
Restart Sampling for Improving Generative Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

Generative processes that involve solving differential equations, such as diffusion models, frequently necessitate balancing speed and quality. ODE-based samplers are fast but plateau in performance while SDE-based samplers deliver higher sample quality at the cost of increased sampling time. We attribute this difference to sampling errors: ODE-samplers involve smaller discretization errors while stochasticity in SDE contracts accumulated errors. Based on these findings, we propose a novel sampling algorithm called *Restart* in order to better balance discretization errors and contraction. The sampling method alternates between adding substantial noise in additional forward steps and strictly following a backward ODE. Empirically, Restart sampler surpasses previous SDE and ODE samplers in both speed and accuracy. Restart not only outperforms the previous best SDE results, but also accelerates the sampling speed by 10-fold / 2-fold on CIFAR-10 / ImageNet 64×64 . In addition, it attains significantly better sample quality than ODE samplers within comparable sampling times. Moreover, Restart better balances text-image alignment/visual quality versus diversity than previous samplers in the large-scale text-to-image Stable Diffusion model pre-trained on LAION 512×512 .

1 Introduction

Deep generative models based on differential equations, such as diffusion models and Poisson flow generative models, have emerged as powerful tools for modeling high-dimensional data, from image synthesis [23, 9, 13, 27, 28] to biological data [10, 26]. These models use iterative backward processes that gradually transform a simple distribution (e.g., Gaussian in diffusion models) into a complex data distribution by solving a differential equations. The associated vector fields (or drifts) driving the evolution of the differential equations are predicted by neural networks. The resulting sample quality can be often improved by enhanced simulation techniques but at the cost of longer sampling times.

Prior samplers for simulating these backward processes can be categorized into two groups: ODE-samplers whose evolution beyond the initial randomization is deterministic, and SDE-samplers where the generation trajectories are stochastic. Several works [23, 12, 13] show that these samplers demonstrate their advantages in different regimes, as depicted in Fig. 1(b). ODE solvers [22, 16, 13] result in smaller discretization errors, allowing for decent sample quality even with larger step sizes (i.e., fewer number of function evaluations (NFE)). However, their generation quality plateaus rapidly. In contrast, SDE achieves better quality in the large NFE regime, albeit at the expense of increased sampling time. To better understand these differences, we theoretically analyze SDE performance: the stochasticity in SDE contracts accumulated error, which consists of both the discretization error along the trajectories as well as the approximation error of the learned neural network relative to the ground truth drift (e.g., score function in diffusion model [23]). The approximation error dominates when NFE is large (small discretization steps), explaining the SDE advantage in this regime. Intuitively, the stochastic nature of SDE helps "forget" accumulated errors from previous time steps.

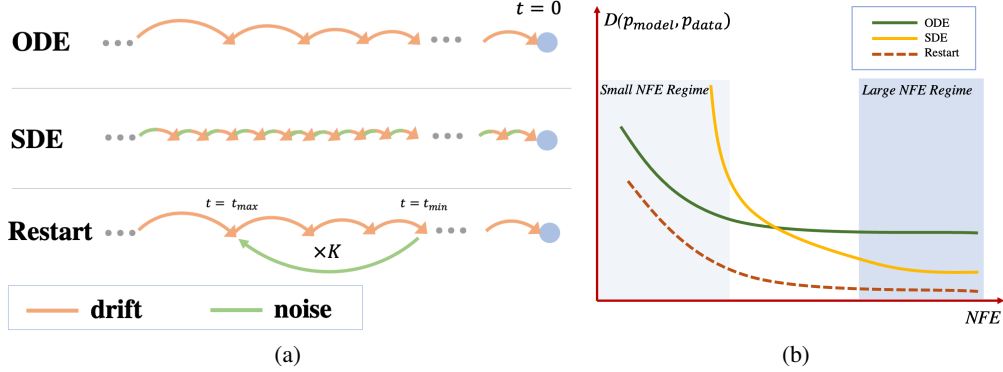


Figure 1: (a) Illustration of the implementation of drift and noise terms in ODE, SDE, and Restart. (b) Sample quality versus number of function evaluations (NFE) for different approaches. ODE (Green) provides fast speeds but attains only mediocre quality, even with a large NFE. SDE (Yellow) obtains good sample quality but necessitates substantial sampling time. In contrast to ODE and SDE, which have their own winning regions, Restart (Red) achieves the best quality across all NFEs.

Inspired by these findings, we propose a novel sampling algorithm called *Restart*, which combines the advantages of ODE and SDE. As illustrated in Fig. 1(a), the Restart sampling algorithm involves K repetitions of two subroutines in a pre-defined time interval: a *Restart forward process* that adds a substantial amount of noise, akin to "restarting" the original backward process, and a *Restart backward process* that runs the backward ODE. The Restart algorithm separates the stochasticity from the drifts, and the amount of added noise in the Restart forward process is significantly larger than the small single-step noise interleaving with drifts in previous SDEs such as [23, 13], thus amplifying the contraction effect on accumulated errors. By repeating the forward-backward cycle K times, the contraction effect introduced in each Restart iteration is further strengthened. The deterministic backward processes allow Restart to reduce discretization errors, thereby enabling step sizes comparable to ODE. To maximize the contraction effects in practice, we typically position the Restart interval towards the end of the simulation, where the accumulated error is larger. Additionally, we apply multiple Restart intervals to further reduce the initial errors in more challenging tasks.

Experimentally, Restart consistently surpasses previous ODE and SDE solvers in both quality and speed over a range of NFEs, datasets, and pre-trained models. Specifically, Restart accelerates the previous best-performing SDEs by $10\times$ fewer steps for the same FID score on CIFAR-10 using VP [23] ($2\times$ fewer steps on ImageNet 64×64 with EDM [13]), and outperforms ODE solvers even in the small NFE regime. When integrated into previous state-of-the-art pre-trained models, Restart further improves performance, achieving FID scores of 1.88 on unconditional CIFAR-10 with PFGM++ [28], and 1.36 on class-conditional ImageNet 64×64 with EDM. To the best of our knowledge, these are the best FID scores obtained on commonly used UNet architectures for diffusion models without additional training. We also apply Restart to the practical application of text-to-image Stable Diffusion model [19] pre-trained on LAION 512×512 . Restart more effectively balances text-image alignment/visual quality (measured by CLIP/Aesthetic scores) and diversity (measured by FID score) with a varying classifier-free guidance strength, compared to previous samplers.

Our contributions can be summarized as follows: (1) We investigate ODE and SDE solvers and theoretically demonstrate the contraction effect of stochasticity via an upper bound on the Wasserstein distance between generated and data distributions (Sec 3); (2) We introduce the Restart sampling, which better harnesses the contraction effect of stochasticity while allowing for fast sampling. The sampler results in a smaller Wasserstein upper bound (Sec 4); (3) Our experiments are consistent with the theoretical bounds and highlight Restart's superior performance compared to previous samplers on standard benchmarks in terms of both quality and speed. Additionally, Restart improves the trade-off between key metrics on the Stable Diffusion model (Sec 5).

2 Background on Generative Models with Differential Equations

Many recent successful generative models have their origin in physical processes, including diffusion models [9, 23, 13] and Poisson flow generative models [27, 28]. These models involve a forward process that transforms the data distribution into a chosen smooth distribution, and a backward process that iteratively reverses the forward process. For instance, in diffusion models, the forward

77 process is the diffusion process with no learned parameters:

$$dx = \sqrt{2\dot{\sigma}(t)\sigma(t)}dW_t,$$

78 where $\sigma(t)$ is a predefined noise schedule increasing with t , and $W_t \in \mathbb{R}^d$ is the standard Wiener
 79 process. For simplicity, we omit an additional scaling function for other variants of diffusion models
 80 as in EDM [13]. Under this notation, the marginal distribution at time t is the convolution of data
 81 distribution $p_0 = p_{\text{data}}$ and a Gaussian kernel, i.e., $p_t = p_0 * \mathcal{N}(\mathbf{0}, \sigma^2(t)\mathbf{I}_{d \times d})$. The prior distribution
 82 is set to $\mathcal{N}(\mathbf{0}, \sigma^2(T)\mathbf{I}_{d \times d})$ since p_T is approximately Gaussian with a sufficiently large T . Sampling
 83 of diffusion models is done via a reverse-time SDE [1] or a marginally-equivalent ODE [23]:

$$(SDE) \quad dx = -2\dot{\sigma}(t)\sigma(t)\nabla_x \log p_t(x)dt + \sqrt{2\dot{\sigma}(t)\sigma(t)}dW_t \quad (1)$$

$$(ODE) \quad dx = -\dot{\sigma}(t)\sigma(t)\nabla_x \log p_t(x)dt \quad (2)$$

84 where $\nabla_x \log p_t(x)$ in the drift term is the score of intermediate distribution at time t . W.l.o.g we set
 85 $\sigma(t) = t$ in the remaining text, as in [13]. Both processes progressively recover p_0 from the prior
 86 distribution p_T while sharing the same time-dependent distribution p_t . In practice, we train a neural
 87 network $s_\theta(x, t)$ to estimate the score field $\nabla_x \log p_t(x)$ by minimizing the denoising score-matching
 88 loss [25]. We then substitute the score $\nabla_x \log p_t(x)$ with $s_\theta(x, t)$ in the drift term of above backward
 89 SDE (Eq. (1))/ODE (Eq. (2)) for sampling.

90 Recent work inspired by electrostatics has not only challenged but also integrated diffusion models,
 91 notably PFGM/PFGM++, enhances performance in both image and antibody generation [27, 28, 10].
 92 They interpret data as electric charges in an augmented space, and the generative processes involve
 93 the simulations of differential equations defined by electric field lines. Similar to diffusion models,
 94 PFGMs train a neural network to approximate the electric field in the augmented space.

95 3 Explaining SDE and ODE performance regimes

96 To sample from the aforementioned generative models, a prevalent approach employs general-purpose
 97 numerical solvers to simulate the corresponding differential equations. This includes Euler and Heun’s
 98 2nd method [2] for ODEs (e.g., Eq. (2)), and Euler-Maruyama for SDEs (e.g., Eq. (1)). Sampling
 99 algorithms typically balance two critical metrics: (1) the quality and diversity of generated samples,
 100 often assessed via the Fréchet Inception Distance (FID) between generated distribution and data
 101 distribution [7] (lower is better), and (2) the sampling time, measured by the number of function
 102 evaluations (NFE). Generally, as the NFE decreases, the FID score tends to deteriorate across all
 103 samplers. This is attributed to the increased discretization error caused by using a larger step size in
 104 numerical solvers.

105 However, as illustrated in Fig. 1(b) and observed in previous works on diffusion models [23, 22, 13],
 106 the typical pattern of the quality vs time curves behaves differently between the two groups of
 107 samplers, ODE and SDE. When employing standard numerical solvers, ODE samplers attain a decent
 108 quality with limited NFEs, whereas SDE samplers struggle in the same small NFE regime. However,
 109 the performance of ODE samplers quickly reaches a plateau and fails to improve with an increase in
 110 NFE, whereas SDE samplers can achieve noticeably better sample quality in the high NFE regime.
 111 This dilemma raises an intriguing question: *Why do ODE samplers outperform SDE samplers in the*
 112 *small NFE regime, yet fall short in the large NFE regime?*

113 The first part of the question is relatively straightforward to address: given the same order of numerical
 114 solvers, simulation of ODE has significantly smaller discretization error compared to the SDE. For
 115 example, the first-order Euler method for ODE results in a local error of $O(\delta^2)$, whereas the first-order
 116 Euler-Maruyama method for SDEs yields a local error of $O(\delta^{\frac{3}{2}})$ (see e.g., Theorem 1 of [4]), where
 117 δ denotes the step size. As $O(\delta^{\frac{3}{2}}) \gg O(\delta^2)$, ODE simulations exhibit lower sampling errors than
 118 SDEs, likely causing the better sample quality with larger step sizes in the small NFE regime.

119 In the large NFE regime the step size δ shrinks and discretization errors become less significant
 120 for both ODEs and SDEs. In this regime it is the *approximation error* — error arising from an
 121 inaccurate estimation of the ground-truth vector field by the neural network s_θ — starts to dominate
 122 the sampling error. We denote the discretized ODE and SDE using the learned field s_θ as ODE_θ and
 123 SDE_θ , respectively. In the following theorem, we evaluate the total errors from simulating ODE_θ
 124 and SDE_θ within the time interval $[t_{\min}, t_{\max}] \subset [0, T]$. This is done via an upper bound on the
 125 Wasserstein-1 distance between the generated and data distributions at time t_{\min} . We characterize the
 126 accumulated initial sampling errors up until t_{\max} by total variation distances. Below we show that the

inherent stochasticity of SDEs aids in contracting these initial errors at the cost of larger additional sampling error in $[t_{\min}, t_{\max}]$. Consequently, SDE results in a smaller upper bound as the step size δ nears 0 (pertaining to the high NFE regime).

Theorem 1 (Informal). *Let t_{\max} be the initial noise level and p_t denote the true distribution at noise level t . Let $p_t^{ODE_\theta}, p_t^{SDE_\theta}$ denote the distributions of simulating ODE_θ, SDE_θ respectively. Assume that $\forall t \in [t_{\min}, t_{\max}], \|x_t\| < B/2$ for any x_t in the support of $p_t, p_t^{ODE_\theta}$ or $p_t^{SDE_\theta}$. Then*

$$\begin{aligned} W_1(p_{t_{\min}}^{ODE_\theta}, p_{t_{\min}}) &\leq B \cdot TV(p_{t_{\max}}^{ODE_\theta}, p_{t_{\max}}) + O(\delta + \epsilon_{\text{approx}}) \cdot (t_{\max} - t_{\min}) \\ \underbrace{W_1(p_{t_{\min}}^{SDE_\theta}, p_{t_{\min}})}_{\text{total error}} &\leq \underbrace{\left(1 - \lambda e^{-U}\right) B \cdot TV(p_{t_{\max}}^{SDE_\theta}, p_{t_{\max}})}_{\text{upper bound on contracted error}} + \underbrace{O(\sqrt{\delta t_{\max}} + \epsilon_{\text{approx}})(t_{\max} - t_{\min})}_{\text{upper bound on additional sampling error}} \end{aligned}$$

In the above, $U = BL_1/t_{\min} + L_1^2 t_{\max}^2/t_{\min}^2$, $\lambda < 1$ is a contraction factor, L_1 and ϵ_{approx} are uniform bounds on $\|ts_\theta(x_t, t)\|$ and the approximation error $\|t\nabla_x \log p_t(x) - ts_\theta(x, t)\|$ for all x_t, t , respectively. $O()$ hides polynomial dependency on various Lipschitz constants and dimension.

We defer the formal version and proof of Theorem 1 to Appendix A.1. As shown in the theorem, the upper bound on the total error can be decomposed into upper bounds on the *contracted error* and *additional sampling error*. $TV(p_{t_{\max}}^{ODE_\theta}, p_{t_{\max}})$ and $TV(p_{t_{\max}}^{SDE_\theta}, p_{t_{\max}})$ correspond to the initial errors accumulated from both approximation and discretization errors during the simulation of the backward process, up until time t_{\max} . In the context of SDE, this accumulated error undergoes contraction by a factor of $1 - \lambda e^{-BL_1/t_{\min} - L_1^2 t_{\max}^2/t_{\min}^2}$ within $[t_{\min}, t_{\max}]$, due to the effect of adding noise. Essentially, the minor additive Gaussian noise in each step can drive the generated distribution and the true distribution towards each other, thereby neutralizing a portion of the initial accumulated error.

