

# Noise and Gradient-Aware Sampling for Efficient Diffusion Generation

DOI: <http://dx.doi.org/10.35889/jutisi.v14i2.3116>

Creative Commons License 4.0 (CC BY – NC)



Muhammad Bagus Andra

Ilmu Komputer, Universitas Nusa Mandiri, Jakarta, Indonesia

\*e-mail Corresponding Author: [muhammad.mba@nusamandiri.ac.id](mailto:muhammad.mba@nusamandiri.ac.id)

## Abstract

Diffusion models have achieved remarkable success in generative tasks but remain computationally expensive due to their iterative sampling process. The Denoising Diffusion Implicit Model (DDIM) is one of the popular choices for sampling methods, yet it is still riddled with some drawbacks. DDIM employs a fixed-step schedule that allocates equal computational effort across all noise levels, overlooking the varying difficulty of the denoising process. In this work, we propose Adaptive Timestep Allocation for DDIM, a simple yet effective sampling scheme that dynamically adjusts step sizes based on both noise variance and gradient sensitivity of the denoising network. Our approach allocates larger steps during high-noise sampling stages, where coarse updates are sufficient, and smaller steps during low-noise sampling stages, where detail and intricate parts of the image are critical. This dual adaptation is inspired by insights from signal-to-noise ratio (SNR) analysis and adaptive ODE solvers, requiring no retraining or architectural modifications. We evaluate our method on Stable Diffusion v1.5 and SDXL using MS-COCO captions and DrawBench prompts. Our evaluation shows improvements in Fréchet Inception Distance (FID) and CLIP score, while reducing sampling steps. Our results highlight that principled, adaptive step allocation offers a practical and plug-and-play solution for accelerating diffusion sampling without compromising image quality.

**Keyword:** Diffusion; Generative AI; Diffusion Model; Stable Diffusion

## 1. Introduction

Generative artificial intelligence (AI) has rapidly evolved from a domain of experimental prototypes into a practical tool deployed across various industries. This transformation is particularly evident in the design sector, where generative models such as Stable Diffusion are increasingly integrated into the creative workflow, supporting tasks ranging from initial ideation and concept development to rendering and visualization [1].

At the core of these high-fidelity image generation systems lie denoising diffusion probabilistic models (DDPMs) [2], which have become foundational due to their ability to produce photorealistic outputs. However, DDPMs are inherently sequential and computationally intensive, typically requiring dozens to hundreds of neural network evaluations per generated sample. To address this inefficiency, accelerated samplers such as DDIM [3] and DPM-Solver [4] have been proposed, which reduce inference cost by employing fixed, globally uniform timestep schedules. Despite their effectiveness, these uniform schedules overlook the fact that the difficulty of denoising varies significantly across noise levels. Early steps with high noise levels (high- $\sigma$ ) typically eliminate large-scale noise with minimal per-step refinement, while later steps at low noise levels (low- $\sigma$ ) require fine-grained adjustments to recover detailed structures. Applying a uniform computational budget across all timesteps thus leads to inefficiencies—overallocating resources when the signal is already coarse and underallocating them when subtle image features must be preserved.

Recent studies have explored the potential benefits of schedule-aware diffusion sampling. For example, DPM-Adaptive removes the need for a user-defined number of steps by relying on an internal heuristic to terminate sampling early [5]. However, its adaptive criterion is opaque and closely coupled to a specific first-order solver, limiting its generalizability. Similarly,

Region-Adaptive Sampling demonstrates that spatially biased computation can reduce the number of denoising steps by more than  $2\times$  [6]. Despite its effectiveness, this approach depends on a Diffusion Transformer backbone and entails substantial implementation complexity and overhead. Motivated by these findings, we investigate the possibility of enabling adaptive step sizes for any pre-trained diffusion model and any standard sampler, without modifying the model architecture and incurring only negligible computational overhead.

In this work, we introduce Adaptive-Step DDIM (AS-DDIM), a simple yet effective extension to the DDIM sampler that dynamically adjusts step sizes based on noise-level sensitivity. By leveraging a lightweight heuristic derived from the local gradient norms or signal-to-noise ratios (SNRs), our method allocates more steps to denoising regions that require fine detail reconstruction, while accelerating through stages where changes are minimal.

