

# IMPROVING DENOISING DIFFUSION WITH EFFICIENT CONDITIONAL ENTROPY REDUCTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Diffusion models (DMs) have achieved significant success in generative modeling, but their iterative denoising process is computationally expensive. Training-free samplers, such as DPM-Solver, accelerate this process through gradient estimation-based numerical iterations. However, the mechanisms behind this acceleration remain insufficiently understood. In this paper, we demonstrate gradient estimation-based iterations enhance the denoising process by effectively *reducing the conditional entropy* of reverse transition distribution. Building on this analysis, we introduce streamlined denoising iterations for DMs that optimize conditional entropy in score-integral estimation to improve the denoising iterations. Experiments on benchmark pre-trained models validate our theoretical insights, demonstrating that numerical iterations based on conditional entropy reduction improve the reverse denoising diffusion process of DMs. [The code will be available.](#)

## 1 INTRODUCTION

It is well established that diffusion models (DMs) (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021b) have achieved significant success across various generative tasks, including image synthesis and editing (Dhariwal & Nichol, 2021; Meng et al., 2022), text-to-image synthesis (Ramesh et al., 2022), voice synthesis (Chen et al., 2021), and video generation (Ho et al., 2022). DMs consist of a forward diffusion process and a reverse denoising diffusion process. In the forward process, Gaussian noise is progressively injected into the data, perturbing the data distribution to collapse towards a standard Gaussian distribution by increasing conditional entropy. During training, the neural network is tasked with learning to reverse this process by minimizing the loss between the predicted and injected noise. Once the model is well-trained, high-quality samples can be synthesized by simulating the reverse-time denoising process associated with the forward noise-adding process.

However, a key limitation of DMs is the slow sequential nature of their iterative denoising process (Song et al., 2021a). To overcome this, training-free methods aim to accelerate denoising process by efficient numerical iterative algorithms without requiring additional training or costly optimization. Many of these methods focus on reformulating the denoising process as the solution of an ODE, allowing for accelerated sampling through numerical techniques. Such examples include PNDM (Liu et al., 2022), EDM Karras et al. (2022), DPM-Solver (Lu et al., 2022a), DEIS (Zhang & Chen, 2023), UniPC (Zhao et al., 2024), and DPM-Solver-v3 (Zheng et al., 2023a).

Despite the success of these numerical discretization techniques, the underlying mechanisms driving their acceleration remain inadequately understood. In particular, the reasons why iterations with similar orders of convergence result in varying levels of acceleration are not well explored. To address this gap, we reexamine the principles driving the accelerated denoising process. Our conditional entropy-based analysis reveals that effective iterations systematically reduce the conditional entropy of denoising transition distributions at each step, thereby directly contributing to a faster denoising process. This insight clarifies the mechanisms of gradient-based acceleration and provides a foundation for designing efficient denoising algorithms. Our main contributions are as follows:

- We introduce a novel perspective on *entropy reduction* in denoising diffusion of DMs, demonstrating that *gradient estimation-based iterations significantly accelerate the denoising process by effectively reducing conditional entropy*. Our theoretical analysis further reveals that denoising iterations *using data-prediction parameterization are more effective* than those using noise-prediction parameterization in minimizing conditional entropy.



(a) DPM-Solver++2m.

(b) DPM-Solver-v3-2m.

(c) Algorithm 1 (2m).

Figure 1: Random samples from Stable-Diffusion Rombach et al. (2022) with a classifier-free guidance scale 7.5, using 10 number of function evaluations (NFE) and text prompt “A beautiful castle beside a waterfall in the woods, by Josef Thoma, matte painting, trending on artstation HQ”.

- Building on our theoretical insights, we propose a denoising iteration method focused on efficient reducing conditional entropy in DMs. Unlike existing training-free methods, our approach improves the denoising process by lowering variance-driven conditional entropy during gradient-based iterations, which provides a simple yet effective improvement.
- Experiments on benchmark pre-trained models in both pixel and latent spaces validate our theoretical insights and demonstrate that our proposed method not only matches but often improves the reverse denoising diffusion process in DMs.

## 2 BACKGROUND

Diffusion models (DMs) define a Markov sequence  $\{x_t\}_{t \in [0, T]}$  in the forward process, starting with  $x_0$ , where  $x_0 \in \mathbb{R}^d$  is drawn from the clean data distribution  $q_0(x_0)$ . This sequence is pushed forward with increasing entropy until it approaches a standard Gaussian distribution via the transition kernel:  $q_t(x_t | x_0) = \mathcal{N}(x_t; \alpha_t x_0, \sigma_t^2 I)$ , where  $\sigma_t$  are smooth monotonic scalar functions w.r.t  $t$ . In DMs,  $\alpha_t$  and  $\sigma_t$  are called as the noise schedules,  $\alpha_t^2 / \sigma_t^2$  is called the signal-to-noise ratio (SNR) function. This transition kernel can be reformulated as the equivalent stochastic differential equation (SDE):

$$dx_t = f(t)x_t dt + g(t)d\omega_t, \quad x_0 \sim q_0(x_0), \quad (2.1)$$

where  $\omega_t$  denotes a standard Wiener process,  $f(t) := \frac{d \log \alpha_t}{dt}$ ,  $g^2(t) := \frac{d\sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2$  (Kingma et al., 2021). The reverse-time SDE of above forward diffusion process can be written as:

$$dx_t = [f(t)x_t - g^2(t)\nabla_x \log q_t(x_t)]dt + g(t)d\bar{\omega}_t, \quad x_T \sim q_T(x_T), \quad (2.2)$$

where  $\bar{\omega}_t$  represents another standard Wiener process. In score-based models (Song et al., 2021b), the diffusion (or probability flow) ordinary differential equation (ODE) used for efficient sampling is derived from the Fokker-Planck evolution equation of the probability density function as follows:

$$\frac{dx_t}{dt} = f(t)x_t - \frac{1}{2}g^2(t)\nabla_x \log q_t(x_t), \quad (2.3)$$

where the marginal distribution  $q_t(x_t)$  of  $x_t$  is equivalent to that of  $x_t$  in the SDE presented by Eq. (2.2). To train DMs, following the practiced in DDPM Ho et al. (2020), a neural network  $\epsilon_\theta(x_t, t)$  is parameterized to predict the noise  $\epsilon$  by minimizing the expectation of mean squared error as follows:

$$\mathbb{E}_{x_0 \sim q_0(x_0), \epsilon \sim \mathcal{N}(0, I), t \sim \mathcal{U}(0, T)} [w(t) \|\epsilon_\theta(\alpha_t x_0 + \sigma_t \epsilon, t) - \epsilon\|_2^2], \quad (2.4)$$

where  $\alpha_t^2 + \sigma_t^2 = 1$ ,  $w(t)$  is a weighting function that depends on the evolution time  $t$ . By substituting the trained noise prediction model  $\epsilon_\theta(x_t, t)$  with the scaled score function:  $-\sigma_t \nabla_x \log q_t(x_t)$ , sampling from DMs can be formulated by solving the diffusion ODE from  $T$  to 0 Song et al. (2021b):

$$\frac{dx_t}{dt} = f(t)x_t + \frac{g^2(t)}{2\sigma_t} \epsilon_\theta(x_t, t), \quad x_T \sim \mathcal{N}(0, \delta^2 I). \quad (2.5)$$