The other term related to additional sampling error includes the accumulation of discretization and approximation errors in $[t_{\min}, t_{\max}]$. Despite the fact that SDE incurs a higher discretization error than ODE ($O(\sqrt{\delta})$ versus $O(\delta)$), the contraction effect on the initial error is the dominant factor impacting the upper bound in the large NFE regime where δ is small. Consequently, the upper bound for SDE is significantly lower. This provides insight into why SDE outperforms ODE in the large NFE regime, where the influence of discretization errors diminishes and the contraction effect dominates. In light of the distinct advantages of SDE and ODE, it is natural to ask whether we can combine their strengths. Specifically, can we devise a sampling algorithm that maintains a comparable level of discretization error as ODE, while also benefiting from, or even amplifying, the contraction effects induced by the stochasticity of SDE? In the next section, we introduce a novel algorithm, termed *Restart*, designed to achieve these two goals simultaneously.

4 Harnessing stochasticity with Restart

In this section, we present the Restart sampling algorithm, which incorporates stochasticity during sampling while enabling fast generation. We introduce the algorithm in Sec 4.1, followed by a theoretical analysis in Sec 4.2. Our analysis shows that Restart achieves a better Wasserstein upper bound compared to those of SDE and ODE in Theorem 1 due to greater contraction effects.

4.1 Method

In the Restart algorithm, simulation performs a few repeated back-and-forth steps within a pre-defined time interval $[t_{\min}, t_{\max}] \subset [0, T]$, as depicted in Figure 1(a). This interval is embedded into the simulation of the original backward ODE referred to as the *main backward process*, which runs from T to 0. In addition, we call the backward process within the Restart interval $[t_{\min}, t_{\max}]$ as the *Restart backward process*, to distinguish it from the main backward process.

Starting with samples at time t_{\min} , which are generated by following the main backward process, the Restart algorithm adds a large noise to transit the samples from t_{\min} to t_{\max} with the help of the forward process. The forward process does not require any evaluation of the neural network $s_\theta(x, t)$, as it is generally defined by an analytical perturbation kernel capable of transporting distributions from t_{\min} to t_{\max} . For instance, in the case of diffusion models, the perturbation kernel is $\mathcal{N}(\mathbf{0}, (\sigma(t_{\max})^2 - \sigma(t_{\min})^2)\mathbf{I}_{d \times d})$. The added noise in this step induces a more significant contraction compared to the small, interleaved noise in SDE. The step acts as if partially restarting the main backward process by increasing the time. Following this step, Restart simulates the backward ODE

174 from t_{\max} back to t_{\min} using the neural network predictions as in regular ODE. We repeat these
 175 forward-backward steps within $[t_{\min}, t_{\max}]$ interval K times in order to further derive the benefit from
 176 contraction. Specifically, the forward and backward processes in the i^{th} iteration ($i \in \{0, \dots, K-1\}$)
 177 proceed as follows:

$$\text{(Restart forward process)} \quad x_{t_{\max}}^{i+1} = x_{t_{\min}}^i + \varepsilon_{t_{\min} \rightarrow t_{\max}} \quad (3)$$

$$\text{(Restart backward process)} \quad x_{t_{\min}}^{i+1} = \text{ODE}_{\theta}(x_{t_{\max}}^{i+1}, t_{\max} \rightarrow t_{\min}) \quad (4)$$

178 where the initial $x_{t_{\min}}^0$ is obtained by simulating the ODE until t_{\min} : $x_{t_{\min}}^0 = \text{ODE}_{\theta}(x_T, T \rightarrow t_{\min})$,
 179 and the noise $\varepsilon_{t_{\min} \rightarrow t_{\max}}$ is sampled from the corresponding perturbation kernel from t_{\min} to t_{\max} .
 180 The Restart algorithm not only adds substantial noise in the Restart forward process (Eq. (3)), but
 181 also separates the stochasticity from the ODE, leading to a greater contraction effect, which we
 182 will demonstrate theoretically in the next subsection. For example, we set $[t_{\min}, t_{\max}] = [0.05, 0.3]$
 183 for VP model [13] on CIFAR-10. Repetitive use of the forward noise effectively mitigates errors
 184 accumulated from the preceding simulation up until t_{\max} . Furthermore, the Restart algorithm does
 185 not suffer from large discretization errors as it is mainly built from following the ODE in the Restart
 186 backward process (Eq. (4)). The effect is that the Restart algorithm is able to reduce the total sampling
 187 errors even in the small NFE regime. Detailed pseudocode for the Restart sampling process can be
 188 found in Algorithm 2, Appendix B.2.

189 4.2 Analysis

190 We provide a theoretical analysis of the Restart algorithm under the same setting as Theorem 1.
 191 In particular, we prove the following theorem, which shows that Restart achieves a much smaller
 192 contracted error in the Wasserstein upper bound than SDE (Theorem 1), thanks to the separation of
 193 the noise from the drift, as well as the large added noise in the Restart forward process (Eq. (3)). The
 194 repetition of the Restart cycle K times further leads to an enhanced reduction in the initial accumulated
 195 error. We denote the intermediate distribution in the i^{th} Restart iteration, following the discretized
 196 trajectories and the learned field s_{θ} , as $p_{t \in [t_{\min}, t_{\max}]}^{\text{Restart}_{\theta}(i)}$.

197 **Theorem 2 (Informal).** *Under the same setting of Theorem 1, assume $K \leq \frac{C}{L_2(t_{\max} - t_{\min})}$ for some*
 198 *universal constant C . Then*

$$\underbrace{W_1(p_{t_{\min}}^{\text{Restart}_{\theta}(K)}, p_{t_{\min}})}_{\text{total error}} \leq \underbrace{B \cdot (1 - \lambda)^K \text{TV}(p_{t_{\max}}^{\text{Restart}_{\theta}(0)}, p_{t_{\max}})}_{\text{upper bound on contracted error}} + \underbrace{(K + 1) \cdot O(\delta + \epsilon_{\text{approx}})(t_{\max} - t_{\min})}_{\text{upper bound on additional sampling error}}$$

199 where $\lambda < 1$ is the same contraction factor as Theorem 1. $O()$ hides polynomial dependency on
 200 various Lipschitz constants, dimension.

201 *Proof sketch.* To bound the total error, we introduce an auxiliary process $q_{t \in [t_{\min}, t_{\max}]}^{\text{Restart}_{\theta}(i)}$, which initiates
 202 from true distribution $p_{t_{\max}}$ and performs the Restart iterations. This process differs from $p_{t \in [t_{\min}, t_{\max}]}^{\text{Restart}_{\theta}(i)}$
 203 only in its initial distribution at t_{\max} ($p_{t_{\max}}$ versus $p_{t_{\max}}^{\text{Restart}_{\theta}(0)}$). We bound the total error by the following
 204 triangular inequality:

$$\underbrace{W_1(p_{t_{\min}}^{\text{Restart}_{\theta}(K)}, p_{t_{\min}})}_{\text{total error}} \leq \underbrace{W_1(p_{t_{\min}}^{\text{Restart}_{\theta}(K)}, q_{t_{\min}}^{\text{Restart}_{\theta}(K)})}_{\text{contracted error}} + \underbrace{W_1(q_{t_{\min}}^{\text{Restart}_{\theta}(K)}, p_{t_{\min}})}_{\text{additional sampling error}}$$

205 To bound the contracted error, we construct a careful coupling process between two individual
 206 trajectories sampled from $p_{t_{\min}}^{\text{Restart}_{\theta}(i)}$ and $q_{t_{\min}}^{\text{Restart}_{\theta}(i)}$, $i = 0, \dots, K-1$. Before these two trajectories
 207 converge, the Gaussian noise added in each Restart iteration is chosen to maximize the probability of
 208 the two trajectories mapping to an identical point, thereby maximizing the mixing rate in TV. After
 209 converging, the two processes evolve under the same Gaussian noise, and will stay converged as their
 210 drifts are the same. Lastly, we convert the TV bound to W_1 bound by multiplying B . The bound on
 211 the additional sampling error echoes the ODE analysis in Theorem 1: since the noise-injection and
 212 ODE-simulation stages are separate, we do not incur the higher discretization error of SDE. \square

213 We defer the formal version and proof of Theorem 2 to Appendix A.1. The first term in RHS bounds
 214 the contraction on the initial error at time t_{\max} and the second term reflects the additional sampling
 215 error of ODE accumulated during repetitive Restart iteration. Comparing the Wasserstein upper
 216 bound of SDE and ODE in Theorem 1, we make the following three observations: (1) Each Restart
 217 iteration has a smaller contraction factor $1 - \lambda$ compared to the one in SDE, since Restart separates

the large additive noise (Eq. (3)) from the ODE (Eq. (4)). (2) Restart backward process (Eq. (4)) has the same order of discretization error $O(\delta)$ as the ODE, compared to $O(\sqrt{\delta})$ in SDE. Hence, the Restart allows for small NFE due to ODE-level discretization error. (3) The contracted error further diminishes exponentially with the number of repetitions K though the additional error increases linearly with K . It suggests that there is a sweet spot of K that strikes a balance between reducing the initial error and increasing additional sampling error. Ideally, one should pick a larger K when the initial error at time t_{\max} greatly outweighs the incurred error in the repetitive backward process from t_{\max} to t_{\min} . We provide empirical evidences in Sec 5.2.

While Theorem 1 and Theorem 2 compare the upper bounds on errors of different methods, we provide empirical validation in Section 5.1 by directly calculating these errors, showing that the Restart algorithm indeed yields a smaller total error due to its superior contraction effects.

4.3 Practical considerations

The Restart algorithm offers several degrees of freedom, including the time interval $[t_{\min}, t_{\max}]$ and the number of restart iterations K . Here we provide a general recipe of parameter selection for practitioners, taking into account factors such as the complexity of the generative modeling tasks and the capacity of the network. Additionally, we discuss a stratified, multi-level Restart approach that further aids in reducing simulation errors along the whole trajectories for more challenging tasks.

Where to Restart? Theorem 2 shows that the Restart algorithm effectively reduces the accumulated error at time t_{\max} by a contraction factor in the Wasserstein upper bound. These theoretical findings inspire us to position the Restart interval $[t_{\min}, t_{\max}]$ towards the end of the main backward process, where the accumulated error is more substantial. In addition, our empirical observations suggest that a larger time interval $t_{\max} - t_{\min}$ is more beneficial for weaker/smaller architectures or more challenging datasets. Even though a larger time interval increases the additional sampling error, the benefits of the contraction significantly outweighs the downside, consistent with our theoretical predictions. We leave the development of principled approaches for optimal time interval selection for future works.

Multi-level Restart For challenging tasks that yield significant approximation errors, the backward trajectories may diverge substantially from the ground truth even at early stage. To prevent the ODE simulation from quickly deviating from the true trajectory, we propose implementing multiple Restart intervals in the backward process, alongside the interval placed towards the end. Empirically, we observe that a 1-level Restart is sufficient for CIFAR-10, while for more challenging datasets such as ImageNet [5], a multi-level Restart results in enhanced performance [5].

5 Experiments

In Sec 5.1, we first empirically verify the theoretical analysis relating to the Wasserstein upper bounds. We then evaluate the performance of different sampling algorithms on standard image generation benchmarks, including CIFAR-10 [14] and ImageNet 64×64 [5] in Sec 5.2. Lastly, we employ Restart on text-to-image generation, using Stable Diffusion model [19] pre-trained on LAION-5B [21] with resolution 512×512 , in Sec 5.3.

5.1 Additional sampling error versus contracted error

Our proposed Restart sampling algorithm demonstrates a higher contraction effect and smaller additional sampling error compared to SDE, according to Theorem 1 and Theorem 2. Although our theoretical analysis compares the upper bounds of the total, contracted and additional sampling errors, we further verify their relative values through a synthetic experiment.

Setup We construct a 20-dimensional dataset with 2000 points sampled from a Gaussian mixture, and train a four-layer MLP to approximate the score field $\nabla_x \log p_t$. We implement the ODE, SDE, and Restart methods within a predefined time range of $[t_{\min}, t_{\max}] = [1.0, 1.5]$, where the process outside this range is conducted via the first-order ODE. To compute various error types, we define the distributions generated by three methods as outlined in the proof of Theorem 2 and directly gauge the errors at end of simulation $t = 0$ instead of $t = t_{\min}$: (1) the generated distribution as p_0^{Sampler} , where $\text{Sampler} \in \{\text{ODE}_\theta, \text{SDE}_\theta, \text{Restart}_\theta(K)\}$; (2) an auxiliary distribution q_0^{Sampler} initiating from true distribution $p_{t_{\max}}$ at time t_{\max} . The only difference between p_0^{Sampler} and q_0^{Sampler} is their initial distribution at t_{\max} ($p_{t_{\max}}^{\text{ODE}_\theta}$ versus $p_{t_{\max}}$); and (3) the true data distribution p_0 . In line with Theorem 2,

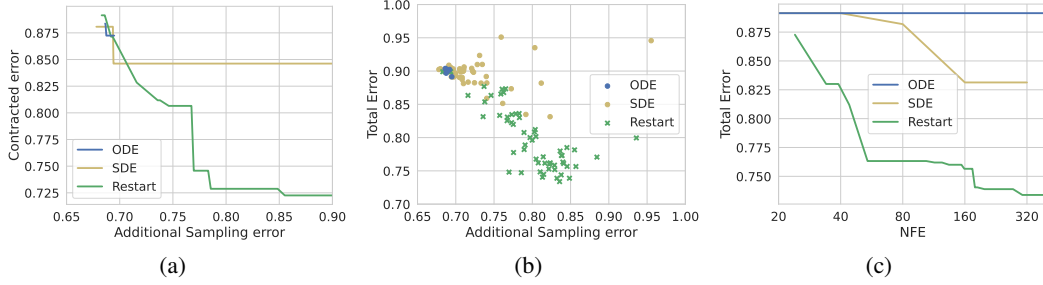


Figure 2: Additional sampling error versus (a) contracted error, where the Pareto frontier is plotted and (b) total error, where the scatter plot is provided. (c) Pareto frontier of NFE versus total error.

we use Wasserstein-1 distance $W_1(p_0^{\text{Sampler}}, q_0^{\text{Sampler}}) / W_1(q_0^{\text{Sampler}}, p_0)$ to measure the contracted error / additional sampling error, respectively. Ultimately, the total error corresponds to $W_1(p_0^{\text{Sampler}}, p_0)$. Detailed information about dataset, metric and model can be found in the Appendix C.5.

Results In our experiment, we adjust the parameters for all three processes and calculate the total, contracted, and additional sampling errors across all parameter settings. Figure 2(a) depicts the Pareto frontier of additional sampling error versus contracted error. We can see that Restart consistently achieves lower contracted error for a given level of additional sampling error, compared to both the ODE and SDE methods, as predicted by theory. In Figure 2(b), we observe that the Restart method obtains a smaller total error within the additional sampling error range of $[0.8, 0.85]$. During this range, Restart also displays a strictly reduced contracted error, as illustrated in Figure 2(a). This aligns with our theoretical analysis, suggesting that the Restart method offers a smaller total error due to its enhanced contraction effects. From Figure 2(c), Restart also strikes an better balance between efficiency and quality, as it achieves a lower total error at a given NFE.