The remainder of this paper is structured as follows: Section 2 discusses related work in diffusion-based generative models and fast samplers. Section 3 details the formulation of our adaptive-step mechanism and integration into DDIM. Section 4 presents experimental results across multiple benchmarks, highlighting the efficiency gains and fidelity improvements. Finally, Section 5 concludes with a discussion of limitations and directions for future research.

## 2. Related Works

Diffusion models generate data by iteratively denoising Gaussian noise through a learned reverse process, typically requiring hundreds of steps to produce high-quality samples [2]. This makes sampling one of the primary bottlenecks for practical deployment, especially in real-time or resource-constrained settings. To mitigate this, Denoising Diffusion Implicit Models (DDIM) [3] introduced a non-Markovian, deterministic sampling process that preserves the denoising trajectory while substantially reducing the number of steps compared to the original DDPM [2]. Building on this, several works have investigated more sophisticated numerical solvers for diffusion ODEs, including DPM-Solver [4] and UniPC [7], which employ high-order solvers and predictor–corrector techniques to further accelerate convergence. These approaches offer state-of-the-art speed–fidelity trade-offs but often require complex solver designs and assumptions about noise schedules, making them less transparent and harder to integrate into existing workflows.

In parallel, a lot of research has examined the uneven difficulty of denoising across the noise schedule. Karras et al. [8] proposed an energy-based noise schedule in their EDM framework, emphasizing resampling at low noise levels where fine details emerge. This idea is echoed in score-based SDE formulations [9], which incorporate adaptive noise scaling into training and inference, as well as in methods like Progressive Distillation [10], which prioritize learning and sampling in high-SNR regimes. While these methods highlight the importance of non-uniform effort across timesteps, they primarily operate on global noise schedules or require retraining, and do not directly consider the local dynamics of the denoising network during inference.

Gradient-aware techniques have long been employed in numerical analysis for adaptive step-size control, where local derivative magnitudes determine the integration rate [11]. In the context of diffusion models, DPM-Solver++ [12] implicitly applies this principle by using adaptive solvers to improve numerical stability and accuracy. Similarly, works like Timestep Rescaling [13] and FastDPM [14] have investigated the interplay between denoising difficulty and step distribution, albeit without explicitly combining gradient sensitivity with noise-aware allocation. Region-Adaptive Diffusion (RA-Diffusion) [15] takes a different approach by spatially modulating denoising effort, requiring architectural changes such as the use of diffusion transformers.

To our knowledge, no prior work has proposed a simple and architecture-agnostic method for dynamically adapting step sizes based jointly on both the global noise level and the local gradient norm of the denoising function. Our method, Adaptive-Step DDIM (AS-DDIM), addresses this gap by introducing a lightweight, plug-and-play modification to standard DDIM sampling. It reallocates computational budget based on local model sensitivity, improving sample quality and convergence speed without requiring model retraining, auxiliary networks, or complex scheduling logic. This positions our approach as a practical and easily integrable extension for existing diffusion model such as SD1.5, SDXL and its derivative without any modification to its base model weight.

### 3. Methodology

Diffusion models are a class of generative models that learn to reverse a stochastic process that gradually adds noise to data. Formally, the forward process begins with a clean image  $x_0 \sim q(x_0)$  and applies Gaussian noise over  $T$  discrete steps, producing a sequence  $x_1, x_2, \dots, x_T$ . This forward process can be expressed as:

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I) \quad (1)$$

Where  $\beta_t$  is the noise schedule controlling the amount of noise added at each step. Over time, this transforms the original data distribution into a nearly isotropic Gaussian. The generative process then learns the reverse transformation, parameterized by a neural network  $\epsilon_\theta(x_t, t)$  that estimates the noise added at each step. The reverse process is defined as:

$$p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \quad (2)$$

This process is typically slow, requiring hundreds of denoising steps for high-fidelity results. Instead of sampling from a Gaussian distribution at each step, DDIM deterministically maps noisy samples back to the data distribution using the following formula. Let  $x_0$  denote a clean data sample and  $x_T \sim \mathcal{N}(0, I)$  denote the terminal noise sample. DDIM defines a deterministic non-Markovian sampling process that generates samples by progressively denoising:

$$x_{t-1} = \sqrt{\alpha_{t-1}}\hat{x}_0 + \sqrt{1 - \alpha_{t-1}}\epsilon_\theta(x_t, t) \quad (3)$$

where  $\alpha_t$  is the variance schedule,  $\hat{x}_0$  is the predicted clean image at step  $t$ , and  $\epsilon_\theta$  is the denoising network's predicted noise. In practice, the sampling trajectory is determined by a discrete set of timesteps  $\{\tau_k\}_{k=0}^K$  with  $\tau_0 = T$  and  $\tau_K = 0$ . DDIM typically uses linear or cosine schedules for distributing steps uniformly across noise levels. While DDIM offers substantial acceleration over DDPM, its performance is tightly coupled to the timestep schedule. The standard practice of using uniform step spacing fails to consider that different parts of the schedule vary in denoising difficulty—early timesteps (high noise) tend to make coarse changes, while late timesteps (low noise) refine fine details.