With a different parameterization, the data prediction model  $\mathbf{x}_\theta(\mathbf{x}_t, t)$  satisfies:  $\mathbf{x}_\theta(\mathbf{x}_t, t) = (\mathbf{x}_t - \sigma_t \epsilon_\theta(\mathbf{x}_t, t)) / \alpha_t$  (Kingma et al., 2021). This results in an equivalent ODE-based diffusion process:

$$\frac{d\mathbf{x}_t}{dt} = \left( f(t) + \frac{g^2(t)}{2\sigma_t^2} \right) \mathbf{x}_t - \alpha_t \frac{g^2(t)}{2\sigma_t^2} \mathbf{x}_\theta(\mathbf{x}_t, t). \quad (2.6)$$

### 3 CONDITIONAL ENTROPY REDUCTION AS A CATALYST FOR DENOISING DIFFUSION

By applying the *variation-of-constants* formula (Hale & Lunel, 2013) to ODEs (2.5) and (2.6), then

$$\mathbf{x}_t = e^{\int_s^t f(r) dr} \left( \int_s^t h_1(r) \epsilon_\theta(\mathbf{x}_r, r) dr + \mathbf{x}_s \right), \mathbf{x}_t = e^{h_2(t)} \left( - \int_s^t e^{-h_2(r)} \frac{\alpha_r g^2(r)}{2\sigma_r^2} \mathbf{x}_\theta(\mathbf{x}_r, r) dr + \mathbf{x}_s \right), \quad (3.1)$$

where  $h_1(r) := e^{-\int_s^r f(z) dz} \frac{g^2(r)}{2\sigma_r^2}$ ,  $h_2(r) := \int_s^r f(z) dz + \frac{g^2(z)}{2\sigma_z^2} dz$ , and  $\mathbf{x}_s$  represents the given initial value. Subsequently, this two diffusion ODEs have a unified semi-linear solution formula.

**Remark 1** Let the noise-prediction and data-prediction diffusion ODEs be defined by equations (2.5) and (2.6), respectively. A unified semi-linear solution formula for both ODEs is then given by:

$$\mathbf{f}(\mathbf{x}_t) - \mathbf{f}(\mathbf{x}_s) = \int_{\kappa(s)}^{\kappa(t)} \mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau)) d\tau, \quad (3.2)$$

where  $\psi(\kappa(t)) := t$ ,  $\{\mathbf{f}(\mathbf{x}_t) := \mathbf{x}_t / \alpha_t, \kappa(t) := \sigma_t / \alpha_t\}$  when  $\mathbf{d}_\theta$  represents the noise prediction model and  $\{\mathbf{f}(\mathbf{x}_t) := \mathbf{x}_t / \sigma_t, \kappa(t) := \alpha_t / \sigma_t\}$  when  $\mathbf{d}_\theta$  represents the data prediction model.

For brevity, we refer to Eq. (3.2) as the *score-integral* process, as the denoiser  $\mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau))$  is often trained to approximate the score function. Note that the semi-linear nature of diffusion ODEs can potentially reduce the sampling error of DMs Lu et al. (2022a;b); Zhang & Chen (2023). Unless otherwise specified, the following discussion defaults to noise prediction models.

#### 3.1 DENOISING ITERATIONS FORMULATED BY SCORE-INTEGRAL ESTIMATION

Denote  $h_{t_i} := \kappa(t_{i-1}) - \kappa(t_i)$ ,  $\mathbf{u}(\mathbf{x}_{t_{i-1}}) := \int_{\kappa(t_i)}^{\kappa(t_{i-1})} \mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau)) d\tau$  and  $\mathbf{d}_\theta^{(k)}(\mathbf{x}_{\psi(\tau)}, \psi(\tau)) := \frac{d^k \mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau))}{d\tau^k}$  as  $k$ -th order total derivative of  $\mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau))$  w.r.t.  $\tau$ . The Taylor expansion of  $\mathbf{d}_\theta(\mathbf{x}_{t_{i-1}}, t_{i-1})$  at  $\tau_{t_i}$  is

$$\mathbf{d}_\theta(\mathbf{x}_{t_{i-1}}, t_{i-1}) = \mathbf{d}_\theta(\mathbf{x}_{t_i}, t_i) + \sum_{k=1}^n \frac{h_{t_i}^k}{k!} \mathbf{d}_\theta^{(k)}(\mathbf{x}_{t_i}, t_i) + O(h_{t_i}^{n+1}). \quad (3.3)$$

Substituting this Taylor expansion into Eq. (3.2) to approximate  $\mathbf{u}(\mathbf{x}_{t_{i-1}})$  yields:

$$\tilde{\mathbf{u}}(\mathbf{x}_{t_{i-1}}) = h_{t_i} \mathbf{d}_\theta(\mathbf{x}_{t_i}, t_i) + \sum_{k=1}^n \frac{h_{t_i}^{k+1}}{(k+1)!} \mathbf{d}_\theta^{(k)}(\mathbf{x}_{t_i}, t_i) + O(h_{t_i}^{n+2}). \quad (3.4)$$

Beyond the transformations within the solving space, this Taylor-based approximation establishes a generalized numerical iterative framework for solving the score-integral in DMs. For instance, when  $n = 1$ , the truncated Taylor approximation reduces to the well-known *DDIM* iterative algorithm Song et al. (2021a), as follows:

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i). \quad (3.5)$$

where  $\tilde{\mathbf{x}}$  is obtained by the definition of  $\mathbf{f}(\tilde{\mathbf{x}})$ . Due to the lack of derivative information, higher-order algorithms can only be formulated by evaluating the derivatives. A widely used technique for evaluating derivatives is the finite difference (FD) method, which approximates  $\mathbf{d}_\theta^{(k)}(\cdot, \cdot)$  as follows ( $k \geq 1$ ):

$$\mathbf{d}_\theta^{(k)}(\mathbf{x}_t, t) = \frac{\mathbf{d}_\theta^{(k-1)}(\mathbf{x}_s, s) - \mathbf{d}_\theta^{(k-1)}(\mathbf{x}_t, t)}{\hat{h}_t} + O(\hat{h}_t). \quad (3.6)$$

Thus, a gradient estimation-based iteration can be obtained by truncating all higher-order derivatives:

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} F_\theta(s_i, t_i), \quad (3.7)$$

where  $F_\theta(s_i, t_i) := \frac{\mathbf{d}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i) - \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)}{\hat{h}_{t_i}}$ ,  $\hat{h}_{t_i} := \kappa(s_i) - \kappa(t_i)$ ,  $\hat{h}_{t_i} \neq 0$  and it is often satisfied that  $\hat{h}_{t_i} / h_{t_i} \leq 1$ .

