section C to align timesteps with the teacher model. We present the teacher model's sampling timesteps corresponding to each few-step sampling strategy in Figure 10.

## E. More results

In this section, we provide more Flux.1-dev [11] generation results using 10 steps. Moreover, we also provide the video generation results in the github page.

A photo of a bear above a spoon

A photo of three laptops

A photo of a bench and a sports ball

A photo of a stop sign above a fork

A photo of a laptop below a sports ball

A photo of a computer mouse and a zebra

A photo of three kites

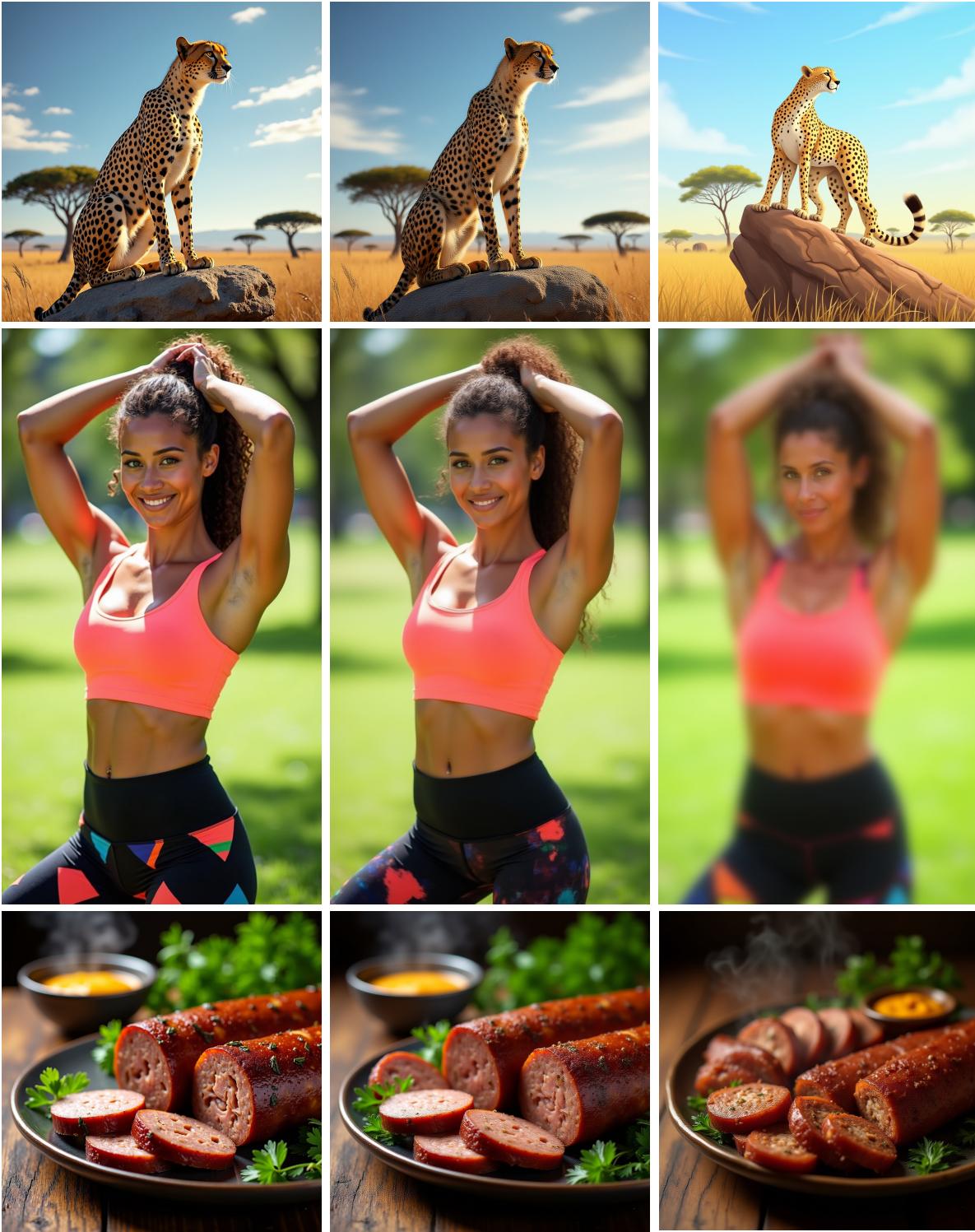| Teacher | FM | DDIM | GITS | DM | Ours |

Figure 11. Flux generation results on Geneval.