5.2 Experiments on standard benchmarks

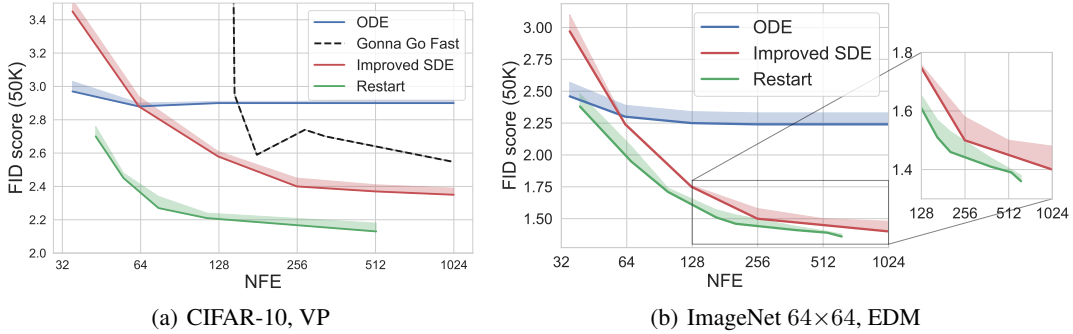


Figure 3: FID versus NFE on (a) unconditional generation on CIFAR-10 with VP; (b) class-conditional generation on ImageNet with EDM.

To evaluate the sample quality and inference speed, we report the FID score [7] (lower is better) on 50K samplers and the number of function evaluations (NFE). We borrow the pretrained VP/EDM/PFGM++ models on CIFAR-10 or ImageNet 64×64 from [13, 28]. We also use the EDM discretization scheme [13] (see Appendix B.1 for details) during sampling.

For the proposed Restart sampler, the hyperparameters include the number of steps in the main/Restart backward processes, the number of Restart iteration K , as well as the time interval $[t_{\min}, t_{\max}]$. We pick the t_{\min} and t_{\max} from the list of time steps in EDM discretization scheme with a number of steps 18. For example, for CIFAR-10 (VP) with $\text{NFE}=75$, we choose $t_{\min}=0.06$, $t_{\max}=0.30$, $K=10$, where $0.30/0.06$ is the 12th/14th time step in the EDM scheme. We also adopt EDM scheme for the Restart backward process in $[t_{\min}, t_{\max}]$. In addition, we apply the multi-level Restart strategy (Sec 4.3) to mitigate the error at early time steps for the more challenging ImageNet 64×64 . We provide the detailed Restart configurations in Appendix C.2.

For SDE, we compare with the previously best-performing stochastic samplers proposed by [13] (Improved SDE). We use their optimal hyperparameters for each dataset. We also report the FID

scores of the adaptive SDE [12] (**Gonna Go Fast**) on CIFAR-10 (VP). Since the vanilla reverse-diffusion SDE [23] has a significantly higher FID score, we omit its results from the main charts and defer them to Appendix D. For ODE samplers, we compare with the Heun’s 2nd order method [2] (**Heun**), which arguably provides an excellent trade-off between discretization errors and NFE [13]. To ensure a fair comparison, we use Heun’s method as the sampler in the main/Restart backward processes in Restart. One could employ faster, albeit more complex, ODE solvers such as the DPM-solver [16] to further accelerate Restart. We reserve this exploration for future works.

We report the FID score versus NFE in Figure 3(a) and Table 1 on CIFAR-10, and Figure 3(b) on ImageNet 64 × 64 with EDM. Our main findings are: **(1)** Restart outperforms other SDE or ODE samplers in balancing quality and speed, across datasets and models. As demonstrated in the figures, Restart achieves a 10-fold / 2-fold acceleration compared to previous best SDE results on CIFAR-10 (VP) / ImageNet 64 × 64 (EDM) at the same FID score. In comparison to ODE sampler (Heun), Restart obtains a better FID score, with the gap increasing significantly with NFE. **(2)** For stronger models such as EDM and PFGM++, Restart further improve over the ODE baseline on CIFAR-10. In contrast, the Improved SDE negatively impacts performance of EDM, as also observed in [13]. It suggests that Restart incorporates stochasticity more effectively. **(3)** Restart establishes new state-of-the-art FID scores for UNet architectures without additional training. In particular, Restart achieves FID scores of 1.36 on class-cond. ImageNet 64 × 64 with EDM, and 1.88 on uncond. CIFAR-10 with PFGM++.

Theorem 4 shows that each Restart iteration reduces the contracted errors while increasing the additional sampling errors in the backward process. In Fig. 4, we explore the choice of the number of Restart iterations K on CIFAR-10. We find that FID score initially improves and later worsens with increasing iterations K , with a smaller turning point for stronger EDM model. This supports the theoretical analysis that sampling errors will eventually outweigh the contraction benefits as K increases, and EDM only permits fewer Restart iterations due to smaller accumulated errors. It also suggests that, as a rule of thumb, we should apply greater Restart strength (*e.g.*, larger K) for weaker or smaller architectures and vice versa.

Table 1: Uncond. CIFAR-10 with EDM and PFGM++

	NFE	FID
<i>EDM-VP [13]</i>		
ODE (Heun)	63	1.97
	35	1.97
Improved SDE	63	2.27
	35	2.45
Restart	43	1.90
<i>PFGM++ [28]</i>		
ODE (Heun)	63	1.91
	35	1.91
Restart	43	1.88

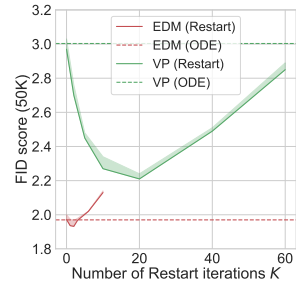
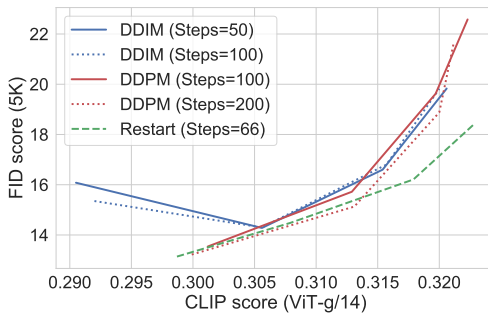
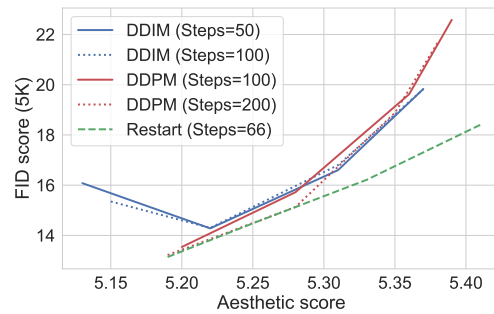


Figure 4: FID score with a varying number of Restart iterations K .

5.3 Experiments on large-scale text-to-image model



(a) FID versus CLIP score



(b) FID versus Aesthetic score

Figure 5: FID score versus (a) CLIP ViT-g/14 score and (b) Aesthetic score for text-to-image generation at 512 × 512 resolution, using Stable Diffusion v1.5 with a varying classifier-free guidance weight $w = 2, 3, 5, 8$.

We further apply Restart to the text-to-image Stable Diffusion v1.5¹ pre-trained on LAION-5B [21] at a resolution of 512 × 512. We employ the commonly used classifier-free guidance [8, 20] for sampling, wherein each sampling step entails two function evaluations – the conditional and unconditional predictions. Following [18, 20], we use the COCO [15] validation set for evaluation. We assess

¹<https://huggingface.co/runwayml/stable-diffusion-v1-5>



Figure 6: Visualization of generated images with classifier-free guidance weight $w = 8$, using four text prompts ("A photo of an astronaut riding a horse on mars.", "A raccoon playing table tennis", "Intricate origami of a fox in a snowy forest" and "A transparent sculpture of a duck made out of glass") and the **same** random seeds.

text-image alignment using the CLIP score [6] with the open-sourced ViT-g/14 [11], and measure diversity via the FID score. We also evaluate visual quality through the Aesthetic score, as rated by the LAION-Aesthetics Predictor V2 [24]. Following [17], we compute all evaluation metrics using 5K captions randomly sampled from the validation set and plot the trade-off curves between CLIP/Aesthetic scores and FID score, with the classifier-free guidance weight w in $\{2, 3, 5, 8\}$.

We compare with commonly used ODE sampler DDIM [22] and the stochastic sampler DDPM [9]. For Restart, we adopt the DDIM solver with 30 steps in the main backward process, and Heun in the Restart backward process, as we empirically find that Heun performs better than DDIM in the Restart. In addition, we select different sets of the hyperparameters for each guidance weight. For instance, when $w = 8$, we use $[t_{\min}, t_{\max}] = [0.1, 2]$, $K=2$ and 10 steps in Restart backward process. We defer the detailed Restart configuration to Appendix C.2, and the results of Heun to Appendix D.1.

As illustrated in Fig. 8(a) and Fig. 8(b), Restart achieves better FID scores in most cases, given the same CLIP/Aesthetic scores, using only 132 function evaluations (*i.e.*, 66 sampling steps). Remarkably, Restart achieves substantially lower FID scores than other samplers when CLIP/Aesthetic scores are high (*i.e.*, with larger w values). Conversely, Restart generally obtains a better text-image alignment/visual quality given the same FID. We also observe that DDPM generally obtains comparable performance with Restart in FID score when CLIP/Aesthetic scores are low, with Restart being more time-efficient. These findings suggest that Restart balances diversity (FID score) against text-image alignment (CLIP score) or visual quality (Aesthetic score) more effectively than previous samplers.

In Fig. 6, we visualize the images generated by Restart, DDIM and DDPM with $w = 8$. Compared to DDIM, the Restart generates images with superior details (*e.g.*, the rendition of duck legs by DDIM is less accurate) and visual quality. Compared to DDPM, Restart yields more photo-realistic images (*e.g.*, the astronaut). We provide extended of text-to-image generated samples in Appendix E.

6 Conclusion and future direction

In this paper, we introduce the Restart sampling for generative processes involving differential equations, such as diffusion models and PFGMs. By interweaving a forward process that adds a significant amount of noise with a corresponding backward ODE, Restart harnesses and even enhances the individual advantages of both ODE and SDE. Theoretically, Restart provides greater contraction effects of stochasticity while maintaining ODE-level discretization error. Empirically, Restart achieves a superior balance between quality and time, and improves the text-image alignment/visual quality and diversity trade-off in the text-to-image Stable Diffusion models.

A current limitation of the Restart algorithm is the absence of a principled way for hyperparameters selection, including the number of iterations K and the time interval $[t_{\min}, t_{\max}]$. At present, we adjust these parameters based on the heuristic that weaker/smaller models, or more challenging tasks, necessitate a stronger Restart strength. In the future direction, we anticipate developing a more principled approach to automating the selection of optimal hyperparameters for Restart based on the error analysis of models, in order to fully unleash the potential of the Restart framework.

References

- [1] Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- [2] Uri M. Ascher and Linda R. Petzold. Computer methods for ordinary differential equations and differential-algebraic equations. 1998.
- [3] Andrei N Borodin and Paavo Salminen. *Handbook of Brownian motion-facts and formulae*. Springer Science & Business Media, 2015.
- [4] Arnak S Dalalyan and Avetik Karagulyan. User-friendly guarantees for the langevin monte carlo with inaccurate gradient. *Stochastic Processes and their Applications*, 129(12):5278–5311, 2019.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [6] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Joseph Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Conference on Empirical Methods in Natural Language Processing*, 2021.
- [7] 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. In *NIPS*, 2017.
- [8] Jonathan Ho. Classifier-free diffusion guidance. *ArXiv*, abs/2207.12598, 2022.
- [9] Jonathan Ho, Ajay Jain, and P. Abbeel. Denoising diffusion probabilistic models. *ArXiv*, abs/2006.11239, 2020.
- [10] Chutian Huang, Zijing Liu, Shengyuan Bai, Linwei Zhang, Chencheng Xu, Zhe Wang, Yang Xiang, and Yuanpeng Xiong. Pf-abgen: A reliable and efficient antibody generator via poisson flow. *Machine Learning for Drug Discovery Workshop, International Conference on Learning Representations*, 2023.
- [11] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021.
- [12] Alexia Jolicoeur-Martineau, Ke Li, Remi Piche-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *ArXiv*, abs/2105.14080, 2021.
- [13] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *ArXiv*, abs/2206.00364, 2022.
- [14] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [15] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014.
- [16] 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. *arXiv preprint arXiv:2206.00927*, 2022.
- [17] Chenlin Meng, Ruiqi Gao, Diederik P. Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. *ArXiv*, abs/2210.03142, 2022.
- [18] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, 2021.
- [19] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021.

- 422 [20] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed
423 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo
424 Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic
425 text-to-image diffusion models with deep language understanding. *ArXiv*, abs/2205.11487,
426 2022.
- 427 [21] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,
428 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick
429 Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk,
430 and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text
431 models. *ArXiv*, abs/2210.08402, 2022.
- 432 [22] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *ArXiv*,
433 abs/2010.02502, 2020.
- 434 [23] Yang Song, Jascha Narain Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Er-
435 mon, and Ben Poole. Score-based generative modeling through stochastic differential equations.
436 *ArXiv*, abs/2011.13456, 2020.
- 437 [24] LAION-AI Team. Laion-aesthetics predictor v2. [https://github.com/
438 christophschuhmann/improved-aesthetic-predictor](https://github.com/christophschuhmann/improved-aesthetic-predictor), 2022.
- 439 [25] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural
440 Computation*, 23:1661–1674, 2011.
- 441 [26] Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E.
442 Eisenach, Woody Ahern, Andrew J. Borst, Robert J. Ragotte, Lukas F. Milles, Basile I. M.
443 Wicky, Nikita Hanikel, Samuel J. Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham
444 Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile
445 Mathieu, Regina Barzilay, T. Jaakkola, Frank DiMaio, Minkyung Baek, and David Baker.
446 Broadly applicable and accurate protein design by integrating structure prediction networks and
447 diffusion generative models. *bioRxiv*, 2022.
- 448 [27] Yilun Xu, Ziming Liu, Max Tegmark, and T. Jaakkola. Poisson flow generative models. *ArXiv*,
449 abs/2209.11178, 2022.
- 450 [28] Yilun Xu, Ziming Liu, Yonglong Tian, Shangyuan Tong, Max Tegmark, and T. Jaakkola.
451 Pfgm++: Unlocking the potential of physics-inspired generative models. *ArXiv*, abs/2302.04265,
452 2023.