### 3.2 CONDITIONAL ENTROPY REDUCTION IN DENOISING ITERATIONS

**Iterative Uncertainty Reduction: Theoretical Insights.** The semi-linear solution formula in Remark 1 provides a structured theoretical framework for analyzing the denoising diffusion process. By iteratively solving this formula, DMs refine noisy latent states closer to the data distribution. From an information-theoretic perspective, *each iteration progressively reduces uncertainty from intermediate representations by leveraging the structured denoising mechanism*. This uncertainty reduction can be formalized through the concept of *mutual information* between consecutive states Jaynes (1957):

$$I_p(\mathbf{x}_{t_i}; \mathbf{x}_{t_{i+1}}) = H_p(\mathbf{x}_{t_i}) - H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}), \quad (3.8)$$

where  $H_p(\mathbf{x}_{t_i})$  is the entropy of state  $\mathbf{x}_{t_i}$ , and  $H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  is the *conditional entropy* of  $\mathbf{x}_{t_i}$  given  $\mathbf{x}_{t_{i+1}}$ . The conditional entropy  $H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  *quantifies the uncertainty in  $\mathbf{x}_{t_i}$  after incorporating information from the subsequent state  $\mathbf{x}_{t_{i+1}}$* . A lower  $H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  indicates that the method effectively utilizes information from  $\mathbf{x}_{t_{i+1}}$  to refine the estimate of  $\mathbf{x}_{t_i}$ , driving the estimate of  $\mathbf{x}_{t_i}$  closer to the target data distribution. Practically, *this conditional entropy reduction aligns with the goal of minimizing reconstruction error during denoising*, improving the quality of generated samples. This theoretical insight not only elucidates the uncertainty reduction mechanism but also provides an optimization criterion for improving the denoising process.

**Conditional Entropy in Gaussian Approximations.** In practical implementations of DMs Ho et al. (2020); Song et al. (2021b), the reverse transition distribution  $p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}, \mathbf{x}_0)$  is commonly approximated as a Gaussian distribution under the *Markov assumption*. For brevity,  $p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}, \mathbf{x}_0)$  is often abbreviated as  $p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$ . Then, this reverse transition distribution can be expressed as

$$p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}) := p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}, \mathbf{x}_0) \approx \mathcal{N}(\boldsymbol{\mu}_{t_i}, \boldsymbol{\Sigma}_{t_i}), \quad (3.9)$$

where  $\boldsymbol{\mu}_{t_i}$  and  $\boldsymbol{\Sigma}_{t_i}$  are derived using Bayes' rule from the forward diffusion process. This Gaussian approximation is widely used for simplifying model training and theoretical analysis, despite potential deviations at extreme steps, as noted in prior works Song et al. (2021b); Luo (2022); Bao et al. (2022); Karras et al. (2022). Under this approximation, the conditional entropy  $H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  simplifies to

$$H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}) \approx \frac{d}{2} (\log 2\pi + 1) + \frac{1}{2} \log |\text{Var}(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})|, \quad (3.10)$$

where  $d$  is the dimensionality of  $\mathbf{x}$ , and  $\text{Var}(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  is the conditional variance. This expression provides a *tractable framework* for analyzing conditional entropy reduction during the denoising iteration, as it establishes a direct relationship between conditional entropy and the variance. Note that the conditional entropy  $H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$  is intrinsically tied to the conditional variance  $\text{Var}(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})$ :

$$H_p(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}}) \propto \log |\text{Var}(\mathbf{x}_{t_i} | \mathbf{x}_{t_{i+1}})|. \quad (3.11)$$

Thus, Eq. (3.11) suggests that minimizing variance directly optimizes conditional entropy reduction.

**Variance-Driven Conditional Entropy Reduction in Gradient-Based Iterations.** Building on the established relationship between conditional variance and entropy, we derive several analytical results that provide insights into the conditional entropy reduction achieved by gradient-based denoising iterations. For instance, under suitable conditions, our analysis suggests that gradient estimation-based iterations (Eq. (3.7)) can effectively drive significant reductions in conditional entropy.

To simplify the analysis, we assume that the estimated noise  $\epsilon_\theta(\cdot)$  at different timesteps is independent. While the forward process has the Markov property, our assumption mainly stems from *practical considerations in training*. Specifically, the training objective of Eq. (2.4) in DDPMs Ho et al. (2020) minimizes the mean squared error at each timestep independently, which aligns with this assumption. Although adopting a parameter-sharing setting across timesteps in the noise prediction model may involve a compromise on the assumption of independence, prior works Song et al. (2021a) indicate that these dependencies have minimal impact on model performance. This makes the independence assumption Ho et al. (2020) a *reasonable and practical surrogate* for theoretical analysis.

Under this independence assumption, we derive the Proposition 3.1, *with the proof in Appendix B.1*.

**Proposition 3.1** *The gradient-based denoising iteration in Eq. (3.7) tends to reduce conditional entropy more efficiently than the first-order iteration in Eq. (3.5) when  $\frac{h_{t_i}}{\tilde{h}_{t_i}} \in \left[1, \frac{4 \cdot \text{Var}(\epsilon_\theta(\tilde{\mathbf{x}}_{t_i}, t_i))}{\text{Var}(\epsilon_\theta(\tilde{\mathbf{x}}_{t_i}, s_i) + \text{Var}(\epsilon_\theta(\tilde{\mathbf{x}}_{t_i}, t_i))}\right]$ .*

Intuitively, this result reveals that gradient-based denoising iterations can achieve greater reductions in uncertainty compared to first-order methods when the step-size ratio is properly chosen. As the reverse process in DMs aims to estimate  $p(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_0)$  Ho et al. (2020); Luo (2022), we examine  $\text{Var}(\epsilon_\theta(\tilde{\mathbf{x}}_t, t) | \mathbf{x}_0)$  to capture the model's uncertainty in noise prediction conditioned on the clean data. For brevity, we denote this variance as  $\text{Var}(\mathbf{d}_\theta(\tilde{\mathbf{x}}_t, t))$  throughout the paper. Based on this consideration, we can establish the practical interval for Proposition 3.1 using the prior-like conditional variance.



**Proposition 3.2** *In the forward process of DMs, the clean data at each step can be expressed by  $x_0 = x_t/\alpha_t - \sigma_t/\alpha_t\epsilon$ . If we assume that  $\text{Var}(\epsilon_\theta(\tilde{x}_t, t)) \propto \sigma_t^2/\alpha_t^2$  to quantify the extent of deviation from the clean data. Under this prior-like assumption, we obtain  $\frac{\text{Var}(\epsilon_\theta(\tilde{x}_{s_i}, s_i))}{\text{Var}(\epsilon_\theta(\tilde{x}_{t_i}, t_i))} = \frac{\text{SNR}(t_i)}{\text{SNR}(s_i)}$ . Then, the relative condition of conditional entropy reduction in Propostion 3.1 is  $h_{t_i}/\hat{h}_{t_i} \in \left[1, \frac{4 \text{SNR}(s_i)}{\text{SNR}(t_i) + \text{SNR}(s_i)}\right]$ .*

Additionally, interpreting the denoising numerical iterative mechanisms through the lens of conditional entropy reduction offers deeper insights into accelerated denoising diffusion solvers, such as the widely recognized accelerated iterations in DPM-Solver Lu et al. (2022a) and EDM Karras et al. (2022). Building on this insight, we present the following proposition, with details in Appendix B.2.

