

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 create 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, we 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>

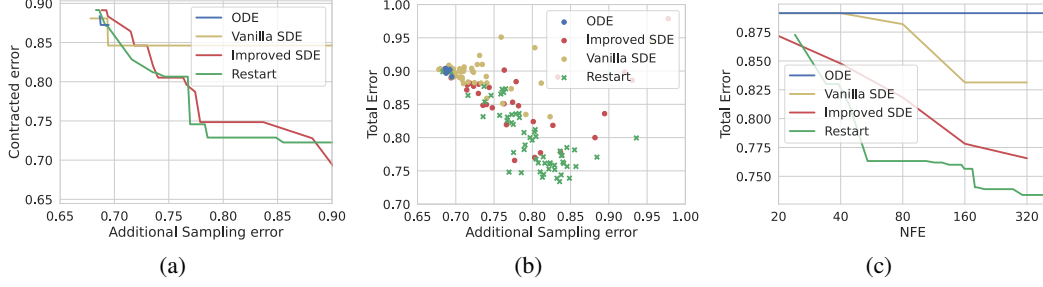


Figure 7: Comparison of additional sampling error versus (a) contracted error (plotting the Pareto frontier) and (b) total error (using a scatter plot). (c) Pareto frontier of NFE versus total error.

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, $\{(2, 10, 0.1, 2)\}$

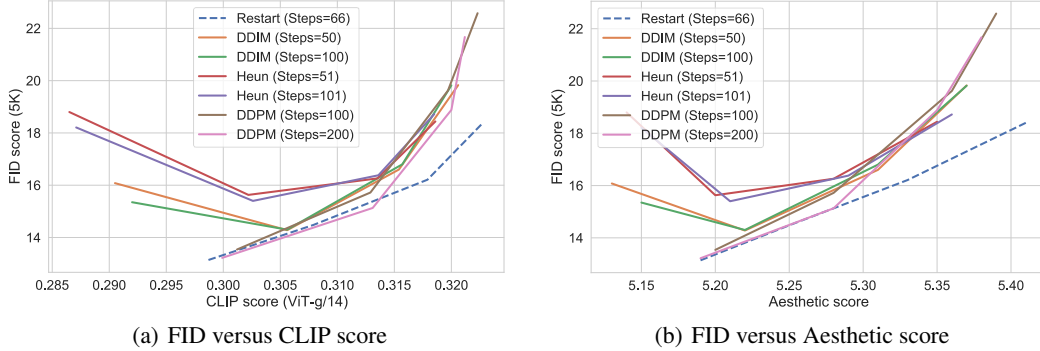


Figure 8: FID score versus (a) CLIP ViT-g/14 score and (b) Aesthetic score for text-to-image generation at 512×512 resolution, using Stable Diffusion v1.5 with varying classifier-free guidance weight $w = 2, 3, 5, 8$.

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 $w = 8$.

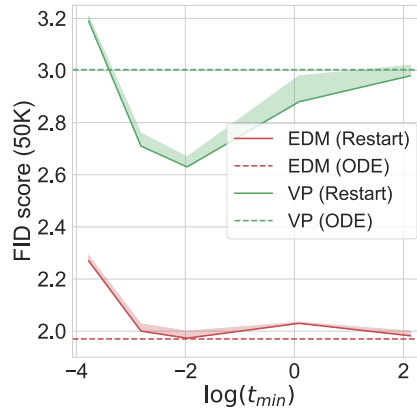


Figure 9: Adjusting t_{\min} in Restart on VP/EDM



(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.



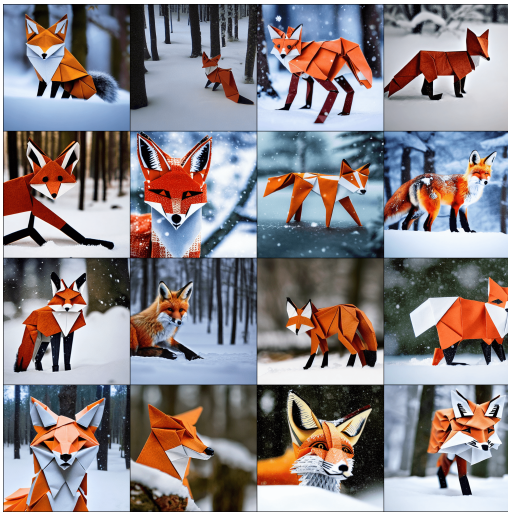
(a) Restart (Steps=66)



(b) DDIM (Steps=100)



(c) Heun (Steps=101)

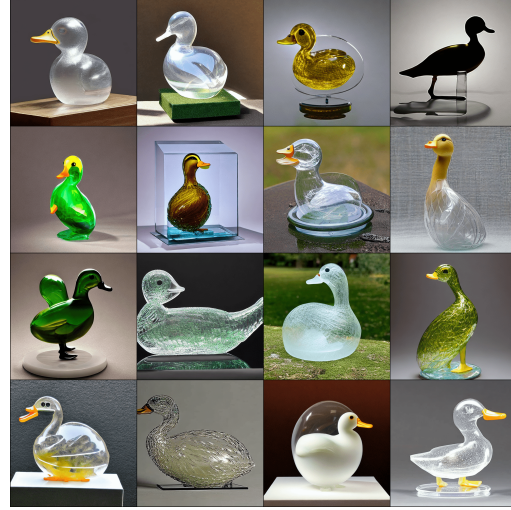


(d) DDPM (Steps=100)

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



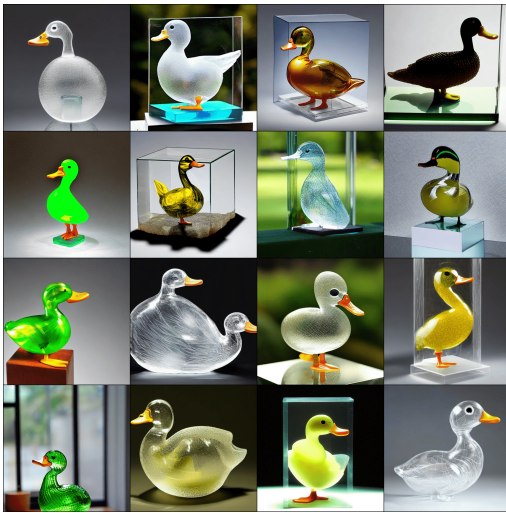
(a) Restart (Steps=66)



(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 $w = 8$.