Appendix

A Proofs of Main Theoretical Results

In this section, we provide proofs of our main results. We define below some crucial notations which we will use throughout. We use $\text{ODE}(\dots)$ to denote the backwards ODE under exact score $\nabla \log p_t(x)$. More specifically, given any $x \in \mathbb{R}^d$ and $s > r > 0$, let x_t denote the solution to the following ODE:

$$dx_t = -t \nabla \log p_t(x_t) dt. \quad (5)$$

$\text{ODE}(x, s \rightarrow r)$ is defined as "the value of x_r when initialized at $x_s = x$ ". It will also be useful to consider a "time-discretized ODE with drift $ts_\theta(x, t)$ ": let δ denote the discretization step size and let k denote any integer. Let δ denote a step size, let \bar{x}_t denote the solution to

$$d\bar{x}_t = -ts_\theta(x_{k\delta}, k\delta) dt, \quad (6)$$

where for any t , k is the unique integer such that $t \in ((k-1)\delta, k\delta]$. We verify that the dynamics of Eq. (6) is equivalent to the following discrete-time dynamics for $t = k\delta, k \in \mathbb{Z}$:

$$\bar{x}_{(k-1)\delta} = \bar{x}_{k\delta} - \frac{1}{2} \left(((k-1)\delta)^2 - (k\delta)^2 \right) s_\theta(x_{k\delta}, k\delta).$$

We similarly denote the value of \bar{x}_r when initialized at $\bar{x}_s = x$ as $\text{ODE}_\theta(x, s \rightarrow r)$. Analogously, we let $\text{SDE}(x, s \rightarrow r)$ and $\text{SDE}_\theta(x, s \rightarrow r)$ denote solutions to

$$\begin{aligned} dy_t &= -2t \nabla \log p_t(y_t) dt + \sqrt{2t} dB_t \\ d\bar{y}_t &= -2ts_\theta(\bar{y}_t, t) dt + \sqrt{2t} dB_t \end{aligned}$$

respectively. Finally, we will define the Restart_θ process as follows:

$$\begin{aligned} (\text{Restart}_\theta \text{ forward process}) \quad & x_{t_{\max}}^{i+1} = x_{t_{\min}}^i + \varepsilon_{t_{\min} \rightarrow t_{\max}}^i \\ (\text{Restart}_\theta \text{ backward process}) \quad & x_{t_{\min}}^{i+1} = \text{ODE}_\theta(x_{t_{\max}}^{i+1}, t_{\max} \rightarrow t_{\min}), \end{aligned} \quad (7)$$

where $\varepsilon_{t_{\min} \rightarrow t_{\max}}^i \sim \mathcal{N}(\mathbf{0}, (t_{\max}^2 - t_{\min}^2) \mathbf{I})$. We use $\text{Restart}_\theta(x, K)$ to denote $x_{t_{\min}}^K$ in the above processes, initialized at $x_{t_{\min}}^0 = x$. In various theorems, we will refer to a function $Q(r) : \mathbb{R}^+ \rightarrow [0, 1/2]$, defined as the Gaussian tail probability $Q(r) = \Pr(a \geq r)$ for $a \sim \mathcal{N}(0, 1)$.

A.1 Main Result

Theorem 3. [Formal version of Theorem 1] Let t_{\max} be the initial noise level. Let the initial random variables $\bar{x}_{t_{\max}} = \bar{y}_{t_{\max}}$, and

$$\begin{aligned} \bar{x}_{t_{\min}} &= \text{ODE}_\theta(\bar{x}_{t_{\max}}, t_{\max} \rightarrow t_{\min}) \\ \bar{y}_{t_{\min}} &= \text{SDE}_\theta(\bar{y}_{t_{\max}}, t_{\max} \rightarrow t_{\min}), \end{aligned}$$

Let p_t denote the true population distribution at noise level t . Let $p_t^{\text{ODE}_\theta}, p_t^{\text{SDE}_\theta}$ denote the distributions for x_t, y_t respectively. Assume that for all x, y, s, t , $s_\theta(x, t)$ satisfies $\|ts_\theta(x, t) - ts_\theta(x, s)\| \leq L_0|s - t|$, $\|ts_\theta(x, t)\| \leq L_1$, $\|ts_\theta(x, t) - ts_\theta(y, t)\| \leq L_2\|x - y\|$, and the approximation error $\|ts_\theta(x, t) - t \nabla \log p_t(x)\| \leq \epsilon_{\text{approx}}$. Assume in addition that $\forall t \in [t_{\min}, t_{\max}]$, $\|x_t\| < B/2$ for any x_t in the support of p_t , $p_t^{\text{ODE}_\theta}$ or $p_t^{\text{SDE}_\theta}$, and $K \leq \frac{C}{L_2(t_{\max} - t_{\min})}$ for some universal constant C . Then

$$\begin{aligned} W_1(p_{t_{\min}}^{\text{ODE}_\theta}, p_{t_{\min}}) &\leq B \cdot \text{TV}(p_{t_{\max}}^{\text{ODE}_\theta}, p_{t_{\max}}) \\ &\quad + e^{L_2(t_{\max} - t_{\min})} \cdot (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min}) \end{aligned} \quad (8)$$

$$\begin{aligned} W_1(p_{t_{\min}}^{\text{SDE}_\theta}, p_{t_{\min}}) &\leq B \cdot \left(1 - \lambda e^{-BL_1/t_{\min} - L_1^2 t_{\max}^2 / t_{\min}^2}\right) \text{TV}(p_{t_{\max}}^{\text{SDE}_\theta}, p_{t_{\max}}) \\ &\quad + e^{2L_2(t_{\max} - t_{\min})} \left(\epsilon_{\text{approx}} + \delta L_0 + L_2 \left(\delta L_1 + \sqrt{2\delta t_{\max}}\right)\right)(t_{\max} - t_{\min}) \end{aligned} \quad (9)$$

where $\lambda := 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right)$.

479 *Proof.* Let us define $x_{t_{\max}} \sim p_{t_{\max}}$, and let $x_{t_{\min}} = \text{ODE}(x_{t_{\max}}, t_{\max} \rightarrow t_{\min})$. We verify that $x_{t_{\min}}$
 480 has density $p_{t_{\min}}$. Let us also define $\hat{x}_{t_{\min}} = \text{ODE}_{\theta}(x_{t_{\max}}, t_{\max} \rightarrow t_{\min})$. We would like to bound
 481 the Wasserstein distance between $\bar{x}_{t_{\min}}$ and $x_{t_{\min}}$ (i.e., $p_{t_{\min}}^{\text{ODE}_{\theta}}$ and $p_{t_{\min}}$), by the following triangular
 482 inequality:

$$W_1(\bar{x}_{t_{\min}}, x_{t_{\min}}) \leq W_1(\bar{x}_{t_{\min}}, \hat{x}_{t_{\min}}) + W_1(\hat{x}_{t_{\min}}, x_{t_{\min}}) \quad (10)$$

483 By Lemma 2, we know that

$$\|\hat{x}_{t_{\min}} - x_{t_{\min}}\| \leq e^{(t_{\max}-t_{\min})L_2} (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}}) (t_{\max} - t_{\min}),$$

484 where we use the fact that $\|\hat{x}_{t_{\max}} - x_{t_{\max}}\| = 0$. Thus we immediately have

$$W_1(\hat{x}_{t_{\min}}, x_{t_{\min}}) \leq e^{(t_{\max}-t_{\min})L_2} (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}}) (t_{\max} - t_{\min}) \quad (11)$$

485 On the other hand,

$$\begin{aligned} W_1(\hat{x}_{t_{\min}}, \bar{x}_{t_{\min}}) &\leq B \cdot TV(\hat{x}_{t_{\min}}, \bar{x}_{t_{\min}}) \\ &\leq B \cdot TV(\hat{x}_{t_{\max}}, \bar{x}_{t_{\max}}) \end{aligned} \quad (12)$$

486 where the last equality is due to the data-processing inequality. Combining Eq. (11), Eq. (12) and the
 487 triangular inequality Eq. (10), we arrive at the upper bound for ODE (Eq. (8)). The upper bound for
 488 SDE (Eq. (9)) shares a similar proof approach. First, let $y_{t_{\max}} \sim p_{t_{\max}}$. Let $\hat{y}_{t_{\min}} = \text{SDE}_{\theta}(y_{t_{\max}}, t_{\max} \rightarrow$
 489 $t_{\min})$. By Lemma 5,

$$TV(\hat{y}_{t_{\min}}, \bar{y}_{t_{\min}}) \leq \left(1 - 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right) \cdot e^{-BL_1/t_{\min} - L_1^2 t_{\max}^2/t_{\min}^2}\right) \cdot TV(\hat{y}_{t_{\max}}, \bar{y}_{t_{\max}})$$

490 On the other hand, by Lemma 4,

$$\mathbb{E}[\|\hat{y}_{t_{\min}} - y_{t_{\min}}\|] \leq e^{2L_2(t_{\max}-t_{\min})} (\epsilon_{\text{approx}} + \delta L_0 + L_2(\delta L_1 + \sqrt{2\delta dt_{\max}})) (t_{\max} - t_{\min}).$$

491 The SDE triangular upper bound on $W_1(\bar{y}_{t_{\min}}, y_{t_{\min}})$ follows by multiplying the first inequality by B (to
 492 bound $W_1(\bar{y}_{t_{\min}}, \hat{y}_{t_{\min}})$) and then adding the second inequality (to bound $W_1(y_{t_{\min}}, \hat{y}_{t_{\min}})$). Notice
 493 that by definition, $TV(\hat{y}_{t_{\max}}, \bar{y}_{t_{\max}}) = TV(y_{t_{\max}}, \bar{y}_{t_{\max}})$. Finally, because of the assumption that
 494 $K \leq \frac{C}{L_2(t_{\max}-t_{\min})}$ for some universal constant, we summarize the second term in the Eq. (8) and
 495 Eq. (9) into the big O in the informal version Theorem 1. \square

496 **Theorem 4.** [Formal version of Theorem 2] Consider the same setting as Theorem 3. Let $p_{t_{\min}}^{\text{Restart}_{\theta}, i}$
 497 denote the distributions after i^{th} Restart iteration, i.e., the distribution of $\bar{x}_{t_{\min}}^i = \text{Restart}_{\theta}(\bar{x}_{t_{\min}}^0, i)$.
 498 Given initial $\bar{x}_{t_{\max}}^0 \sim p_{t_{\max}}^{\text{Restart}_{\theta}, 0}$, let $\bar{x}_{t_{\min}}^0 = \text{ODE}_{\theta}(\bar{x}_{t_{\max}}^0, t_{\max} \rightarrow t_{\min})$. Then

$$\begin{aligned} W_1(p_{t_{\min}}^{\text{Restart}_{\theta}, K}, p_{t_{\min}}) &\leq \underbrace{B \cdot (1 - \lambda)^K TV(p_{t_{\max}}^{\text{Restart}_{\theta}, 0}, p_{t_{\max}})}_{\text{upper bound on contracted error}} \\ &\quad + \underbrace{e^{(K+1)L_2(t_{\max}-t_{\min})} (K+1) (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}}) (t_{\max} - t_{\min})}_{\text{upper bound on additional sampling error}} \end{aligned} \quad (13)$$

499 where $\lambda = 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right)$.

500 *Proof.* Let $x_{t_{\max}}^0 \sim p_{t_{\max}}$. Let $x_{t_{\min}}^K = \text{Restart}(x_{t_{\min}}^0, K)$. We verify that $x_{t_{\min}}^K$ has density $p_{t_{\min}}$. Let us
 501 also define $\hat{x}_{t_{\min}}^0 = \text{ODE}_{\theta}(x_{t_{\max}}^0, t_{\max} \rightarrow t_{\min})$ and $\hat{x}_{t_{\min}}^K = \text{Restart}_{\theta}(\hat{x}_{t_{\min}}^0, K)$.

502 By Lemma 1,

$$\begin{aligned} TV(\bar{x}_{t_{\min}}^K, \hat{x}_{t_{\min}}^K) &\leq \left(1 - 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right)\right)^K TV(\bar{x}_{t_{\min}}^0, \hat{x}_{t_{\min}}^0) \\ &\leq \left(1 - 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right)\right)^K TV(\bar{x}_{t_{\max}}^0, \hat{x}_{t_{\max}}^0) \\ &= \left(1 - 2Q\left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}}\right)\right)^K TV(\bar{x}_{t_{\max}}^0, x_{t_{\max}}^0) \end{aligned}$$

503 The second inequality holds by data processing inequality. The above can be used to bound the
 504 1-Wasserstein distance as follows:

$$W_1(\bar{x}_{t_{\min}}^K, \hat{x}_{t_{\min}}^K) \leq B \cdot TV(\bar{x}_{t_{\min}}^K, \hat{x}_{t_{\min}}^K) \leq \left(1 - 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right)\right)^K TV(\bar{x}_{t_{\max}}^0, x_{t_{\max}}^0) \quad (14)$$

505 On the other hand, using Lemma 3,

$$W_1(x_{t_{\min}}^K, \hat{x}_{t_{\min}}^K) \leq \|x_{t_{\min}}^K - \hat{x}_{t_{\min}}^K\| \leq e^{(K+1)L_2(t_{\max}-t_{\min})} (K+1) (\delta(L_2L_1 + L_0) + \epsilon_{approx}) (t_{\max} - t_{\min}) \quad (15)$$

506 We arrive at the result by combining the two bounds above (Eq. (14), Eq. (15)) with the following
 507 triangular inequality,

$$W_1(\bar{x}_{t_{\min}}^K, x_{t_{\min}}^K) \leq W_1(\bar{x}_{t_{\min}}^K, \hat{x}_{t_{\min}}^K) + W_1(\hat{x}_{t_{\min}}^K, x_{t_{\min}}^K)$$

508 □

509 A.2 Mixing under Restart with exact ODE

510 **Lemma 1.** Consider the same setup as Theorem 4. Consider the Restart_θ process defined in
 511 equation 7. Let

$$\begin{aligned} x_{t_{\min}}^i &= \text{Restart}_\theta(x_{t_{\min}}^0, i) \\ y_{t_{\min}}^i &= \text{Restart}_\theta(y_{t_{\min}}^0, i). \end{aligned}$$

512 Let $p_t^{\text{Restart}_\theta(i)}$ and $q_t^{\text{Restart}_\theta(i)}$ denote the densities of x_t^i and y_t^i respectively. Then

$$TV(p_{t_{\min}}^{\text{Restart}_\theta(K)}, q_{t_{\min}}^{\text{Restart}_\theta(K)}) \leq (1 - \lambda)^K TV(p_{t_{\min}}^{\text{Restart}_\theta(0)}, q_{t_{\min}}^{\text{Restart}_\theta(0)}),$$

513 where $\lambda = 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right)$.

514 *Proof.* Conditioned on $x_{t_{\min}}^i, y_{t_{\min}}^i$, let $x_{t_{\max}}^{i+1} = x_{t_{\min}}^i + \sqrt{t_{\max}^2 - t_{\min}^2} \xi_i^x$ and $y_{t_{\max}}^{i+1} = y_{t_{\min}}^i +$
 515 $\sqrt{t_{\max}^2 - t_{\min}^2} \xi_i^y$. We now define a coupling between $x_{t_{\min}}^{i+1}$ and $y_{t_{\min}}^{i+1}$ by specifying the joint dis-
 516 tribution over ξ_i^x and ξ_i^y .