**Proposition 3.3** *The exponential integrator-based iterations in DPM-Solver and the Heun iterations in EDM can be interpreted as specific instances of accelerated denoising mechanisms driven by conditional entropy reduction, thereby distinguishing them from conventional gradient-based methods.*

Finally, based on our comprehensive analysis of the differences in conditional entropy reduction between denoising iterations using data-prediction and noise-prediction parameterization, we derive the following conclusion. The detailed proof is provided in Appendix B.4.

**Proposition 3.4** *Assuming that the injected noise at different time steps in DM is mutually independent, denoising iterations using data-prediction parameterization are more effective at reducing conditional entropy than those using noise-prediction parameterization in a well-trained DM.*

Proposition 3.4 highlights the key advantage of data-prediction: it directly aligns with the target distribution  $x_0$ , bypassing the intermediate noise-to-data mapping  $\epsilon_t \mapsto x_t \mapsto x_0$ , which can accumulate errors, especially in late timesteps with high noise variance (or few-step sampling). By minimizing conditional entropy more effectively, data-prediction reduces uncertainty in  $x_0$  without relying on intermediate transformations. Nonetheless, this advantage is contingent on the training quality. If the model struggles to accurately predict  $x_0$ , noise-prediction parameterization, which treats timesteps more uniformly, may perform better in practice.

In summary, the perspective of conditional entropy reduction deepens our understanding of denoising mechanisms in diffusion model sampling, while the variance-driven approach provides valuable insights into the design of efficient denoising algorithms.

## 4 VARIANCE-DRIVEN EFFICIENT CONDITIONAL ENTROPY REDUCTION ITERATION

In this section, we elucidate the approach for improving both single-step and multi-step numerical iterations through conditional entropy reduction. Building on prior-like model variance assumptions, we derive several efficient iteration rules for conditional entropy reduction and establish the convergence orders of these iterations. Finally, we further optimize these iteration rules by refining the conditional variance with the actual state differences observed during the iterative process.

### 4.1 SINGLE-STEP ITERATION WITH CONDITIONAL ENTROPY REDUCTION

One key insight is that the model parameter  $\epsilon_\theta(\tilde{x}_{s_i}, s_i)$  can be further leveraged to enhance gradient estimation-based iteration, as observed in Eq. (3.7) and supported by conditional entropy analysis, without additional model parameters. Formally, the improvement iteration can be defined as follows:

$$\mathbf{f}(\tilde{x}_{t_{i-1}}) = \mathbf{f}(\tilde{x}_{t_i}) + h_{t_i}(\gamma_i \mathbf{d}_\theta(\tilde{x}_{s_i}, s_i) + (1 - \gamma_i) \mathbf{d}_\theta(\tilde{x}_{t_i}, t_i)) + \frac{h_{t_i}^2}{2} F_\theta(s_i, t_i), \quad (4.1)$$

where  $\gamma_i \in (0, 1]$ . This improved iteration shares the same limit state as the vanilla gradient estimation-based denoising iteration in Eq. (3.7) when  $s_i \rightarrow t_i$ . For convenience, we refer to the standard gradient estimation-based iteration as the **FD-based** iteration. In the analysis of conditional entropy, we can compare the different components of Eq. (3.7) and Eq. (4.1). Therefore, the variance of the key distinct components in each conditional distribution is as follows:

$$\text{Var}_{p_1} = h_{t_i}^2 \cdot \text{Var}(\epsilon_\theta(\tilde{x}_{t_i}, t_i)), \quad \text{Var}_{p_2}(\gamma_i) = h_{t_i}^2 \left( \gamma_i^2 \cdot \text{Var}(\epsilon_\theta(\tilde{x}_{s_i}, s_i)) + (1 - \gamma_i)^2 \cdot \text{Var}(\epsilon_\theta(\tilde{x}_{t_i}, t_i)) \right). \quad (4.2)$$

Then, the difference in conditional entropy between the two gradient estimation-based iterations is

$$\Delta H(p) = \frac{1}{2} \log \frac{\text{Var}_{p_2}(\gamma_i)}{\text{Var}_{p_1}} = \frac{1}{2} \log \left( 1 - 2\gamma_i + \gamma_i^2 + \gamma_i^2 \frac{\text{Var}(\epsilon_\theta(\tilde{x}_{s_i}, s_i))}{\text{Var}(\epsilon_\theta(\tilde{x}_{t_i}, t_i))} \right). \quad (4.3)$$

Due to  $\gamma_i \in (0, 1]$  and  $\text{SNR}(t_i) \leq \text{SNR}(s_i)$ ,  $\Delta H(p) \leq 0$  consistently holds under the assumption that  $\text{Var}(\epsilon_\theta(\tilde{x}_t, t)) \propto \sigma_t^2/\alpha_t^2$ . Therefore, this improved iteration can more efficiently reduce conditional entropy compared to the vanilla iteration by using subsequent model parameters in lower-variance regions as guidance. Consequently, based on  $\Delta H(p) \leq 0$ , we have the following proposition.

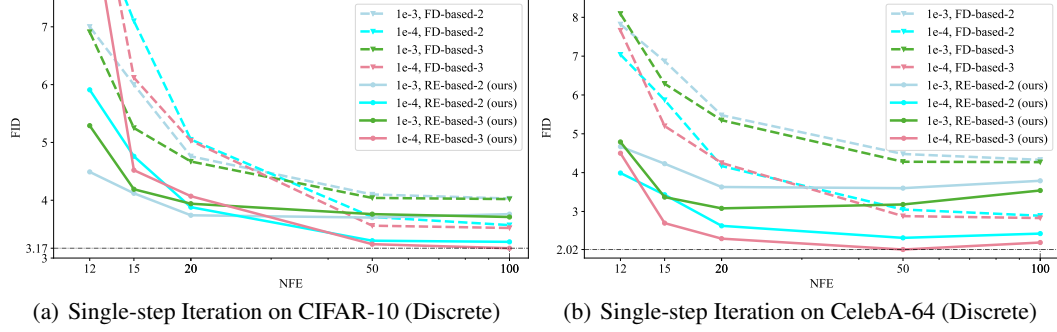


Figure 2: Comparisons of FID ↓ computed by **RE-based** and FD-based iterations demonstrate that efficient entropy reduction consistently enhances image quality across various ablation scenarios.