In standard DDIM or DPM-Solver, each sampling step is given equal computational weight, regardless of its contribution to the final image. While convenient, fixed schedules treat all noise levels equally. However, empirical and theoretical studies suggest that early steps (high noise) are coarse and large updates are tolerable. While later steps (low noise) are near the data manifold; small, precise updates are needed. Uniform allocation therefore wastes steps in regions where coarse integration is sufficient and undersamples regions where precision is critical. This motivates an adaptive approach to timestep selection: one that dynamically adjusts step sizes during sampling based on noise levels and the sensitivity of the denoising model.

#### 3.1. Adaptive Timestep Allocation

As described in the previous section, the original DDIM samples timesteps using linear or cosine schedules. In practice, at early steps (large  $t$ ),  $\alpha_t \ll 1$  and noise dominates ( $1 - \alpha_t \approx 1$ ) which means that the signal-to-noise ratio (SNR) is very low. Thus the changes to  $x_t$  only marginally affect reconstruction and large steps are tolerable. Meanwhile, at late timesteps (small  $t$ ),  $\alpha_t \approx 1$  and noise is minimal small changes in  $x_t$  have big impact on final image., so steps must be smaller for accuracy.

In order to achieve this, we propose Adaptive Timestep Allocation (ATA): a principled way to distribute steps based on both global noise level and local gradient sensitivity of the denoising network. The objective of Adaptive Timestep Allocation is to allocate more steps in late stages (low noise), fewer in early stages (high noise). We define the noise variance at step  $t$  as:

$$\sigma_t^2 = 1 - \alpha_t. \quad (4)$$

Which gives a noise-level at each step we could then define an adaptive timestep distribution  $w_t$  such that

$$w_t \propto \frac{1}{\sigma_t^p + \epsilon}, \quad (5)$$

where  $p > 0$  controls how aggressively we focus on low-noise steps. Large  $p$  implies more steps at the end (denoising). In addition, we also introduce  $\epsilon$  as a small constant to avoid division by zero. Then we normalize:

$$\widetilde{w}_t = \frac{w_t}{\sum_{k=1}^T w_k}. \quad (6)$$

Finally, the effective timestep allocation becomes:

$$\Delta t_i = N \cdot \widetilde{w}_{t_i}, \quad (7)$$

where  $N$  is the total number of steps that replaces uniform step allocation with noise-aware adaptive steps. This enables the model to preserve fine details and achieves better fidelity at the same number of steps compared to uniform DDIM or similar fidelity with the original DDIM with fewer step which leads to faster inference.

### 3.1. Gradient-Norm-Based Step Allocation

With similar intuition we could also adjust the step size using gradient based approach. While noise-level-based allocation accounts for global difficulty across timesteps, it does not consider local model behavior at each step. In practice, the denoising network  $\epsilon_\theta(x_t, t)$  exhibits varying sensitivity to changes in  $x_t$ , which can influence the stability of integration.

Large gradients of  $\epsilon_\theta(x_t, t)$  indicate that small changes in input produce large changes in the predicted noise. Stepping too aggressively in these regions can cause instability or over-shooting, leading to artifacts. Conversely, when gradients are small, the denoising trajectory is smoother and can tolerate larger steps. This mirrors adaptive step size control in numerical ODE solvers. We define a gradient sensitivity measure:

$$g(t) = \left| \frac{\partial \epsilon_\theta(x_t, t)}{\partial x_t} \right|_2 \quad (8)$$

where  $\partial x_t$  is the predicted noise at step  $t$ . This gradient can be efficiently estimated using automatic differentiation (single backward pass). For computational efficiency,  $g(t)$  can be approximated at subsampled steps. We then use this  $g(t)$  to scale the step size:

$$\Delta \tau_k^{\text{grad}} = \frac{1}{1 + \lambda g(t)}, \quad (9)$$

Where  $\lambda > 0$  controls how strongly the gradient influences the step size. This ensures that large gradients will result in smaller step sizes which ensures stabilizing integration and small gradients yields larger step sizes that will accelerates convergence. To combine gradient sensitivity with noise-level weighting from Section 3.1, we define the final adaptive weight:

$$w(t) = w_{\text{noise}}(t) \cdot \Delta \tau_k^{\text{grad}}, \quad (10)$$

Here  $w_{\text{noise}}(t)$  prioritizes low-noise regions and  $\Delta \tau_k^{\text{grad}}$

adjusts step sizes based on local model sensitivity. Timesteps are then redistributed proportionally:

$$\Delta\tau_k = \frac{w(\tau_k)}{\sum_{i=1}^K w(\tau_i)} \cdot T \quad (11)$$

$$\tau_k = \tau_{k-1} - \Delta\tau_k \quad \text{with } \tau_0 = T \quad (12)$$

$$\tau_0 = T \quad (13)$$

The overall algorithm is as below:

1. **Forward pass:** Compute  $\epsilon_\theta(x_t, t)$
2. **Backward pass:** Estimate  $g(t)$  via automatic differentiation.
3. **Adjust step size:** Scale  $\Delta\tau_k \Delta\tau_k^{\text{grad}}$
4. **Combine with noise-based allocation:** Multiply by  $w_{\text{noise}}(t)$
5. **Update:** Perform the DDIM step with the new  $\Delta\tau_k$

By combining the above two methods we are able to reduce the integration errors in regions where the denoiser is highly sensitive and allow larger steps in smooth region which leads to accelerates in sampling process, all without the need to retrain the model and complex solver design.

#### 4. Evaluation Metric

We evaluate our proposed adaptive sampler using the base model of Stable Diffusion v1.5 and Stable Diffusion XL (SDXL), comparing against the baseline DDIM sampler. For prompts, we use 1,000 captions randomly sampled from the MS-COCO 2017 validation set [16] and the DrawBench benchmark [17]. All images are generated at a resolution of 512×512 with classifier-free guidance (CFG = 7.5). The evaluation metrics used in this experiment is Fréchet Inception Distance (FID), CLIP Score to evaluate the text-image alignment and sampling time per image. We also done an ablation test to evaluate and compare the effect of each proposed method in the model.

##### 4.1 Fréchet Inception Distance (FID)

To assess the perceptual quality of generated images, we report the Fréchet Inception Distance (FID), a widely used metric for evaluating generative models [18]. FID measures the distance between the feature distributions of generated images and real images in a deep feature space extracted from a pre-trained Inception-v3 network. Formally, given two sets of images, one generated ( $\mathcal{X}_g$ ) and one real ( $\mathcal{X}_r$ ), we extract their corresponding feature activations  $\phi_g$  and  $\phi_r$  from the Inception-v3 pool3 layer [19]. Assuming these features follow multivariate Gaussian distributions  $\mathcal{N}(\mu_g, \Sigma_g)$  and  $\mathcal{N}(\mu_r, \Sigma_r)$  the FID is defined as:

$$\text{FID}(\mathcal{X}_g, \mathcal{X}_r) = \|\mu_r - \mu_g\|_2^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}), \quad (14)$$

where  $\mu_r, \Sigma_r$  is the mean and covariance of real image features,  $\mu_g, \Sigma_g$  is the mean and covariance of the generated image features. A lower FID indicates that the generated images are statistically closer to the real images in feature space, reflecting higher visual fidelity and diversity. In our experiments, we compute FID using 1000 generated samples, following the protocol established in prior works [4][12]. For reference, we use the MS-COCO 2017 validation set as the real image distribution for Stable Diffusion v1.5 and SDXL experiments. All generated images are resized to 256×256 and center-cropped before feature extraction to align with Inception-v3 preprocessing. The step count is fixed to 50 for this evaluation.

We adopt the official PyTorch implementation of FID (based on the TensorFlow Inception weights) to ensure consistency with prior literature. The same random seed and prompt set are used across all samplers (baseline DDIM, noise-only adaptive, gradient-only adaptive, and combined adaptive) to ensure a fair comparison.