517 If $x_{t_{\min}}^i = y_{t_{\min}}^i$, let $\xi_i^x = \xi_i^y$, so that $x_{t_{\min}}^{i+1} = y_{t_{\min}}^{i+1}$. On the other hand, if $x_{t_{\min}}^i \neq y_{t_{\min}}^i$, let $x_{t_{\max}}^{i+1}$ and $y_{t_{\max}}^{i+1}$
 518 be coupled as described in the proof of Lemma 7, with $x' = x_{t_{\max}}^{i+1}, y' = y_{t_{\max}}^{i+1}, \sigma = \sqrt{t_{\max}^2 - t_{\min}^2}$.
 519 Under this coupling, we verify that,

$$\begin{aligned} & \mathbb{E} [\mathbb{1} \{x_{t_{\min}}^{i+1} \neq y_{t_{\min}}^{i+1}\}] \\ & \leq \mathbb{E} [\mathbb{1} \{x_{t_{\max}}^{i+1} \neq y_{t_{\max}}^{i+1}\}] \\ & \leq \mathbb{E} \left[\left(1 - 2Q \left(\frac{\|x_{t_{\min}}^i - y_{t_{\min}}^i\|}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right) \right) \mathbb{1} \{x_{t_{\min}}^i \neq y_{t_{\min}}^i\} \right] \\ & \leq \left(1 - 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right) \right) \mathbb{E} [\mathbb{1} \{x_{t_{\min}}^i \neq y_{t_{\min}}^i\}]. \end{aligned}$$

520 Applying the above recursively,

$$\mathbb{E} [\mathbb{1} \{x_{t_{\min}}^K \neq y_{t_{\min}}^K\}] \leq \left(1 - 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right) \right)^K \mathbb{E} [\mathbb{1} \{x_{t_{\min}}^0 \neq y_{t_{\min}}^0\}].$$

521 The conclusion follows by noticing that $TV(p_{t_{\min}}^{\text{Restart}_\theta(K)}, q_{t_{\min}}^{\text{Restart}_\theta(K)}) \leq Pr(x_{t_{\min}}^K \neq y_{t_{\min}}^K) =$
 522 $\mathbb{E} [\mathbb{1} \{x_{t_{\min}}^K \neq y_{t_{\min}}^K\}]$, and by selecting the initial coupling so that $Pr(x_{t_{\min}}^0 \neq y_{t_{\min}}^0) =$
 523 $TV(p_{t_{\min}}^{\text{Restart}_\theta(0)}, q_{t_{\min}}^{\text{Restart}_\theta(0)})$. □

524 A.3 W_1 discretization bound

525 **Lemma 2** (Discretization bound for ODE). *Let $x_{t_{\min}} = \text{ODE}(x_{t_{\max}}, t_{\max} \rightarrow t_{\min})$ and let $\bar{x}_{t_{\min}} =$*
 526 *$\text{ODE}_\theta(\bar{x}_{t_{\max}}, t_{\max} \rightarrow t_{\min})$. Assume that for all x, y, s, t , $s_\theta(x, t)$ satisfies $\|ts_\theta(x, t) - ts_\theta(x, s)\| \leq$*
 527 *$L_0|s - t|$, $\|ts_\theta(x, t)\| \leq L_1$ and $\|ts_\theta(x, t) - ts_\theta(y, t)\| \leq L_2\|x - y\|$. Then*

$$\|x_{t_{\min}} - \bar{x}_{t_{\min}}\| \leq e^{(t_{\max} - t_{\min})L_2} (\|x_{t_{\max}} - \bar{x}_{t_{\max}}\| + (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min}))$$

528 *Proof.* Consider some fixed arbitrary k , and recall that δ is the step size. Recall that by definition of
 529 ODE and ODE_θ , for $t \in ((k-1)\delta, k\delta]$,

$$\begin{aligned} dx_t &= -t \nabla \log p_t(x_t) dt \\ d\bar{x}_t &= -ts_\theta(\bar{x}_{k\delta}, k\delta) dt. \end{aligned}$$

530 For $t \in [t_{\min}, t_{\max}]$, let us define a time-reversed process $x_t^\leftarrow := x_{-t}$. Let $v(x, t) := \nabla \log p_{-t}(x)$.
 531 Then for $t \in [-t_{\max}, -t_{\min}]$

$$dx_t^\leftarrow = tv(x_t^\leftarrow, t) ds.$$

532 Similarly, define $\bar{x}_t^\leftarrow := \bar{x}_{-t}$ and $\bar{v}(x, t) := s_\theta(x, -t)$. It follows that

$$d\bar{x}_t^\leftarrow = t\bar{v}(\bar{x}_{k\delta}^\leftarrow, k\delta) ds,$$

533 where k is the unique (negative) integer satisfying $t \in [k\delta, (k+1)\delta)$. Following these definitions,

$$\begin{aligned} & \frac{d}{dt} \|x_t^\leftarrow - \bar{x}_t^\leftarrow\| \\ & \leq \|tv(x_t^\leftarrow, t) - t\bar{v}(\bar{x}_t^\leftarrow, t)\| \\ & \quad + \|t\bar{v}(\bar{x}_t^\leftarrow, t) - t\bar{v}(\bar{x}_t^\leftarrow, k\delta)\| \\ & \quad + \|t\bar{v}(\bar{x}_t^\leftarrow, k\delta) - t\bar{v}(\bar{x}_{k\delta}^\leftarrow, k\delta)\| \\ & \leq \epsilon_{\text{approx}} + L_2 \|x_t^\leftarrow - \bar{x}_t^\leftarrow\| + \delta L_0 + L_2 \|\bar{x}_t^\leftarrow - \bar{x}_{k\delta}^\leftarrow\| \\ & \leq \epsilon_{\text{approx}} + L_2 \|x_t^\leftarrow - \bar{x}_t^\leftarrow\| + \delta L_0 + \delta L_2 L_1. \end{aligned}$$

534 Applying Gronwall's Lemma over the interval $t \in [-t_{\max}, -t_{\min}]$,

$$\begin{aligned} & \|x_{t_{\min}} - \bar{x}_{t_{\min}}\| \\ & = \|x_{-t_{\min}}^\leftarrow - \bar{x}_{-t_{\min}}^\leftarrow\| \\ & \leq e^{L_2(t_{\max} - t_{\min})} (\|x_{-t_{\max}}^\leftarrow - \bar{x}_{-t_{\max}}^\leftarrow\| + (\epsilon_{\text{approx}} + \delta L_0 + \delta L_2 L_1)(t_{\max} - t_{\min})) \\ & = e^{L_2(t_{\max} - t_{\min})} (\|x_{t_{\max}} - \bar{x}_{t_{\max}}\| + (\epsilon_{\text{approx}} + \delta L_0 + \delta L_2 L_1)(t_{\max} - t_{\min})). \end{aligned}$$

535 □

536 **Lemma 3.** *Given initial $x_{t_{\max}}^0$, let $x_{t_{\min}}^0 = \text{ODE}(x_{t_{\max}}^0, t_{\max} \rightarrow t_{\min})$, and let $\hat{x}_{t_{\min}}^0 =$*
 537 *$\text{ODE}_\theta(x_{t_{\max}}^0, t_{\max} \rightarrow t_{\min})$. We further denote the variables after K Restart iterations as $x_{t_{\min}}^K =$*
 538 *$\text{Restart}(x_{t_{\min}}^0, K)$ and $\hat{x}_{t_{\min}}^K = \text{Restart}_\theta(\hat{x}_{t_{\min}}^0, K)$, with true field and learned field respectively. Then*
 539 *there exists a coupling between $x_{t_{\min}}^K$ and $\hat{x}_{t_{\min}}^K$ such that*

$$\|x_{t_{\min}}^K - \hat{x}_{t_{\min}}^K\| \leq e^{(K+1)L_2(t_{\max} - t_{\min})} (K+1) (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min}).$$

540 *Proof.* We will couple $x_{t_{\min}}^i$ and $\hat{x}_{t_{\min}}^i$ by using the same noise $\varepsilon_{t_{\min} \rightarrow t_{\max}}^i$ in the Restart forward process
 541 for $i = 0 \dots K-1$ (see Eq. (7)). For any i , let us also define $y_{t_{\min}}^{i,j} := \text{Restart}_\theta(x_{t_{\min}}^i, j-i)$, and this
 542 process uses the same noise $\varepsilon_{t_{\min} \rightarrow t_{\max}}^i$ as previous ones. From this definition, $y_{t_{\min}}^{K,K} = x_{t_{\min}}^K$. We can
 543 thus bound

$$\|x_{t_{\min}}^K, \hat{x}_{t_{\min}}^K\| \leq \|y_{t_{\min}}^{0,K} - \hat{x}_{t_{\min}}^K\| + \sum_{i=0}^{K-1} \|y_{t_{\min}}^{i,K} - y_{t_{\min}}^{i+1,K}\| \quad (16)$$

544 Using the assumption that $ts_\theta(\cdot, t)$ is L_2 Lipschitz,

$$\begin{aligned} & \left\| y_{t_{\min}}^{0,i+1} - \hat{x}_{t_{\min}}^{i+1} \right\| \\ &= \left\| \text{ODE}_\theta(y_{t_{\max}}^{0,i}, t_{\max} \rightarrow t_{\min}) - \text{ODE}_\theta(\hat{x}_{t_{\max}}^i, t_{\max} \rightarrow t_{\min}) \right\| \\ &\leq e^{L_2(t_{\max}-t_{\min})} \left\| y_{t_{\max}}^{0,i} - \hat{x}_{t_{\max}}^i \right\| \\ &= e^{L_2(t_{\max}-t_{\min})} \left\| y_{t_{\min}}^{0,i} - \hat{x}_{t_{\min}}^i \right\|, \end{aligned}$$

545 where the last equality is because we add the same additive Gaussian noise $\varepsilon_{t_{\min} \rightarrow t_{\max}}^i$ to $y_{t_{\min}}^{0,i}$ and $\hat{x}_{t_{\min}}^i$
546 in the Restart forward process. Applying the above recursively, we get

$$\begin{aligned} \left\| y_{t_{\min}}^{0,K} - \hat{x}_{t_{\min}}^K \right\| &\leq e^{KL_2(t_{\max}-t_{\min})} \left\| y_{t_{\min}}^{0,0} - \hat{x}_{t_{\min}}^0 \right\| \\ &\leq e^{KL_2(t_{\max}-t_{\min})} \left\| x_{t_{\min}}^0 - \hat{x}_{t_{\min}}^0 \right\| \\ &\leq e^{(K+1)L_2(t_{\max}-t_{\min})} (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min}), \end{aligned} \quad (17)$$

547 where the last line follows by Lemma 2 when setting $x_{t_{\max}} = \bar{x}_{t_{\max}}$. We will now bound

548 $\left\| y_{t_{\min}}^{i,K} - y_{t_{\min}}^{i+1,K} \right\|$ for some $i \leq K$. It follows from definition that

$$\begin{aligned} y_{t_{\min}}^{i,i+1} &= \text{ODE}_\theta(x_{t_{\max}}^i, t_{\max} \rightarrow t_{\min}) \\ y_{t_{\min}}^{i+1,i+1} &= x_{t_{\min}}^{i+1} = \text{ODE}(x_{t_{\max}}^i, t_{\max} \rightarrow t_{\min}). \end{aligned}$$

549 By Lemma 2,

$$\left\| y_{t_{\min}}^{i,i+1} - y_{t_{\min}}^{i+1,i+1} \right\| \leq e^{L_2(t_{\max}-t_{\min})} (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min})$$

550 For the remaining steps from $i+2 \dots K$, both $y^{i,\cdot}$ and $y^{i+1,\cdot}$ evolve with ODE_θ in each step. Again
551 using the assumption that $ts_\theta(\cdot, t)$ is L_2 Lipschitz,

$$\left\| y_{t_{\min}}^{i,K} - y_{t_{\min}}^{i+1,K} \right\| \leq e^{(K-i)L_2(t_{\max}-t_{\min})} (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min})$$

552 Summing the above for $i = 0 \dots K-1$, and combining with Eq. (16) and Eq. (17) gives

$$\left\| x_{t_{\min}}^K - \hat{x}_{t_{\min}}^K \right\| \leq e^{(K+1)L_2(t_{\max}-t_{\min})} (K+1) (\delta(L_2L_1 + L_0) + \epsilon_{\text{approx}})(t_{\max} - t_{\min}).$$

553 □

554 **Lemma 4.** Consider the same setup as Theorem 3. Let $x_{t_{\min}} = \text{SDE}(x_{t_{\max}}, t_{\max} \rightarrow t_{\min})$ and let
555 $\bar{x}_{t_{\min}} = \text{SDE}(\bar{x}_{t_{\max}}, t_{\max} \rightarrow t_{\min})$. Then there exists a coupling between x_t and \bar{x}_t such that

$$\begin{aligned} \mathbb{E}[\|x_{t_{\min}} - \bar{x}_{t_{\min}}\|] &\leq e^{2L_2(t_{\max}-t_{\min})} \mathbb{E}[\|x_{t_{\max}} - \bar{x}_{t_{\max}}\|] \\ &\quad + e^{2L_2(t_{\max}-t_{\min})} \left(\epsilon_{\text{approx}} + \delta L_0 + L_2 \left(\delta L_1 + \sqrt{2\delta dt_{\max}} \right) \right) (t_{\max} - t_{\min}) \end{aligned}$$

556 *Proof.* Consider some fixed arbitrary k , and recall that δ is the stepsize. By definition of SDE and
557 SDE_θ , for $t \in ((k-1)\delta, k\delta]$,

$$\begin{aligned} dx_t &= -2t \nabla \log p_t(x_t) dt + \sqrt{2t} dB_t \\ d\bar{x}_t &= -2ts_\theta(\bar{x}_{k\delta}, k\delta) dt + \sqrt{2t} dB_t. \end{aligned}$$

558 Let us define a coupling between x_t and \bar{x}_t by identifying their respective Brownian motions. It
559 will be convenient to define the time-reversed processes $x_t^\leftarrow := x_{-t}$, and $\bar{x}_t^\leftarrow := \bar{x}_{-t}$, along with
560 $v(x, t) := \nabla \log p_{-t}(x)$ and $\bar{v}(x, t) := s_\theta(x, -t)$. Then there exists a Brownian motion B_t^\leftarrow , such
561 that for $t \in [-t_{\max}, -t_{\min}]$,

$$\begin{aligned} dx_t^\leftarrow &= -2tv(x_t^\leftarrow, t) dt + \sqrt{-2t} dB_t^\leftarrow \\ d\bar{x}_t^\leftarrow &= -2t\bar{v}(\bar{x}_{k\delta}^\leftarrow, k\delta) dt + \sqrt{-2t} dB_t^\leftarrow \\ \Rightarrow d(x_t^\leftarrow - \bar{x}_t^\leftarrow) &= -2t(v(x_t^\leftarrow, t) - \bar{v}(\bar{x}_{k\delta}^\leftarrow, k\delta)) dt, \end{aligned}$$

562 where k is the unique negative integer such that $t \in [k\delta, (k+1)\delta)$. Thus

$$\begin{aligned}
& \frac{d}{dt} \mathbb{E} [\|x_t^{\leftarrow} - \bar{x}_t^{\leftarrow}\|] \\
& \leq 2 (\mathbb{E} [\|tv(x_t^{\leftarrow}, t) - t\bar{v}(x_t^{\leftarrow}, t)\|] + \mathbb{E} [\|t\bar{v}(x_t^{\leftarrow}, t) - t\bar{v}(\bar{x}_t^{\leftarrow}, t)\|]) \\
& \quad + 2 (\mathbb{E} [\|t\bar{v}(\bar{x}_t^{\leftarrow}, t) - t\bar{v}(\bar{x}_t^{\leftarrow}, k\delta)\|] + \mathbb{E} [\|t\bar{v}(\bar{x}_t^{\leftarrow}, k\delta) - t\bar{v}(\bar{x}_{k\delta}^{\leftarrow}, k\delta)\|]) \\
& \leq 2 (\epsilon_{approx} + L_2 \mathbb{E} [\|x_t^{\leftarrow} - \bar{x}_t^{\leftarrow}\|] + \delta L_0 + L_2 \mathbb{E} [\|\bar{x}_t^{\leftarrow} - \bar{x}_{k\delta}^{\leftarrow}\|]) \\
& \leq 2 \left(\epsilon_{approx} + L_2 \mathbb{E} [\|x_t^{\leftarrow} - \bar{x}_t^{\leftarrow}\|] + \delta L_0 + L_2 \left(\delta L_1 + \sqrt{2\delta dt_{\max}} \right) \right).
\end{aligned}$$