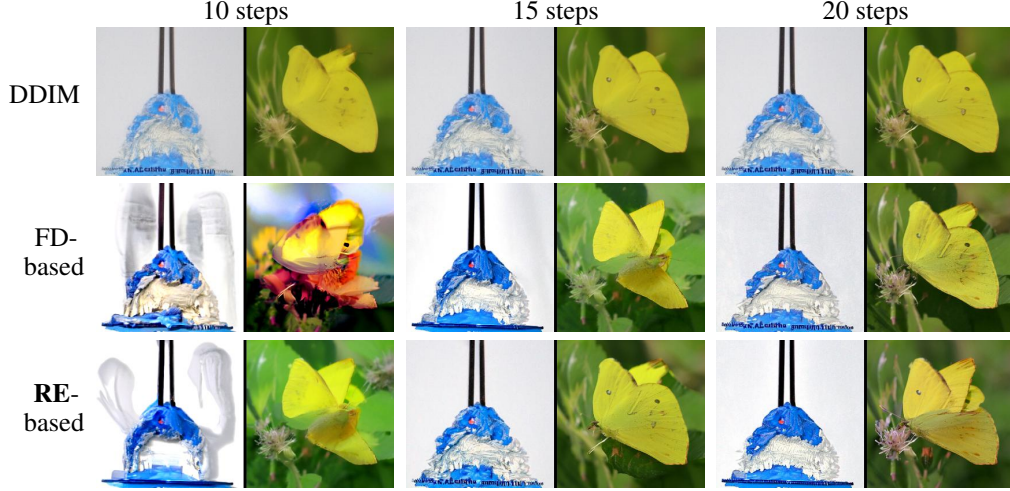


Figure 3: Samples were generated from a pre-trained DM on the ImageNet 256×256 dataset using **noise-prediction parameterization** with 10-20 single-step iterations. The sample results indicate that **RE-based** iterations can improve sample quality by reducing the conditional variance.

**Proposition 4.1** *The iteration specified in Eq. (4.1) consistently achieves a more efficient reduction in conditional entropy than the FD-based iteration. Then, an efficient improvement interval for  $\gamma_i$  is recommended as  $\left[ \frac{\text{SNR}(t_i)}{\text{SNR}(t_i) + \text{SNR}(s_i)}, \frac{\max[2 \cdot \text{SNR}(t_i), \text{SNR}(s_i)]}{\text{SNR}(t_i) + \text{SNR}(s_i)} \right]$ . For clarity, we identify this iteration that enhances denoising efficiency by reducing conditional entropy as the **RE-based** iteration.*

Accordingly, Proposition 4.1 demonstrates that the RE-based iteration can consistently surpass the FD-based iteration in reducing conditional entropy. In the following, we provide the convergence guarantees in Theorem 4.1 for the RE-based iteration, the proof is provided in Appendix C.2.

**Theorem 4.1** *If  $d_\theta(x_1, t)$  satisfies Assumption C.1, the RE-based iteration constitutes a globally convergent second-order iterative algorithm.*

Consequently, although the RE-based iteration in Eq. (4.1) shares the same order of convergence as the FD-based iteration, the primary distinction between the RE-based and FD-based iterations lies in their handling of conditional variance, which improves the denoising diffusion process by enabling more efficient conditional entropy reduction with the same model parameters.

#### 4.2 MULTI-STEP ITERATION WITH CONDITIONAL ENTROPY REDUCTION

This section focuses on the multi-step iteration with a step size determined by the two adjacent time points. Our discussion focuses on the conditional entropy reduction in multi-step iterations that leverage *data-prediction parameterization*, as this approach has demonstrated its superiority through our theoretical result presented in Proposition 3.4 and the empirical evidence from the earlier study in Lu et al. (2022b). The difference analysis of multi-step iterations is provided in Appendix B.3. Formally, the iteration with a step size determined by two adjoint time points can be written as:

$$f(\tilde{x}_{t_{i+1}}) = f(\tilde{x}_{t_i}) + h_{t_i} d_\theta(\tilde{x}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} B_\theta(t_i, t_{i+1}), \quad (4.4)$$

where  $B_\theta(t_i, t_{i+1}) := \frac{d_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) - d_\theta(\tilde{\mathbf{x}}_{t_{i+1}}, t_{i+1})}{h_{t_{i+1}}}$ . In this iteration, the step size  $|h_{t_i}|$  is smaller than the step size  $|h_{t_i} - h_{t_{i+1}}|$  used in gradient estimation-based iterations for the single-step case. As smaller step sizes reduce conditional entropy according to Eq. (3.10), the iteration in Eq. (4.4) offers greater potential for improving the denoising process compared to single-step counterparts. A straightforward improvement of the iteration in Eq. (4.4) can be formulated as follows:

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} B_\theta(t_i, l_i), \quad (4.5)$$

where  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, l_i) = \zeta_i \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + (1 - \zeta_i) \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_{i+1}}, t_{i+1})$  represents a linear interpolation of the model parameters between times  $t_i$  and  $t_{i+1}$ . Similarly, the implicit improvement approach is as follows:

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} B_\theta(s_i, t_i), \quad (4.6)$$

where  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i) = \zeta_i \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}) + (1 - \zeta_i) \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)$ . Note that both iterations can be unified as

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} \zeta_i \bar{B}_\theta(t_i; u_i), \quad (4.7)$$

where  $\bar{B}_\theta(t_i; u_i) = B_\theta(s_i, t_i)$  when  $u_i = s_i$ , and  $\bar{B}_\theta(t_i; u_i) = B_\theta(t_i, l_i)$  when  $u_i = l_i$ . Similar to the case of single-step iterations, the following conditional entropy relation also holds for multi-step iterations.

**Remark 2** The improved multi-step iterations in Eq. (4.7) reduce the conditional entropy of the vanilla multi-step iterations in Eq. (4.4) by leveraging model parameters from low-variance regions.

However, a key question arises: how should  $\zeta_i$  and  $\hat{h}_{t_i}$  be determined? In the data prediction model,  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)$  is designed to directly predict the clean data  $\mathbf{x}_0$  from the intermediate noisy data  $\tilde{\mathbf{x}}_{t_i}$ . Since  $\tilde{\mathbf{x}}_{t_i}$  is perturbed by Gaussian noise with a standard deviation  $\sigma_{t_i}$ ,  $\sigma_{t_i}$  reflects the amount of noise present at time step  $t_i$ . Then, for the interpolation of  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i)$ , we have the following proposition.

**Proposition 4.2** If assume that  $\text{Var}(\mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)) \propto \sigma_{t_i}^2$ , the minimizing variance can be achieved when  $\zeta_i = \frac{\sigma_{t_{i-1}}^2}{\sigma_{t_i}^2 + \sigma_{t_{i-1}}^2}$  for  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i)$ . For  $\mathbf{d}_\theta(\tilde{\mathbf{x}}_{l_i}, l_i)$ , the optimal choice of lower variance is  $\zeta_i = \frac{\sigma_{t_i}^2}{\sigma_{t_i}^2 + \sigma_{t_{i+1}}^2}$ .