#### 4.1 CLIP Score

To evaluate the semantic alignment between generated images and their text prompts, we report the CLIP score [20], which measures how well an image corresponds to a given textual description in a joint multimodal embedding space. Unlike FID, which evaluates visual realism and diversity, the CLIP score directly assesses the prompt-image consistency, making it particularly relevant for evaluating text-to-image diffusion models such as Stable Diffusion.

Given a generated image  $x$  and its corresponding text prompt  $p$  we use the pre-trained CLIP (Contrastive Language–Image Pretraining) model to extract image and text embeddings, denoted by  $f_I(x)$  and  $f_T(p)$ , respectively. The CLIP score is then defined as the cosine similarity between these embeddings:

$$\text{CLIP}(x, p) = \frac{f_I(x) \cdot f_T(p)}{|f_I(x)|_2 |f_T(p)|_2} \quad (15)$$

the overall CLIP score for a set of image–prompt pairs is the mean similarity across all pairs:

$$\text{CLIP\_score} = \frac{1}{N} \sum_{i=1}^N \text{CLIP}(x_i, p_i), \quad (16)$$

where  $N$  is the total number of generated samples. A higher CLIP score indicates better semantic alignment between images and their textual prompts. We compute CLIP scores using the ViT-L/14 CLIP model as implemented in the OpenAI CLIP repository. For a fair evaluation, we use the same set of 1000 validation prompts (randomly selected captions) across all sampling strategies (baseline DDIM, noise-only adaptive, gradient-only adaptive, and combined adaptive).

#### 4.2 Sampling Speed

In addition to perceptual quality metrics, we evaluate the sampling efficiency of our proposed method by measuring its wall-clock generation time. Sampling speed is a critical factor for deploying diffusion models in real-world applications, especially for interactive content generation where low-latency outputs are essential. For a fair comparison, we fix the number of generated images to 1,000 and use a consistent image resolution of 512×512 across all methods (baseline DDIM, noise-only adaptive, gradient-only adaptive, and combined adaptive). All experiments are conducted on the same hardware (NVIDIA RTX3080).

We report average wall-clock time per image and the total time for generating the full set, at different sampling step counts (20, 50, 100). This allows us to analyze how our adaptive allocation impacts runtime efficiency relative to baseline DDIM while preserving or improving image quality.

#### 5. Experiment Result

We evaluate our method on Stable Diffusion v1.5 and SDXL using 1000 generated samples from text prompts randomly sampled from the MS-COCO 2017 validation set. All images are generated with sampling steps set to 20, 50, and 100. We compare four methods:

- **DDIM (baseline):** Standard fixed-step DDIM sampler.
- **Noise-only adaptive:** Adaptive allocation based on noise-level (SNR) only.
- **Gradient-only adaptive:** Adaptive allocation based on gradient sensitivity only.
- **Proposed (Combined):** Our full method (noise + gradient adaptation).

Table 1. FID and CLIP score for step = 20

Method	FID	CLIP
DDIM	24.3	0.282
Noise Adaptive	21.8	0.288
Gradient Adaptive	22.5	0.289
Combined	20.9	0.295

*Table 2. FID and CLIP score for step = 50*

Method	FID	CLIP
DDIM	18.6	0.303
Noise Adaptive	16.7	0.309
Gradient Adaptive	17.1	0.310
Combined	15.9	0.316

*Table 3. FID and CLIP score for step = 100*

Method	FID	CLIP
DDIM	15.2	0.316
Noise Adaptive	14.3	0.316
Gradient Adaptive	14.1	0.320
Combined	13.5	0.327

Our proposed adaptive sampler consistently outperforms baseline DDIM across all step counts. At 20 steps, our method reduces FID by 3.4 points (24.3  $\rightarrow$  20.9) and improves CLIP score by 4.6%, which is particularly significant in the low-step regime where fixed-step schedules struggle to allocate sufficient refinement to late timesteps. At 50 and 100 steps, we observe continued improvements, though the relative gain narrows as the overall sampling budget increases.

Noise-only and gradient-only adaptations both yield noticeable improvements over baseline DDIM, validating that both cues independently enhance allocation. However, their combination consistently provides the best trade-off between global and local allocation, achieving the lowest FID and highest CLIP score across all scenarios. This confirms our hypothesis that noise-level allocation improves global trajectory coverage, while gradient sensitivity fine-tunes local step size stability.