563 By Gronwall's Lemma,

$$\begin{aligned}
& \mathbb{E} [\|x_{t_{\min}} - \bar{x}_{t_{\min}}\|] \\
& = \mathbb{E} [\|x_{-t_{\min}}^{\leftarrow} - \bar{x}_{-t_{\min}}^{\leftarrow}\|] \\
& \leq e^{2L_2(t_{\max} - t_{\min})} \left(\mathbb{E} [\|x_{-t_{\max}}^{\leftarrow} - \bar{x}_{-t_{\max}}^{\leftarrow}\|] + \left(\epsilon_{approx} + \delta L_0 + L_2 \left(\delta L_1 + \sqrt{2\delta dt_{\max}} \right) \right) (t_{\max} - t_{\min}) \right) \\
& = e^{2L_2(t_{\max} - t_{\min})} \left(\mathbb{E} [\|x_{t_{\max}} - \bar{x}_{t_{\max}}\|] + \left(\epsilon_{approx} + \delta L_0 + L_2 \left(\delta L_1 + \sqrt{2\delta dt_{\max}} \right) \right) (t_{\max} - t_{\min}) \right)
\end{aligned}$$

564 □

565 A.4 Mixing Bounds

566 **Lemma 5.** Consider the same setup as Theorem 3. Assume that $\delta \leq t_{\min}$. Let

$$\begin{aligned}
x_{t_{\min}} &= SDE_{\theta}(x_{t_{\max}}, t_{\max} \rightarrow t_{\min}) \\
y_{t_{\min}} &= SDE_{\theta}(y_{t_{\max}}, t_{\max} \rightarrow t_{\min}).
\end{aligned}$$

567 Then there exists a coupling between x_s and y_s such that

$$TV(x_{t_{\min}}, y_{t_{\min}}) \leq \left(1 - 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right) \cdot e^{-BL_1/t_{\min} - L_1^2 t_{\max}^2 / t_{\min}^2} \right) TV(x_{t_{\max}}, y_{t_{\max}})$$

568 *Proof.* We will construct a coupling between x_t and y_t . First, let $(x_{t_{\max}}, y_{t_{\max}})$ be sampled from the
569 optimal TV coupling, i.e., $Pr(x_{t_{\max}} \neq y_{t_{\max}}) = \frac{1}{2}TV(x_{t_{\max}}, y_{t_{\max}})$. Recall that by definition of SDE_{θ} ,
570 for $t \in ((k-1)\delta, k\delta]$,

$$dx_t = -2ts_{\theta}(x_{k\delta}, k\delta)dt + \sqrt{2t}dB_t.$$

571 Let us define a time-rescaled version of x_t : $\bar{x}_t := x_{t^2}$. We verify that

$$d\bar{x}_t = -s_{\theta}(\bar{x}_{(k\delta)^2}, k\delta)dt + dB_t,$$

572 where k is the unique integer satisfying $t \in [((k-1)\delta)^2, k^2\delta^2)$. Next, we define the time-reversed
573 process $\bar{x}_t^{\leftarrow} := \bar{x}_{-t}$, and let $v(x, t) := s_{\theta}(x, -t)$. We verify that there exists a Brownian motion B_t^x
574 such that, for $t \in [-t_{\max}^2, -t_{\min}^2]$,

$$d\bar{x}_t^{\leftarrow} = v_t^x dt + dB_t^x,$$

575 where $v_t^x = s_{\theta}(\bar{x}_{-(k\delta)^2}^{\leftarrow}, -k\delta)$, where k is the unique positive integer satisfying $-t \in (((k-1)\delta)^2, (k\delta)^2]$. Let $d\bar{y}_t^{\leftarrow} = v_t^y dt + dB_t^y$, be defined analogously. For any positive integer k and for
576 any $t \in [-(k\delta)^2, -((k-1)\delta)^2]$, let us define
577

$$z_t = \bar{x}_{-k^2\delta^2}^{\leftarrow} - \bar{y}_{-k^2\delta^2}^{\leftarrow} + (2k-1)\delta^2 \left(v_{-(k\delta)^2}^x - v_{-(k\delta)^2}^y \right) + \left(B_t^x - B_{-(k\delta)^2}^x \right) - \left(B_t^y - B_{-(k\delta)^2}^y \right).$$

578 Let $\gamma_t := \frac{z_t}{\|z_t\|}$. We will now define a coupling between dB_t^x and dB_t^y as

$$dB_t^y = (I - 2\mathbb{1}\{t \leq \tau\}\gamma_t\gamma_t^T) dB_t^x,$$

where $\mathbb{1}\{\cdot\}$ denotes the indicator function, i.e. $\mathbb{1}\{t \leq \tau\} = 1$ if $t \leq \tau$, and τ is a stopping time given by the first hitting time of $z_t = 0$. Let $r_t := \|z_t\|$. Consider some $t \in (-i^2\delta^2, -(i-1)^2\delta^2)$, and Let $j := \frac{t_{\max}}{\delta}$ (assume w.l.o.g that this is an integer), then

$$\begin{aligned} r_t - r_{-t_{\max}^2} &\leq \sum_{k=i}^j (2k-1)\delta^2 \left\| (v_{-(k\delta)^2}^x - v_{-(k\delta)^2}^y) \right\| + \int_{-t_{\max}^2}^t \mathbb{1}\{t \leq \tau\} 2dB_s^1 \\ &\leq \sum_{k=i}^j (k^2 - (k-1)^2) \delta^2 2L_1 / (t_{\min}) + \int_{-t_{\max}^2}^t \mathbb{1}\{t \leq \tau\} 2dB_t^1 \\ &= \int_{-t_{\max}^2}^{-(i-1)\delta^2} \frac{2L_1}{t_{\min}} ds + \int_{-t_{\max}^2}^t \mathbb{1}\{t \leq \tau\} 2dB_s^1, \end{aligned}$$

where $dB_s^1 = \langle \gamma_t, dB_s^x - dB_s^y \rangle$ is a 1-dimensional Brownian motion. We also verify that

$$\begin{aligned} r_{-t_{\max}^2} &= \|z_{-t_{\max}^2}\| \\ &= \left\| \bar{x}_{-t_{\max}^2}^{\leftarrow} - \bar{y}_{-t_{\max}^2}^{\leftarrow} + (2j-1)\delta^2 (v_{-t_{\max}^2}^x - v_{-t_{\max}^2}^y) + (B_t^x - B_{-t_{\max}^2}^x) - (B_t^y - B_{-t_{\max}^2}^y) \right\| \\ &\leq \left\| \bar{x}_{-t_{\max}^2}^{\leftarrow} + (2j-1)\delta^2 v_{-t_{\max}^2}^x + (B_{-(j-1)^2\delta^2}^x - B_{-t_{\max}^2}^x) \right\| \\ &\quad + \left\| \bar{y}_{-t_{\max}^2}^{\leftarrow} + (2j-1)\delta^2 v_{-t_{\max}^2}^y + (B_{-(j-1)^2\delta^2}^y - B_t^y + B_t^x - B_{-t_{\max}^2}^y) \right\| \leq B \end{aligned}$$

where the third relation is by adding and subtracting $B_{-(j-1)^2\delta^2}^x - B_t^x$ and using triangle inequality.

The fourth relation is by noticing that $\bar{x}_{-t_{\max}^2}^{\leftarrow} + (2j-1)\delta^2 v_{-t_{\max}^2}^x + (B_{-(j-1)^2\delta^2}^x - B_{-t_{\max}^2}^x) = \bar{x}_{-(j-1)^2\delta^2}^{\leftarrow}$ and that $\bar{y}_{-t_{\max}^2}^{\leftarrow} + (2j-1)\delta^2 v_{-t_{\max}^2}^y + (B_{-(j-1)^2\delta^2}^y - B_t^y + B_t^x - B_{-t_{\max}^2}^y) \stackrel{d}{=} \bar{y}_{-(j-1)^2\delta^2}^{\leftarrow}$, and then using our assumption in the theorem statement that all processes are supported on a ball of radius $B/2$.

We now define a process s_t defined by $ds_t = 2L_1/t_{\min}dt + 2dB_t^1$, initialized at $s_{-t_{\max}^2} = B \geq r_{-t_{\max}^2}$. We can verify that, up to time τ , $r_t \leq s_t$ with probability 1. Let τ' denote the first-hitting time of s_t to 0, then $\tau \leq \tau'$ with probability 1. Thus

$$Pr(\tau \leq -t_{\min}^2) \geq Pr(\tau' \leq -t_{\min}^2) \geq 2Q \left(\frac{B}{2\sqrt{t_{\max}^2 - t_{\min}^2}} \right) \cdot e^{-BL_1/t_{\min} - L_1^2 t_{\max}^2/t_{\min}^2}$$

where we apply Lemma 6. The proof follows by noticing that, if $\tau \leq -t_{\min}^2$, then $x_{t_{\min}} = y_{t_{\min}}$. This is because if $\tau \in [-k^2\delta^2, -(k-1)^2\delta^2]$, then $\bar{x}_{-(k-1)^2\delta^2}^{\leftarrow} = \bar{y}_{-(k-1)^2\delta^2}^{\leftarrow}$, and thus $\bar{x}_t^{\leftarrow} = \bar{y}_t^{\leftarrow}$ for all $t \geq -(k-1)^2\delta^2$, in particular, at $t = -t_{\min}^2$.

594

□

Lemma 6. Consider the stochastic process

$$dr_t = dB_t^1 + cdt.$$

Assume that $r_0 \leq B/2$. Let τ denote the hitting time for $r_t = 0$. Then for any $T \in \mathbb{R}^+$,

$$Pr(\tau \leq T) \geq 2Q \left(\frac{B}{2\sqrt{T}} \right) \cdot e^{-ac - \frac{c^2 T}{2}},$$

where Q is the tail probability of a standard Gaussian defined in Definition 1.

Proof. We will use the following facts in our proof:

1. For $x \sim \mathcal{N}(0, \sigma^2)$, $Pr(x > r) = \frac{1}{2} \left(1 - \operatorname{erf} \left(\frac{r}{\sqrt{2}\sigma} \right) \right) = \frac{1}{2} \operatorname{erfc} \left(\frac{r}{\sqrt{2}\sigma} \right)$.
2. $\int_0^T \frac{a \exp\left(-\frac{a^2}{2t}\right)}{\sqrt{2\pi t^3}} dt = \operatorname{erfc} \left(\frac{a}{\sqrt{2T}} \right) = 2Pr(\mathcal{N}(0, T) > a) = 2Q \left(\frac{a}{\sqrt{T}} \right)$ by definition of Q .

601 Let $dr_t = dB_t^1 + cdt$, with $r_0 = a$. The density of the hitting time τ is given by

$$p(\tau = t) = f(a, c, t) = \frac{a \exp\left(-\frac{(a+ct)^2}{2t}\right)}{\sqrt{2\pi t^3}}. \quad (18)$$

602 (see e.g. [3]). From item 2 above,

$$\int_0^T f(a, 0, t) dt = 2Q\left(\frac{a}{\sqrt{T}}\right).$$

603 In the case of a general $c \neq 0$, we can bound $\frac{(a+ct)^2}{2t} = \frac{a^2}{2t} + ac + \frac{c^2 t}{2}$. Consequently,

$$f(a, c, t) \geq f(a, 0, t) \cdot e^{-ac - \frac{c^2 t}{2}}.$$

604 Therefore,

$$\Pr(\tau \leq T) = \int_0^T f(a, c, t) dt \geq \int_0^T f(a, 0, t) dt e^{-c} = 2Q\left(\frac{B}{2\sqrt{T}}\right) \cdot e^{-ac - \frac{c^2 T}{2}}.$$

605 □

606 A.5 TV Overlap

607 **Definition 1.** Let x be sampled from standard normal distribution $\mathcal{N}(0, 1)$. We define the Gaussian
608 tail probability $Q(a) := \Pr(x \geq a)$.

609 **Lemma 7.** We verify that for any two random vectors $\xi_x \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and $\xi_y \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, each
610 belonging to \mathbb{R}^d , the total variation distance between $x' = x + \xi_x$ and $y' = y + \xi_y$ is given by

$$TV(x', y') = 1 - 2Q(r) \leq 1 - \frac{2r}{r^2 + 1} \frac{1}{\sqrt{2\pi}} e^{-r^2/2},$$

611 where $r = \frac{\|x-y\|}{2\sigma}$, and $Q(r) = \Pr(\xi \geq r)$, when $\xi \sim \mathcal{N}(0, 1)$.

612 *Proof.* Let $\gamma := \frac{x-y}{\|x-y\|}$. We decompose x', y' into the subspace/orthogonal space defined by γ :

$$\begin{aligned} x' &= x^\perp + \xi_x^\perp + x^\parallel + \xi_x^\parallel \\ y' &= y^\perp + \xi_y^\perp + y^\parallel + \xi_y^\parallel \end{aligned}$$

613 where we define

$$\begin{aligned} x^\parallel &:= \gamma \gamma^T x & x^\perp &:= x - x^\parallel \\ y^\parallel &:= \gamma \gamma^T y & y^\perp &:= y - y^\parallel \\ \xi_x^\parallel &:= \gamma \gamma^T \xi_x & \xi_x^\perp &:= \xi_x - \xi_x^\parallel \\ \xi_y^\parallel &:= \gamma \gamma^T \xi_y & \xi_y^\perp &:= \xi_y - \xi_y^\parallel \end{aligned}$$

614 We verify the independence $\xi_x^\perp \perp \xi_x^\parallel$ and $\xi_y^\perp \perp \xi_y^\parallel$ as they are orthogonal decompositions of the
615 standard Gaussian. We will define a coupling between x' and y' by setting $\xi_x^\perp = \xi_y^\perp$. Under this
616 coupling, we verify that

$$(x^\perp + \xi_x^\perp) - (y^\perp + \xi_y^\perp) = x - y - \gamma \gamma^T (x - y) = 0$$

617 Therefore, $x' = y'$ if and only if $x^\parallel + \xi_x^\parallel = y^\parallel + \xi_y^\parallel$. Next, we draw (a, b) from the optimal coupling
618 between $\mathcal{N}(0, 1)$ and $\mathcal{N}(\frac{\|x-y\|}{\sigma}, 1)$. We verify that $x^\parallel + \xi_x^\parallel$ and $y^\parallel + \xi_y^\parallel$ both lie in the span of
619 γ . Thus it suffices to compare $\langle \gamma, x^\parallel + \xi_x^\parallel \rangle$ and $\langle \gamma, y^\parallel + \xi_y^\parallel \rangle$. We verify that $\langle \gamma, x^\parallel + \xi_x^\parallel \rangle =$

620 $\langle \gamma, y^\parallel \rangle + \langle \gamma, x^\parallel - y^\parallel \rangle + \langle \gamma, \xi_x^\parallel \rangle \sim \mathcal{N}(\langle \gamma, y^\parallel \rangle + \|x - y\|, \sigma^2) \stackrel{d}{=} \langle \gamma, y^\parallel \rangle + \sigma b$. We similarly verify
 621 that $\langle \gamma, y^\parallel + \xi_y^\parallel \rangle = \langle \gamma, y^\parallel \rangle + \langle \gamma, \xi_y^\parallel \rangle \sim \mathcal{N}(\langle \gamma, y^\parallel \rangle, \sigma^2) \stackrel{d}{=} \langle \gamma, y^\parallel \rangle + \sigma a$.
 622 Thus $TV(x', y') = TV(\sigma a, \sigma b) = 1 - 2Q\left(\frac{\|x - y\|}{2\sigma}\right)$. The last inequality follows from

$$Pr(\mathcal{N}(0, 1) \geq r) \geq \frac{r}{r^2 + 1} \frac{1}{\sqrt{2\pi}} e^{-r^2/2}$$

623

□

624 B More on Restart Algorithm

625 B.1 EDM Discretization Scheme

626 [13] proposes a discretization scheme for ODE given the starting t_{\max} and end time t_{\min} . Denote the
 627 number of steps as N , then the EDM discretization scheme is:

$$t_{i < N} = \left(t_{\max}^\rho + \frac{i}{N-1} (t_{\min}^\rho - t_{\max}^\rho) \right)^\rho$$

628 with $t_0 = t_{\max}$ and $t_{N-1} = t_{\min}$. ρ is a hyperparameter that determines the extent to which steps near
 629 t_{\min} are shortened. We adopt the value $\rho = 7$ suggested by [13] in all of our experiments. We apply
 630 the EDM scheme to creates a time discretization in each Restart interval $[t_{\max}, t_{\min}]$ in the Restart
 631 backward process, as well as the main backward process between $[0, T]$ (by additionally setting
 632 $t_{\min} = 0.002$ and $t_N = 0$ as in [13]). It is important to note that t_{\min} should be included within the
 633 list of time steps in the main backward process to seamlessly incorporate the Restart interval into the
 634 main backward process. We summarize the scheme as a function in Algorithm 1.