One key insight is that we can further improve the denoising iteration in Eq. (4.4) with gradient estimation by incorporating  $B_\theta(t_i, s_i)$  and  $B_\theta(s_i, t_i)$  as follows:

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} (\eta_i B_\theta(s_i, t_i) + (1 - \eta_i) B_\theta(t_i, l_i)). \quad (4.8)$$

In Eq. (4.8),  $\eta_i$  determines the variance of the gradient term. From the perspective of conditional entropy reduction, we can reduce this variance by establishing an optimization objective that measures the differences between the corresponding states. Thus, in the next section, we will discuss the optimized  $\eta_i$  and  $\zeta_i$  based on the actual state differences observed during the iterative process.

### 4.3 IMPROVING RE-BASED ITERATIONS WITH ACTUAL STATE DIFFERENCES

In the previous sections, we derived the RE-based numerical iteration to reduce conditional entropy, grounded in the model's prior-like variance. In this section, the RE-based iteration is further optimized by refining the model variance with the actual state differences observed during the iterative process.

**Improving  $\zeta_i$  with Evolution State Differences.** Our goal is to refine  $\zeta_i$  in the RE-based iteration. What follows is the optimized  $\zeta_i$  by formulating an optimization objective. Denote  $G(\zeta_i) := \zeta_i \mathbf{d}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i) + (1 - \zeta_i) \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)$ , where  $\zeta_i \in (0, 1]$ . On one hand, we can rewrite the RE-based iteration as

$$\mathbf{f}(\tilde{\mathbf{x}}_{t_i}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) - h_{t_i} G(\zeta_i) - \frac{h_{t_i}^2}{2} F_\theta(s_i, t_i). \quad (4.9)$$

Notice that  $\tilde{\mathbf{x}}_{t_i}$  in Eq. (4.9) is determined by  $\zeta_i$ . On the other hand, we can consider  $\tilde{\mathbf{x}}_{t_{i-1}}$  as a starting point and perform an inverse iterative from  $t_{i-1}$  to  $t_i$  to approximate  $\tilde{\mathbf{x}}_{t_i}$ . The inverse iterative formula is

$$\mathbf{f}(\mathbf{x}_s) - \mathbf{f}(\mathbf{x}_t) = \int_{\kappa(t)}^{\kappa(s)} \mathbf{d}_\theta(\mathbf{x}_{\psi(\tau)}, \psi(\tau)) d\tau. \quad (4.10)$$

Similar to the score-integral iteration in Eq. (3.7), this inverse integral can be estimated by

$$\tilde{\Delta}_{t_i}^{\text{reverse}} = -h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}) + \frac{h_{t_i}^2}{2} F_\theta(s_i, t_{i-1}). \quad (4.11)$$

Based on equations (4.10) and (4.11), we obtain a new estimation  $\hat{\mathbf{x}}_{t_i}$  for the state  $\mathbf{x}_{t_i}$  as follows:

$$\mathbf{f}(\hat{\mathbf{x}}_{t_i}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) - h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}) + \frac{h_{t_i}^2}{2} F_\theta(s_i, t_{i-1}). \quad (4.12)$$

Drawing inspiration from equations (4.9) and (4.12), we can determine  $\zeta_i$  by minimizing the differences between two estimations. Then, the optimization objective for  $\zeta_i$  is defined as follows:

$$\min_{\zeta_i \in (0,1]} \mathcal{L}_1(\zeta_i) := \|(\tilde{\mathbf{x}}_{t_i} - \mathbf{x}_{t_i}) + (\hat{\mathbf{x}}_{t_i} - \mathbf{x}_{t_i})\|_F, \quad (4.13)$$

Directly solving this objective is challenging, as  $\mathbf{x}_{t_i}$  is unknown in practice. Fortunately, there exists an tractable error upper bound (EUB) for  $\mathcal{L}_1(\zeta_i)$ . Specifically, denote  $\mathcal{L}_{1s}(\zeta_i) := \|\tilde{\mathbf{x}}_{t_i} + \hat{\mathbf{x}}_{t_i}\|_F$ , we have

$$\mathcal{L}_1(\zeta_i) = \|\tilde{\mathbf{x}}_{t_i} + \hat{\mathbf{x}}_{t_i} - 2\mathbf{x}_{t_i}\|_F \leq \|\tilde{\mathbf{x}}_{t_i} + \hat{\mathbf{x}}_{t_i}\|_F + \|2\mathbf{x}_{t_i}\|_F = \mathcal{L}_{1s}(\zeta_i) + \|2\mathbf{x}_{t_i}\|_F, \quad (4.14)$$

where  $\|2\mathbf{x}_{t_i}\|_F$  can be viewed as a specific regularization term. Since  $\|2\mathbf{x}_{t_i}\|_F$  is independent of the target  $\zeta_i$ , we can optimize the vanilla  $\mathcal{L}_1(\zeta_i)$  by minimizing  $\mathcal{L}_{1s}(\zeta_i)$  according to Eq. (4.14). Then, the optimized  $\zeta_i$  can be obtained by solving  $\min \mathcal{L}_{1s}(\zeta_i)$  with a small regularization using  $\tilde{\mathbf{x}}_{t_i}$ . For example, denote  $P(\tilde{\mathbf{x}}_{t_{i-1}}^p) := \hat{\mathbf{x}}_{t_i} + \frac{\sigma_{t_i}}{\sigma_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}}^p - \sigma_{t_i} \frac{h_{t_i}^2}{2} F_\theta(s_i, t_i)$ , where  $\tilde{\mathbf{x}}_{t_{i-1}}^p$  can be obtained by prior RE-based iteration. Then, the simplified optimization objective  $\mathcal{L}_{1s}(\zeta_i)$  can be rewritten as:

$$\mathcal{L}_{1s}(\zeta_i) = \|P(\tilde{\mathbf{x}}_{t_{i-1}}^p) - \sigma_{t_i} h_{t_i} G(\zeta_i)\|_F. \quad (4.15)$$

**Practical Considerations.** In practice, the constraints on  $\zeta_i$  can hinder its computational efficiency. To address this, we adopt an optimization-guided streamlined approach for determining  $\zeta_i$ . Specifically, we observe that the optimization objective admits a closed-form solution, as presented in Lemma 4.1, when the constraints on  $\zeta_i$  are relaxed. These constraints are then satisfied by applying an activation function to map the closed-form solutions. This optimization-driven streamlined approach not only captures the actual differences in states during the iterative process, but also circumvents the computational cost of solving constrained optimization problems iteratively Boyd et al. (2011).