Figure 12. Flux generation results. Left: Original sampling result using 200 steps. Middle: Optimal stepsize sampling result within 10 steps. Right: Naively reducing sampling steps to 10.
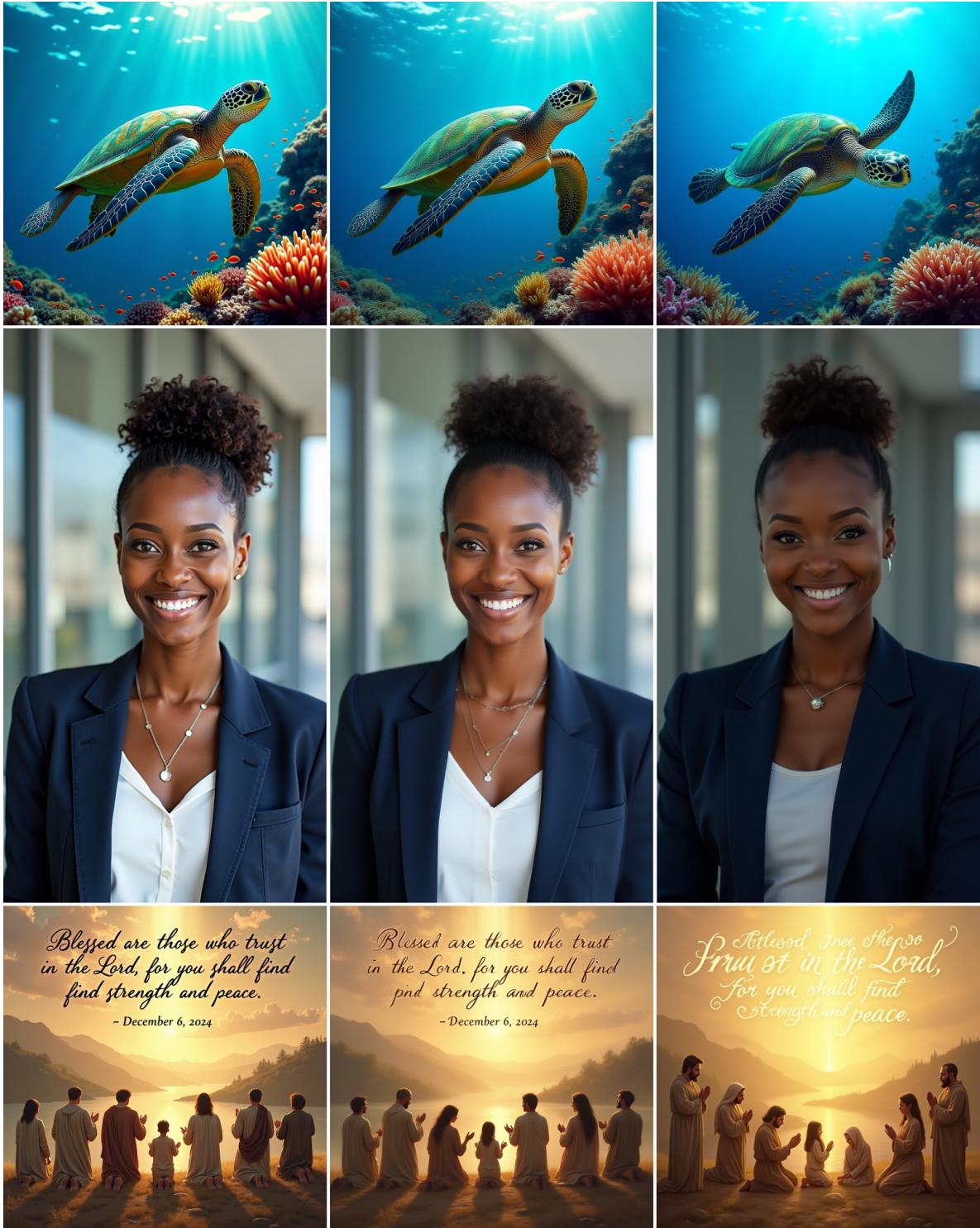
Figure 13. Flux generation results. Left: Original sampling result using 200 steps. Middle: Optimal stepsize sampling result within 10 steps. Right: Naively reducing sampling steps to 10.
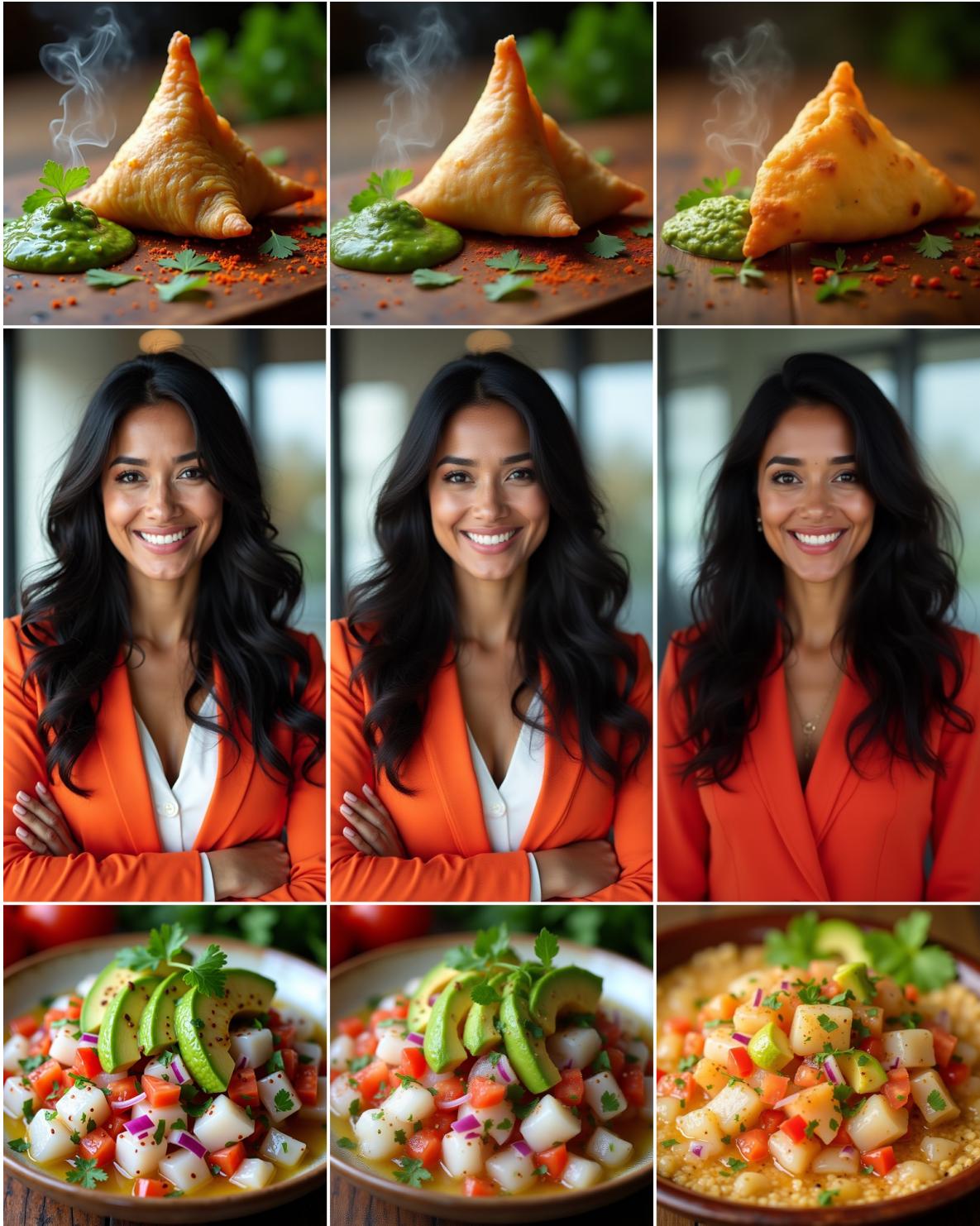
Figure 14. Flux generation results. Left: Original sampling result using 200 steps. Middle: Optimal stepsize sampling result within 10 steps. Right: Naively reducing sampling steps to 10.