Algorithm 1 EDM_Scheme($t_{\min}, t_{\max}, N, \rho = 7$)

1: **return** $\left\{ \left(t_{\max}^\rho + \frac{i}{N-1} (t_{\min}^\rho - t_{\max}^\rho) \right)^\rho \right\}_{i=0}^{N-1}$

635 B.2 Restart Algorithm

636 We present the pseudocode for the Restart algorithm in Algorithm 2. In this pseudocode, we describe
 637 a more general case that applies l -level Restarting strategy. For each Restart segment, the include
 638 the number of steps in the Restart backward process N_{Restart} , the Restart interval $[t_{\min}, t_{\max}]$ and the
 639 number of Restart iteration K . We further denote the number of steps in the main backward process
 640 as N_{main} . We use the EDM discretization scheme (Algorithm 1) to construct time steps for the main
 641 backward process ($t_0 = T, t_{N_{\text{main}}} = 0$) as well as the Restart backward process, when given the
 642 starting/end time and the number of steps.

643 Although Heun's 2nd order method [2] (Algorithm 3) is the default ODE solver in the pseudocode, it
 644 can be substituted with other ODE solvers, such as Euler's method or the DPM solver [16].

645 The provided pseudocode in Algorithm 2 is tailored specifically for diffusion models [13]. To
 646 adapt Restart for other generative models like PFGM++ [28], we only need to modify the Gaussian
 647 perturbation kernel in the Restart forward process (line 10 in Algorithm 2) to the one used in
 648 PFGM++.

649 C Experimental Details

650 In this section, we discuss the configurations for different samplers in details. All the experiments are
 651 conducted on eight NVIDIA A100 GPUs.

652 C.1 Configurations for Baselines

653 We select **Vanilla SDE** [23], **Improved SDE** [13], **Gonna Go Fast** [12] as SDE baselines and
 654 the **Heun's** 2nd order method [2] (Alg 3) as ODE baseline on standard benchmarks CIFAR-10 and

Algorithm 2 Restart sampling

```
1: Input: Score network  $s_\theta$ , time steps in main backward process  $t_{i \in \{0, N_{\text{main}}\}}$ , Restart parameters  
    $\{(N_{\text{Restart},j}, K_j, t_{\min,j}, t_{\max,j})\}_{j=1}^l$   
2: Round  $t_{\min,j \in \{1,l\}}$  to its nearest neighbor in  $t_{i \in \{0, N_{\text{main}}\}}$   
3: Sample  $x_0 \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$   
4: for  $i = 0 \dots N_{\text{main}} - 1$  do ▷ Main backward process  
5:    $x_{t_{i+1}} = \text{OneStep\_Heun}(s_\theta, t_i, t_{i+1})$  ▷ Running single step ODE  
6:   if  $\exists j \in \{1, \dots, l\}, t_{i+1} = t_{\min,j}$  then  
7:      $t_{\min} = t_{\min,j}, t_{\max} = t_{\max,j}, K = K_j, N_{\text{Restart}} = N_{\text{Restart},j}$   
8:      $x_{t_{\min}}^0 = x_{t_{i+1}}$   
9:     for  $k = 0 \dots K - 1$  do ▷ Restart for  $K$  iterations  
10:       $\varepsilon_{t_{\min} \rightarrow t_{\max}} \sim \mathcal{N}(\mathbf{0}, (t_{\max}^2 - t_{\min}^2) \mathbf{I})$   
11:       $x_{t_{\max}}^{k+1} = x_{t_{\min}}^k + \varepsilon_{t_{\min} \rightarrow t_{\max}}$  ▷ Restart forward process  
12:       $\{\bar{t}_m\}_{m=0}^{N_{\text{Restart}}-1} = \text{EDM\_Scheme}(t_{\min}, t_{\max}, N_{\text{Restart}})$   
13:      for  $m = 0 \dots N_{\text{Restart}} - 1$  do ▷ Restart backward process  
14:         $x_{\bar{t}_{m+1}}^{k+1} = \text{OneStep\_Heun}(s_\theta, \bar{t}_m, \bar{t}_{m+1})$   
15:      end for  
16:    end for  
17:  end if  
18: end for  
19: return  $x_{t_{N_{\text{main}}}}$ 
```

Algorithm 3 OneStep_Heun($s_\theta, x_{t_i}, t_i, t_{i+1}$)

```
1:  $d_i = t_i s_\theta(x_{t_i}, t_i)$   
2:  $x_{t_{i+1}} = x_{t_i} - (t_{i+1} - t_i) d_i$   
3: if  $t_{i+1} \neq 0$  then  
4:    $d'_i = t_{i+1} s_\theta(x_{t_{i+1}}, t_{i+1})$   
5:    $x_{t_{i+1}} = x_{t_i} - (t_{i+1} - t_i) (\frac{1}{2} d_i + \frac{1}{2} d'_i)$   
6: end if  
7: return  $x_{t_{i+1}}$ 
```

655 ImageNet 64×64 . We choose **DDIM** [22], **Heun**’s 2nd order method, and **DDPM** [9] for comparison
656 on Stable Diffusion model.

657 Vanilla SDE denotes the reverse-time SDE sampler in [23]. For Improved SDE, we use the recom-
658 mended dataset-specific hyperparameters (*e.g.*, $S_{\max}, S_{\min}, S_{\text{churn}}$) in Table 5 of the EDM paper [13].
659 They obtained these hyperparameters by grid search. Gonna Go Fast [12] applied an adaptive step
660 size technique based on Vanilla SDE and we directly report the FID scores listed in [12] for Gonna
661 Go Fast on CIFAR-10 (VP). For fair comparison, we use the EDM discretization scheme [13] for
662 Vanilla SDE, Improved SDE, Heun as well as Restart.

663 We borrow the hyperparameters such as discretization scheme or initial noise scale on Stable Diffusion
664 models in the diffuser² code repository. We directly use the DDIM and DDPM samplers implemented
665 in the repo. We apply the same set of hyperparameters to Heun and Restart.

666 C.2 Configurations for Restart

667 We report the configurations for Restart for different models and NFE on standard benchmarks
668 CIFAR-10 and ImageNet 64×64 . The hyperparameters of Restart include the number of steps
669 in the main backward process N_{main} , the number of steps in the Restart backward process N_{Restart} ,
670 the Restart interval $[t_{\min}, t_{\max}]$ and the number of Restart iteration K . In Table 3 (CIFAR-10, VP)
671 we provide the quintuplet $(N_{\text{main}}, N_{\text{Restart}}, t_{\min}, t_{\max}, K)$ for each experiment. Since we apply the
672 multi-level Restart strategy for ImageNet 64×64 , we provide N_{main} as well as a list of quadruple
673 $\{(N_{\text{Restart},i}, K_i, t_{\min,i}, t_{\max,i})\}_{i=1}^l$ (l is the number of Restart interval depending on experiments) in
674 Table 5. In order to integrate the Restart time interval to the main backward process, we round $t_{\min,i}$

²<https://github.com/huggingface/diffusers>

to its nearest neighbor in the time steps of main backward process, as shown in line 2 of Algorithm 2. We apply Heun method for both main/backward process. The formula for NFE calculation is
$$\text{NFE} = \underbrace{2 \cdot N_{\text{main}} - 1}_{\text{main backward process}} + \sum_{i=1}^l \underbrace{K_i}_{\text{number of repetitions}} \cdot \underbrace{(2 \cdot (N_{\text{Restart},i} - 1))}_{\text{per iteration in } i^{\text{th}} \text{ Restart interval}}$$
 in this case. Inspired by [13], we inflate the additive noise in the Restart forward process by multiplying $S_{\text{noise}} = 1.003$ on ImageNet 64×64 , to counteract the over-denoising tendency of neural networks. We also observe that setting $\gamma = 0.05$ in Algorithm 2 of EDM [13] would slightly boost the Restart performance on ImageNet 64×64 when $t \in [0.01, 1]$.

We further include the configurations for Restart on Stable Diffusion models in Table 10, with a varying guidance weight w . Similar to ImageNet 64×64 , we use multi-level Restart with a fixed number of steps $N_{\text{main}} = 30$ in the main backward process. We utilize the Euler method for the main backward process and the Heun method for the Restart backward process, as our empirical observations indicate that the Heun method doesn't yield significant improvements over the Euler method, yet necessitates double the steps. The number of steps equals to $N_{\text{main}} + \sum_{i=1}^l K_i \cdot (2 \cdot (N_{\text{Restart},i} - 1))$ in this case. We set the total number of steps to 66, including main backward process and Restart backward process.

Given the prohibitively large search space for each Restart quadruple, a comprehensive enumeration of all possibilities is impractical due to computational limitations. Instead, we adjust the configuration manually, guided by the heuristic that weaker/smaller models or more challenging tasks necessitate a stronger Restart strength (e.g., larger K , wider Restart interval, etc). On average, we select the best configuration from 5 sets for each experiment; these few trials have empirically outperformed previous SDE/ODE samplers. We believe that developing a systematic approach for determining Restart configurations could be of significant value in the future.

C.3 Pre-trained Models

For CIFAR-10 dataset, we use the pre-trained VP and EDM models from the EDM repository³, and PFGM++ ($D = 2048$) model from the PFGM++ repository⁴. For ImageNet 64×64 , we borrow the pre-trained EDM model from EDM repository as well.

C.4 Classifier-free Guidance

We follow the convention in [20], where each step in classifier-free guidance is as follows:

$$\tilde{s}_{\theta}(x, c, t) = w s_{\theta}(x, c, t) + (1 - w) s_{\theta}(x, t)$$

where c is the conditions, and $s_{\theta}(x, c, t)/s_{\theta}(x, t)$ is the conditional/unconditional models, sharing parameters. Increasing w would strengthen the effect of guidance, usually leading to a better text-image alignment [20].

C.5 More on the Synthetic Experiment

C.5.1 Discrete Dataset

We generate the underlying discrete dataset S with $|S| = 2000$ as follows. Firstly, we sample 2000 points, denoted as S_1 , from a mixture of two Gaussians in \mathbb{R}^4 . Next, we project these points onto \mathbb{R}^{20} . To ensure a variance of 1 on each dimension, we scale the coordinates accordingly. This setup aims to simulate data points that primarily reside on a lower-dimensional manifold with multiple modes.

The specific details are as follows: $S_1 \sim 0.3N(a, s^2 I) + 0.7(-a, s^2 I)$, where $a = (3, 3, 3, 3) \subset \mathbb{R}^4$ and $s = 1$. Then, we randomly select a projection matrix $P \in \mathbb{R}^{20 \times 4}$, where each entry is drawn from $N(0, 1)$, and compute $S_2 = PS_1$. Finally, we scale each coordinate by a constant factor to ensure a variance of 1.

³<https://github.com/NVlabs/edm>

⁴<https://github.com/Newbeeer/pfgmpp>

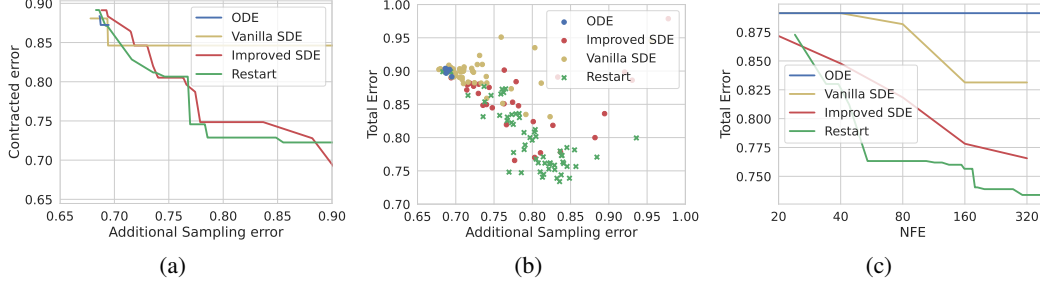


Figure 7: Comparison of additional sampling error versus (a) contracted error (plotting the Pareto frontier) and (b) total error (using a scatter plot). (c) Pareto frontier of NFE versus total error.

C.5.2 Model Architecture

We employ a common MLP architecture with a latent size of 64 to learn the score function. The training method is adapted from [13], which includes the preconditioning technique and denoising score-matching objective [25].

C.5.3 Varying Hyperparameters

To achieve the best trade-off between contracted error and additional sampling error, and optimize the NFE versus FID (Fréchet Inception Distance) performance, we explore various hyperparameters. [13] shows that the Vanilla SDE can be endowed with additional flexibility by varying the coefficient $\beta(t)$ (Eq.(6) in [13]). Hence, regarding SDE, we consider NFE values from $\{20, 40, 80, 160, 320\}$, and multiply the original $\beta(t) = \dot{\sigma}(t)/\sigma(t)$ [13] with values from $\{0, 0.25, 0.5, 1, 1.5, 2, 4, 8\}$. It is important to note that larger NFE values do not lead to further performance improvements. For restarts, we tried the following two settings: first we set the number of steps in Restart backward process to 40 and vary the number of Restart iterations K in the range $\{0, 5, 10, 15, 20, 25, 30, 35\}$. We also conduct a grid search with the number of Restart iterations K ranging from 5 to 25 and the number of steps in Restart backward process varying from 2 to 7. For ODE, we experiment with the number of steps set to $\{20, 40, 80, 160, 320, 640\}$.

Additionally, we conduct an experiment for Improved SDE in EDM. We try different values of S_{churn} in the range of $\{0, 1, 2, 4, 8, 16, 32, 48, 64\}$. We also perform a grid search where the number of steps ranged from 20 to 320 and S_{churn} takes values of $[0.2 \times \text{steps}, 0.5 \times \text{steps}, 20, 60]$. The plot combines the results from SDE and is displayed in Figure 7.

To mitigate the impact of randomness, we collect the data by averaging the results from five runs with the same hyperparameters. To compute the Wasserstein distance between two discrete distributions, we use minimum weight matching.