---

**Algorithm 1** Denoising Diffusion Sampling by Variance-Driven Conditional Entropy Reduction.

---

**Require:** initial value  $\mathbf{x}_T$ , time schedule  $\{t_i\}_{i=0}^N$ , model  $\mathbf{d}_\theta$ .

```

1:  $\tilde{\mathbf{x}}_{t_N} \leftarrow \mathbf{x}_T$ ,  $h_{t_i} \leftarrow \kappa(t_{i-1}) - \kappa(t_i)$ 
2: for  $i \leftarrow N$  to 1 do
3:    $\mathbf{f}(\tilde{\mathbf{x}}_{t_i}) \leftarrow \mathbf{f}(\tilde{\mathbf{x}}_{t_{i+1}}) + h_{t_{i+1}} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_{i+1}}, t_{i+1})$ 
4:    $\zeta_i = r_i h_{t_i}$ , where  $r_i$  is used to balance the prior-like variance, such as the log-SNR ratio.
5:    $\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) = \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} B_\theta(t_i, l_i)$ 
6:    $\eta_i = \text{Sigmoid}(|\eta_i^*|)$ , where  $\eta_i^*$  is computed using Eq. (4.19).
7:    $B_\theta(t_i) \leftarrow \frac{\eta_i}{2} B_\theta(s_i, t_i) + (1 - \frac{\eta_i}{2}) B_\theta(t_i, l_i)$ 
8:    $\zeta_i = \text{Sigmoid}(|\zeta_i^*| - \mu)$ , where  $\mu$  is the shift parameter, and  $\zeta_i^*$  is computed using Eq. (4.16).
9:    $\mathbf{f}(\tilde{\mathbf{x}}_{t_{i-1}}) \leftarrow \mathbf{f}(\tilde{\mathbf{x}}_{t_i}) + h_{t_i} \mathbf{d}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i) + \frac{h_{t_i}^2}{2} \zeta_i B_\theta(t_i)$ 
10: end for
return :  $\tilde{\mathbf{x}}_0$ .
```

---

**Lemma 4.1** The minimizing problem  $\min_{\zeta_i} \mathcal{L}_{1s}^2(\zeta_i)$  possesses the following closed-form solution:

$$\zeta_i^* = -\frac{\text{vec}^T(D_i) \text{vec}(\tilde{P}_i)}{\sigma_{t_i} h_{t_i} \text{vec}^T(D_i) \text{vec}(D_i)}, \quad (4.16)$$

where  $\tilde{P}_i := P(\tilde{\mathbf{x}}_{t_{i-1}}^p) - \sigma_{t_i} h_{t_i} \mathbf{x}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)$ ,  $D_i := \mathbf{x}_\theta(\tilde{\mathbf{x}}_{s_i}, s_i) - \mathbf{x}_\theta(\tilde{\mathbf{x}}_{t_i}, t_i)$ , and  $\text{vec}(\cdot)$  denotes the vectorization operation. The proof details can be found in Appendix C.3.

**Improving  $\eta_i$  with Balanced Difference Techniques.** Our goal is to refine the RE-based iteration by optimize  $\eta_i$  with the available information at current step. Denote  $\tilde{\Delta}_{t_i}^s = \eta_i F_\theta(t_{i-1}, t_i) + (1 - \eta_i) F_\theta(t_{i+1}, t_i)$ . We can define the estimation error of derivative at point  $\tau_{t_i}$  as  $E(t_{i-1}, t_i) := F_\theta(t_{i-1}, t_i) - \mathbf{d}_\theta^{(1)}(\mathbf{x}_{t_i}, t_i)$ . For balancing the estimation errors, we formulate the following optimization objective:

$$\min_{\eta_i \in (0,1]} \mathcal{L}_2(\eta_i) := \|\eta_i E(t_{i-1}, t_i) + (1 - \eta_i) E(t_{i+1}, t_i)\|_F. \quad (4.17)$$

We can rewrite  $\mathcal{L}_2(\eta_i)$  as  $\mathcal{L}_2(\eta_i) = \|\tilde{\Delta}_t^g - d_\theta^{(1)}(x_{t_i}, t_i)\|_F$ . Denote  $\mathcal{L}_{2s}(\eta_i) := \|\tilde{\Delta}_t^g\|_F$ . Then, we have

$$\mathcal{L}_2(\eta_i) = \|\tilde{\Delta}_t^g - \epsilon_\theta^{(1)}(x_{t_i}, t_i)\|_F \leq \mathcal{L}_{2s}(\eta_i) + \|d_\theta^{(1)}(x_{t_i}, t_i)\|_F, \quad (4.18)$$

where  $d_\theta^{(1)}(x_{t_i}, t_i)$  can be regarded as a specific regularization term independent of the target  $\eta_i$ . The optimized  $\eta_i$  can be obtained by minimizing the tractable EUB term  $\mathcal{L}_{2s}(\eta_i)$ . Similarly, for practical considerations, we employ optimization-guided streamlined approach for determining  $\eta_i$ . We first calculate the closed-form solution outlined in Lemma 4.2. The refined  $\eta_i$  is then obtained by mapping these solutions into the constrained space using an activation function, such as the Sigmoid function.

**Lemma 4.2** *The minimizing problem  $\min_{\eta_i} \mathcal{L}_{2s}^2(\eta_i)$  possesses the following closed-form solution:*

$$\eta_i^* = -\frac{\text{vec}^T(\tilde{F}_i)\text{vec}(F_\theta(t_{i+1}, t_i))}{\text{vec}^T(\tilde{F}_i)\text{vec}(\tilde{F}_i)}, \quad (4.19)$$

where  $\tilde{F}_i := F_\theta(t_{i+1}, t_i) - F_\theta(t_{i-1}, t_i)$ . The proof process is similar to that of Lemma 4.1.

Consequently, by integrating the optimized  $\zeta_i$  and  $\eta_i$  into the iterations of Eq. (4.8), we obtain the refined RE iterations. Algorithm 1 outlines this improved iteration process, which exhibits second-order global convergence, and the proof details are provided in Appendix C.4.

## 5 EXPERIMENTS

In this section, we experimentally validate our approach in both single-step and multi-step scenarios, demonstrating that variance-driven conditional entropy reduction improves the denoising process of pre-trained diffusion model in both pixel and latent spaces. This method effectively extends the capabilities of existing training-free ODE samplers without incurring additional computational overhead. We compare Algorithm 1 against the baseline methods on Stable Diffusion, as illustrated in Figure 1. More implementation details and additional results are provided in Appendix D.

### 5.1 SINGLE-STEP ITERATIONS

In the single-step iterations, we adopt DPM-Solver Lu et al. (2022a) as our baseline, focusing on denoising iterations based on noise prediction parameterization. Each step of the single-step mechanism only requires information from the starting point and prior to the endpoint. To ensure variance reduction in each step, we configure the step size ratio  $r_i$  following the effective interval defined in Proposition 4.1 and  $\gamma_i$  as specified in Proposition 3.2. As a specific instance of RE-based iterations (Proposition 3.3), DPM-Solver has demonstrated its effectiveness over traditional gradient-based iterations. We further validate RE-based iterations through experiments on CIFAR-10 Krizhevsky (2009), CelebA 64 Liu et al. (2015), and ImageNet 256 Deng et al. (2009), comparing them against solvers such as DDPM Ho et al. (2020), DDIM Song et al. (2021a), and Analytic-DDPM Bao et al. (2022). Results (Figures 2, 3, and 1) consistently show improved performance due to improved variance reduction. On CIFAR-10, the RE-based iteration achieves a 3.15 FID with only 84 NFEs, surpassing DDPM’s Ho et al. (2020) 3.17 FID with 1000 NFEs, improving quality while realizing approximately 10× acceleration. Additional comparisons are provided in Figure 4.