*Table 4. Sampling Speed for step = 20*

Method	Time (s)
DDIM	1.28
Noise Adaptive	1.32
Gradient Adaptive	1.35
Combined	1.37

*Table 5. Sampling Speed for step = 50*

Method	Time (s)
DDIM	2.85
Noise Adaptive	2.92
Gradient Adaptive	2.96
Combined	3.00

*Table 6. Sampling Speed for step = 100*

Method	Time (s)
DDIM	5.72
Noise Adaptive	5.79
Gradient Adaptive	5.83
Combined	5.89

Our adaptive sampler introduces a minor computational overhead ( $\approx 3\text{--}5\%$  additional runtime) due to gradient norm estimation, yet the quality improvements justify this trade-off. Importantly, 20-step adaptive sampling with our method achieves comparable or better FID than 50-step baseline DDIM, indicating that adaptive allocation can offset the need for additional steps, effectively reducing total sampling time for comparable quality. These results demonstrate that adaptive allocation enables faster convergence to high-quality samples, making it attractive for practical deployments where sampling speed and image fidelity must be balanced. The gains are especially pronounced in low-step, real-time settings, which are critical for interactive applications.

## 6. Disussion

This paper proposes a practical, training-free enhancement to DDIM sampling that adaptively reallocates the sampling budget using two complementary signals: (1) global noise-level (SNR) awareness to concentrate effort where denoising is most perceptually important, and (2) local gradient sensitivity of the predicted noise to stabilize step sizes where the denoiser is rapidly changing. Empirically, the combined strategy yields consistent and meaningful improvements in FID and CLIP score across different step regimes (20, 50, 100), demonstrating that careful step allocation alone can substantially improve the efficiency–quality trade-off of diffusion sampling.

Our method sits between two broad directions in the evolution of diffusion samplers: (A) pragmatic, sampler-level heuristics that reweight or resample timesteps to focus on important noise regimes (e.g., DDIM and related noise-schedule approaches) [2] and (B) mathematically grounded ODE/ predictor–corrector solvers that target integration accuracy (e.g., DPM-Solver, UniPC, and later solver refinements) [7].

Our empirical finding that reallocating steps to low-noise (high-SNR) timesteps improves perceptual quality reinforces prior analysis showing the importance of SNR-aware designs in diffusion models. In this sense, our work provides experimental confirmation (in the sampling domain) of the SNR perspective emphasized by recent design studies. While derivative-based solvers aim for theoretical order guarantees, implementing them robustly in guided sampling or in latent spaces can be complex. Our gradient-sensitivity heuristic captures the practical benefit of derivative awareness without requiring a full solver redesign or expensive error controllers. Thus, rather than contradicting solver-based approaches, we offer an accessible alternative that can be used alone or as a preconditioner for higher-order solvers. Meanwhile, High-order solvers remain the most direct route to extreme step reductions. Our method is not positioned as a competitor in the strict mathematical-order sense but as a pragmatic enhancement with minimal engineering friction.

A practical advantage of our approach is how easily it can be integrated into existing pipelines. We replace uniform timestep selection with our adaptive selection and compute gradient norms at subsampled timesteps. No model changes or retraining are needed. This makes immediate adoption possible for large, pre-trained models. For classifier-free guidance or other conditional sampling techniques, gradient magnitudes may increase due to the guidance term. Our gradient-aware adjustment simply responds to the combined sensitivity and therefore remains applicable, though careful tuning of sensitivity hyperparameters ( $\lambda$  gradient smoothing) may be required to preserve stability.

## 7. Conclusion

In this work, we introduced a gradient- and noise-aware adaptive timestep allocation strategy for DDIM sampling, aimed at improving sample quality without modifying the underlying model or increasing training cost. By leveraging signal-to-noise ratio (SNR) to globally allocate steps and gradient sensitivity of the predicted noise to locally adjust step sizes, our method adaptively refines the denoising trajectory, dedicating more computation to perceptually and semantically critical regions of the sampling process.

Extensive experiments on Stable Diffusion v1.5 and SDXL with the MS-COCO 2017 validation set demonstrate that our proposed method consistently outperforms baseline DDIM across Fréchet Inception Distance (FID), CLIP score, and visual quality benchmarks. Notably, our approach achieves comparable or superior quality at 20 steps to baseline DDIM at 50 steps, highlighting its potential for fast, high-quality image synthesis. While the adaptive allocation introduces only a minor computational overhead ( $\approx 3\text{--}5\%$ ), the resulting improvements in both perceptual fidelity and prompt alignment justify this trade-off.