C.5.4 Plotting the Pareto frontier

We generate the Pareto frontier plots as follows. For the additional sampling error versus contracted error plot, we first sort all the data points based on their additional sampling error and then connect the data points that represent prefix minimums of the contracted error. Similarly, for the NFE versus FID plot, we sort the data points based on their NFE values and connect the points where the FID is a prefix minimum.

D Extra Experimental Results

D.1 Numerical Results

In this section, we provide the corresponding numerical results of Fig. 3(a) and Fig. 3(b), in Table 2, 3 (CIFAR-10 VP) and Table 4, 5 (ImageNet 64×64 EDM), respectively. We also include the performance of Vanilla SDE in those tables. For the evaluation, we compute the Fréchet distance between 50000 generated samples and the pre-computed statistics of CIFAR-10 and ImageNet 64×64 . We follow the evaluation protocol in EDM [13] that calculates each FID scores three times with different seeds and report the minimum.

We also provide the numerical results on the Stable Diffusion model [19], with a classifier guidance weight $w = 2, 3, 5, 8$ in Table 6, 7, 8, 9. As in [17], we report the zero-shot FID score on 5K random prompts sampled from the COCO validation set. We evaluate CLIP score [6] with the open-sourced ViT-g/14 [11], Aesthetic score by the more recent LAION-Aesthetics Predictor V2⁵. We average the CLIP and Aesthetic scores over 5K generated samples. The number of function evaluations is two times the sampling steps in Stable Diffusion model, since each sampling step involves the evaluation of the conditional and unconditional model.

Table 2: CIFAR-10 sample quality (FID score) and number of function evaluations (NFE) on VP [23] for baselines

	NFE	FID
<i>ODE (Heun) [13]</i>	1023	2.90
	511	2.90
	255	2.90
	127	2.90
	63	2.89
	35	2.97
<i>Vanilla SDE [23]</i>	1024	2.79
	512	4.01
	256	4.79
	128	12.57
<i>Gonna Go Fast [12]</i>	1000	2.55
	329	2.70
	274	2.74
	179	2.59
	147	2.95
	49	72.29
<i>Improved SDE [13]</i>	1023	2.35
	511	2.37
	255	2.40
	127	2.58
	63	2.88
	35	3.45

Table 3: CIFAR-10 sample quality (FID score), number of function evaluations (NFE) and configurations on VP [23] for Restart

NFE	FID	Configuration
		$(N_{\text{main}}, N_{\text{Restart}, i}, K_i, t_{\text{min}, i}, t_{\text{max}, i})$
519	2.11	(20, 9, 30, 0.06, 0.20)
115	2.21	(18, 3, 20, 0.06, 0.30)
75	2.27	(18, 3, 10, 0.06, 0.30)
55	2.45	(18, 3, 5, 0.06, 0.30)
43	2.70	(18, 3, 2, 0.06, 0.30)

⁵<https://github.com/christophschuhmann/improved-aesthetic-predictor>

Table 4: ImageNet 64×64 sample quality (FID score) and number of function evaluations (NFE) on EDM [13] for baselines

	NFE	FID (50k)
<i>ODE (Heun)</i> [13]	1023	2.24
	511	2.24
	255	2.24
	127	2.25
	63	2.30
	35	2.46
<i>Vanilla SDE</i> [23]	1024	1.89
	512	3.38
	256	11.91
	128	59.71
<i>Improved SDE</i> [13]	1023	1.40
	511	1.45
	255	1.50
	127	1.75
	63	2.24
	35	2.97

Table 5: ImageNet 64×64 sample quality (FID score), number of function evaluations (NFE) and configurations on EDM [13] for Restart

NFE	FID (50k)	Configuration $N_{\text{main}}, \{(N_{\text{Restart},i}, K_i, t_{\min,i}, t_{\max,i})\}_{i=1}^l$
623	1.36	36, {(10, 3, 19.35, 40.79), (10, 3, 1.09, 1.92), (7, 6, 0.59, 1.09), (7, 6, 0.30, 0.59), (7, 25, 0.06, 0.30)}
535	1.39	36, {(6, 1, 19.35, 40.79), (6, 1, 1.09, 1.92), (7, 6, 0.59, 1.09), (7, 6, 0.30, 0.59), (7, 25, 0.06, 0.30)}
385	1.41	36, {(3, 1, 19.35, 40.79), (6, 1, 1.09, 1.92), (6, 5, 0.59, 1.09), (6, 5, 0.30, 0.59), (6, 20, 0.06, 0.30)}
203	1.46	36, {(4, 1, 19.35, 40.79), (4, 1, 1.09, 1.92), (4, 5, 0.59, 1.09), (4, 5, 0.30, 0.59), (6, 6, 0.06, 0.30)}
165	1.51	18, {(3, 1, 19.35, 40.79), (4, 1, 1.09, 1.92), (4, 5, 0.59, 1.09), (4, 5, 0.30, 0.59), (4, 10, 0.06, 0.30)}
99	1.71	18, {(3, 1, 19.35, 40.79), (4, 1, 1.09, 1.92), (4, 4, 0.59, 1.09), (4, 1, 0.30, 0.59), (4, 4, 0.06, 0.30)}
67	1.95	18, {(5, 1, 19.35, 40.79), (5, 1, 1.09, 1.92), (5, 1, 0.59, 1.09), (5, 1, 0.06, 0.30)}
39	2.38	14, {(3, 1, 19.35, 40.79), (3, 1, 1.09, 1.92), (3, 1, 0.06, 0.30)}

Table 6: Numerical results on Stable Diffusion v1.5 with a classifier-free guidance weight $w = 2$

	Steps	FID (5k) ↓	CLIP score ↑	Aesthetic score ↑
<i>DDIM</i> [22]	50	16.08	0.2905	5.13
	100	15.35	0.2920	5.15
<i>Heun</i>	51	18.80	0.2865	5.14
	101	18.21	0.2871	5.15
<i>DDPM</i> [9]	100	13.53	0.3012	5.20
	200	13.22	0.2999	5.19
<i>Restart</i>	66	13.16	0.2987	5.19

Table 7: Numerical results on Stable Diffusion v1.5 with a classifier-free guidance weight $w = 3$

	Steps	FID (5k) ↓	CLIP score ↑	Aesthetic score ↑
<i>DDIM</i> [22]	50	14.28	0.3056	5.22
	100	14.30	0.3056	5.22
<i>Heun</i>	51	15.63	0.3022	5.20
	101	15.40	0.3026	5.21
<i>DDPM</i> [9]	100	15.72	0.3129	5.28
	200	15.13	0.3131	5.28
<i>Restart</i>	66	14.48	0.3079	5.25

Table 8: Numerical results on Stable Diffusion v1.5 with a classifier-free guidance weight $w = 5$

	Steps	FID (5k) ↓	CLIP score ↑	Aesthetic score ↑
<i>DDIM</i> [22]	50	16.60	0.3154	5.31
	100	16.80	0.3157	5.31
<i>Heun</i>	51	16.26	0.3135	5.28
	101	16.38	0.3136	5.29
<i>DDPM</i> [9]	100	19.62	0.3197	5.36
	200	18.88	0.3200	5.35
<i>Restart</i>	66	16.21	0.3179	5.33

Table 9: Numerical results on Stable Diffusion v1.5 with a classifier-free guidance weight $w = 8$

	Steps	FID (5k) ↓	CLIP score ↑	Aesthetic score ↑
<i>DDIM</i> [22]	50	19.83	0.3206	5.37
	100	19.82	0.3200	5.37
<i>Heun</i>	51	18.44	0.3186	5.35
	101	18.72	0.3185	5.36
<i>DDPM</i> [9]	100	22.58	0.3223	5.39
	200	21.67	0.3212	5.38
<i>Restart</i>	47	18.40	0.3228	5.41

760 D.2 Study on Adjusting t_{\min}

761 We also investigate the impact of varying t_{\min} when $t_{\max} = t_{\min} + 0.3$. Fig. ?? reveals that FID scores
762 achieve a minimum at a t_{\min} close to 0 on VP, indicating higher accumulated errors at the end of

Table 10: Restart (Steps=66) configurations on Stable Diffusion v1.5

w	Configuration $N_{\text{main}}, \{(N_{\text{Restart},i}, K_i, t_{\min,i}, t_{\max,i})\}_{i=1}^l$
2	30, $\{(5, 2, 1, 9), (5, 2, 5, 10)\}$
3	30, $\{(2, 10, 0.1, 3)\}$
5	30, $\{(2, 10, 0.1, 2)\}$
8	30, $\{(2, 10, 0.1, 2)\}$

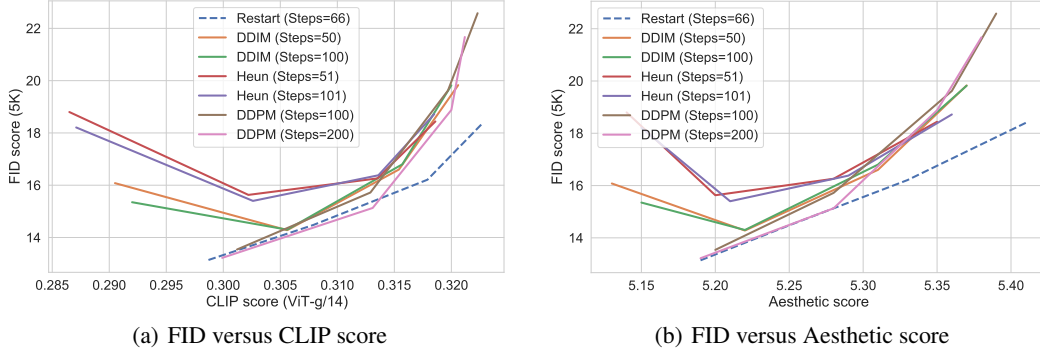


Figure 8: FID score versus (a) CLIP ViT-g/14 score and (b) Aesthetic score for text-to-image generation at 512×512 resolution, using Stable Diffusion v1.5 with varying classifier-free guidance weight $w = 2, 3, 5, 8$.

sampling and poor neural estimations at small t . Note that the Restart interval 0.3 is about twice the length of the one in Table 1 and Restart does not outperform the ODE baseline on EDM. This suggests that, as a rule of thumb, we should apply greater Restart strength (e.g., larger K , $t_{\max} - t_{\min}$) for weaker or smaller architectures and vice versa.

E Extended Generated Images

In this section, we provide extended generated images by Restart, DDIM, Heun and DDPM on text-to-image Stable Diffusion v1.5 model [19]. We showcase the samples of four sets of text prompts in Fig. 10, Fig. 11, Fig. 12, Fig. 13, with a classifier-guidance weight $w = 8$.

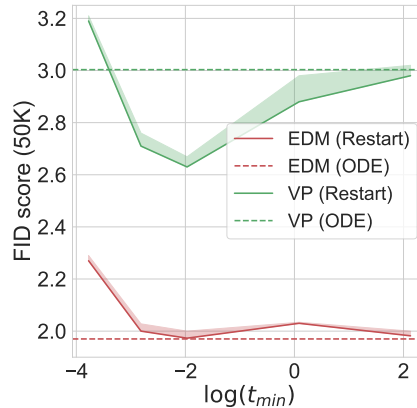


Figure 9: Adjusting t_{\min} in Restart on VP/EDM



(a) Restart (Steps=66)



(b) DDIM (Steps=100)



(c) Heun (Steps=101)



(d) DDPM (Steps=100)

Figure 10: Generated images with text prompt="A photo of an astronaut riding a horse on mars" and $w = 8$.



Figure 11: Generated images with text prompt="A raccoon playing table tennis" and $w = 8$.

F Broader Impact

The field of deep generative models incorporating differential equations is rapidly evolving and holds significant potential to shape our society. Nowadays, a multitude of photo-realistic images generated by text-to-image Stable Diffusion models populate the internet. Our work introduces Restart, a novel sampling algorithm that outperforms previous samplers for diffusion models and PFGM++. With applications extending across diverse areas, the Restart sampling algorithm is especially suitable for generation tasks demanding high quality and rapid speed. Yet, it is crucial to recognize that the utilization of such algorithms can yield both positive and negative repercussions, contingent on their specific applications. On the one hand, Restart sampling can facilitate the generation of highly realistic images and audio samples, potentially advancing sectors such as entertainment, advertising, and education. On the other hand, it could also be misused in *deepfake* technology, potentially leading to social scams and misinformation. In light of these potential risks, further research is required to develop robustness guarantees for generative models, ensuring their use aligns with ethical guidelines and societal interests.



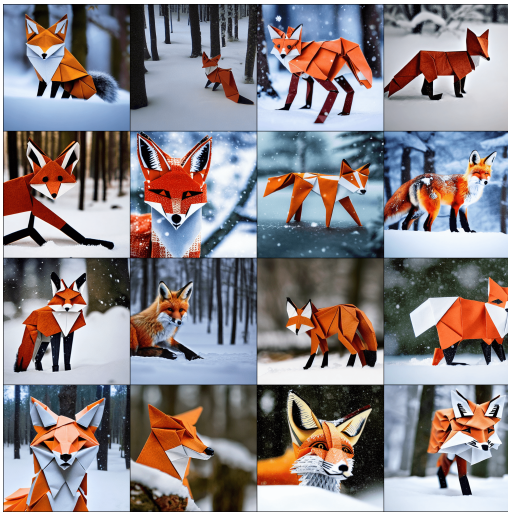
(a) Restart (Steps=66)



(b) DDIM (Steps=100)



(c) Heun (Steps=101)

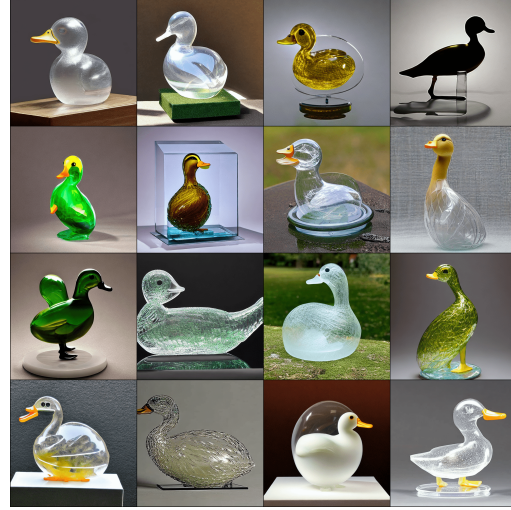


(d) DDPM (Steps=100)

Figure 12: Generated images with text prompt="Intricate origami of a fox in a snowy forest" and $w = 8$.



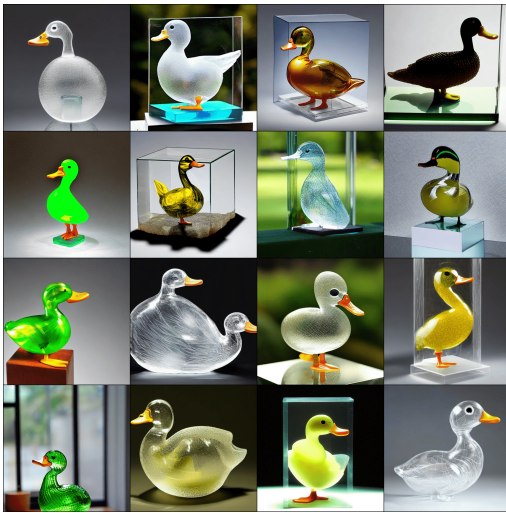
(a) Restart (Steps=66)



(b) DDIM (Steps=100)



(c) Heun (Steps=101)



(d) DDPM (Steps=100)

Figure 13: Generated images with text prompt="A transparent sculpture of a duck made out of glass" and $w = 8$.