### 5.2 MULTI-STEP ITERATIONS

In the multi-step iterations, we primarily adopt DPM-Solver++ Lu et al. (2022b) as our baseline, focusing on denoising iterations based on data-prediction parameterization. Leveraging the multi-step mechanism enables us to utilize marginally more model information compared to single-step approaches. This allows us to optimize conditional variance through actual state differences, thereby circumventing the limitations imposed by prior-like variance assumptions. To demonstrate the effectiveness of efficient conditional entropy reduction in improving the denoising process of DMs, we propose Algorithm 1, which improves denoising diffusion sampling by leveraging the variance-driven approach aimed at minimizing actual state differences. Notably, DPM-Solver-v3 Zheng et al. (2023b) recently introduced a novel optimization-based parameterization scheme, distinct from data-prediction and noise-prediction parameterizations, achieving impressive sampling performance, particularly on CIFAR-10. Therefore, we adopt DPM-Solver-v3 as our baseline method for CIFAR-10 experiments, considering its demonstrated advantages in optimized parameterization on this dataset. We evaluated the RE-based iterations against widely-recognized benchmark solvers, including DPM-Solver++ Lu et al. (2022b), DEIS Zhang & Chen (2023), UniPC Zhao et al. (2024), and DPM-Solver-v3 Zheng et al. (2023a) on both CIFAR-10 and ImageNet 256 datasets. The experimental



results (Tables 3, 2) demonstrate that the variance-driven conditional entropy reduction consistently improves sampling performance. Furthermore, we validated the effectiveness of our approach on pre-trained models in the latent space, such as Stable Diffusion, with results illustrated in Figure 1.

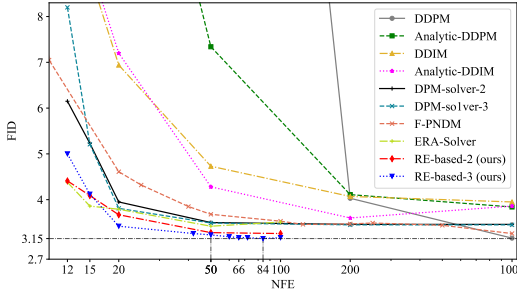


Figure 4: Comparisons of FID  $\downarrow$  for single-step RE iterations on discrete DMs in CIFAR-10.

Discrete	Continuous	Cond. EDM
3.17	2.55	1.79
DDPM	Hybrid PC	EDM
3.26	2.64	1.79
F-PNDM	DPM-Solver-v3	Heun's 2nd
<b>3.15</b>	<b>2.41</b>	<b>1.76</b>
RE-based	RE-based	RE-based

Table 1: The comparison for the performance limits of sampling methods on CIFAR-10 Krizhevsky (2009) indicates that RE-based iterations can further improve the denoising process.

Table 2: Quantitative results of the gradient estimation-based denoising iterations on ImageNet-256 Deng et al. (2009). We report the FID  $\downarrow$  evaluated on 10k samples for various NFEs.

Method	Model	NFE						
		5	6	8	10	12	15	20
DPM-Solver++	Guided-Diffusion ( $s = 2$ )	15.69	11.65	9.06	8.29	7.94	7.70	7.48
UniPC		15.03	11.30	9.07	8.36	8.01	7.71	7.47
DPM-Solver-v3		14.92	11.13	8.98	8.14	7.93	7.70	7.42
RE-based		<b>13.98</b>	<b>10.98</b>	<b>8.84</b>	<b>8.14</b>	<b>7.79</b>	<b>7.48</b>	<b>7.25</b>

Table 3: Quantitative results of the gradient estimation-based denoising iterations on CIFAR10. We report the FID  $\downarrow$  evaluated on 50k samples for the different NFEs. We directly borrow the results of reported in the original paper of other methods.

Method	Model	NFE						
		5	6	8	10	12	15	20
DEIS	ScoreSDE	15.37	\	\	4.17	\	3.37	2.86
DPM-Solver++		28.53	13.48	5.34	4.01	4.04	3.32	2.90
UniPC		23.71	10.41	5.16	3.93	3.88	3.05	2.73
DPM-Solver-v3		<b>12.76</b>	<b>7.40</b>	<b>3.94</b>	3.40	3.24	2.91	2.71
RE-based		13.54	8.56	4.11	<b>3.38</b>	<b>3.22</b>	<b>2.76</b>	<b>2.42</b>
Heun's 2nd	EDM	320.80	103.86	39.66	16.57	7.59	4.76	2.51
DPM-Solver++		24.54	11.85	4.36	2.91	2.45	2.17	2.05
UniPC		23.52	11.10	3.86	2.85	2.38	2.08	2.01
DPM-Solver-v3		12.21	8.56	3.50	2.51	2.24	2.10	2.02
RE-based		<b>11.82</b>	<b>8.30</b>	<b>3.46</b>	<b>2.48</b>	<b>2.21</b>	<b>2.07</b>	<b>2.01</b>

## CONCLUSIONS

In this paper, we introduce a novel framework that leverages variance-driven conditional entropy reduction to improve the sampling performance of pre-trained diffusion model without the need for retraining. Our theoretical analysis establishes that minimizing conditional entropy in the reverse process of diffusion models leads to more accurate and efficient denoising, which provides a principled foundation for optimizing this process. Building on these insights, we propose a Reduced Entropy (RE) approach for sampling of diffusion models, which improves the denoising process through efficient conditional variance minimization. Our method achieves state-of-the-art performance across multiple benchmark training-free methods, demonstrating promising improvements in both sampling speed and generation quality. While our approach yields promising results in image generation tasks, the full potential of conditional entropy reduction-based sampling methods remains to be explored.

## REFERENCES

- Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-DPM: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=0xiJLKH-ufZ>.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=NsMLjcFa080>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- Jack K Hale and Sjoerd M Verduyn Lunel. *Introduction to functional differential equations*, volume 99. Springer Science & Business Media, 2013.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- 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.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- A Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University of Tront*, 2009.
- Paul Langevin et al. Sur la théorie du mouvement brownien. *CR Acad. Sci. Paris*, 146(530-533):530, 1908.
- Shengmeng Li, Luping Liu, Zenghao Chai, Runnan Li, and Xu Tan. Era-solver: Error-robust adams solver for fast sampling of diffusion probabilistic models. *arXiv preprint arXiv:2301.12935*, 2023.
- Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=P1KWVd2yBkY>.
- Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=XVjTT1nw5z>.

- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pp. 3730–3738, 2015.
- 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*, 2022a.
- 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*, 2022b.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11461–11471, 2022.
- Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=aBsCjcPu\\_tE](https://openreview.net/forum?id=aBsCjcPu_tE).
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=FjNys5c7VyY>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- 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*, pp. 10684–10695, 2022.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your steps: Optimizing sampling schedules in diffusion models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=nBGBzV4It3>.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=TIIdIXIpzhoI>.