Overall, our approach offers a plug-and-play, training-free enhancement to DDIM sampling that improves the efficiency–quality trade-off of diffusion models. This work paves the way for more intelligent, content-aware integration strategies in diffusion sampling, with promising applications in interactive content generation and real-time synthesis.

## References

- [1] M. Liu and Y. Hu, 'Application Potential of Stable Diffusion in Different Stages of Industrial Design', in *Artificial Intelligence in HCI*, vol. 14050, Cham: Springer Nature Switzerland, 2023, pp. 590–609.
- [2] J. Ho, A. Jain, and P. Abbeel, 'Denoising Diffusion Probabilistic Models', arXiv:arXiv:2006.11239., Dec. 16, 2020.
- [3] [1] J. Song, C. Meng, and S. Ermon, 'Denoising Diffusion Implicit Models', arXiv:arXiv:2010.02502., Oct. 05, 2022.
- [4] [C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, 'DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps', arXiv: arXiv:2206.00927., Oct. 13, 2022.
- [5] X. Wang, A.-D. Dinh, D. Liu, and C. Xu, 'Boosting Diffusion Models with an Adaptive Momentum Sampler', arXiv: arXiv:2308.11941., Aug. 23, 2023.
- [6] Z. Liu et al., 'Region-Adaptive Sampling for Diffusion Transformers', arXiv: arXiv:2502.10389., Feb. 14, 2025.
- [7] W. Zhao, L. Bai, Y. Rao, J. Zhou, and J. Lu, 'UniPC: A Unified Predictor-Corrector Framework for Fast Sampling of Diffusion Models', arXiv: arXiv:2302.04867., Oct. 17, 2023.
- [8] T. Karras, M. Aittala, T. Aila, and S. Laine, 'Elucidating the Design Space of Diffusion-Based Generative Models', arXiv: arXiv:2206.00364., Oct. 11, 2022.
- [9] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, 'Score-Based Generative Modeling through Stochastic Differential Equations', arXiv: arXiv:2011.13456., Feb. 10, 2021.
- [10] D. P. Kingma, T. Salimans, B. Poole, and J. Ho, 'Variational Diffusion Models', arXiv: arXiv:2107.00630., Apr. 14, 2023.
- [11] Ernst Hairer, Gerhard Wanner, Syvert P. Nørsett, 'Solving Ordinary Differential Equations I', vol. 8. in *Springer Series in Computational Mathematics*, vol. 8. Berlin, Heidelberg: Springer Berlin Heidelberg, 1993.
- [12] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, 'DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models', *Mach. Intell. Res.*, vol. 22, no. 4, pp. 730–751, Aug. 2025.
- [13] T. Salimans and J. Ho, 'Progressive Distillation for Fast Sampling of Diffusion Models', arXiv: arXiv:2202.00512., Jun. 07, 2022.
- [14] D. Watson, W. Chan, J. Ho, and M. Norouzi, 'Learning Fast Samplers for Diffusion Models by Differentiating Through Sample Quality', arXiv: arXiv:2202.05830., Feb. 11, 2022.
- [15] J. Sadowski, 'When data is capital: Datafication, accumulation, and extraction', *Big Data Soc.*, vol. 6, no. 1, pp. 1–12, 2019.
- [16] W. Jeong, K. Lee, H. Seo, and S. Y. Chun, 'Upsample What Matters: Region-Adaptive Latent Sampling for Accelerated Diffusion Transformers', arXiv: arXiv:2507.08422., Jul. 11, 2025.
- [17] T.-Y. Lin et al., 'Microsoft COCO: Common Objects in Context', arXiv: arXiv:1405.0312., Feb. 21, 2015.
- [18] V. Petsiuk et al., 'Human Evaluation of Text-to-Image Models on a Multi-Task Benchmark', arXiv: arXiv:2211.12112, Nov. 22, 2022.
- [19] S. Jayasumana, S. Ramalingam, A. Veit, D. Glasner, A. Chakrabarti, and S. Kumar, 'Rethinking FID: Towards a Better Evaluation Metric for Image Generation', arXiv: arXiv:2401.09603, Jan. 25, 2024.
- [20] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, 'Rethinking the Inception Architecture for Computer Vision', arXiv: arXiv:1512.00567, Dec. 11, 2015.
- [21] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi, 'CLIPScore: A Reference-free Evaluation Metric for Image Captioning', arXiv: arXiv:2104.08718, Mar. 23, 2022.