

# Towards Faster Non-Asymptotic Convergence for Diffusion Generative Models

Gen Li\*      Yuting Wei\*      Yuxin Chen\*      Yuejie Chi<sup>†</sup>

May 23, 2023

## Abstract

Diffusion models, which convert noise into new data instances by learning to reverse a Markov diffusion process, have become a cornerstone in contemporary generative modeling. While their practical power has now been widely recognized, the theoretical underpinnings remain far from mature. In this work, we develop a suite of non-asymptotic theory towards understanding the data generation process of diffusion models, assuming access to reliable estimates of the (Stein) score functions. For a popular deterministic sampler (based on the probability flow ODE), we establish a convergence rate proportional to  $1/T$  (with  $T$  the total number of steps), improving upon past results; for another mainstream stochastic sampler (i.e., a type of the denoising diffusion probabilistic model (DDPM)), we derive a convergence rate proportional to  $1/\sqrt{T}$ , matching the state-of-the-art theory. Our theory imposes only minimal assumptions on the target data distribution, and is developed based on an elementary yet versatile non-asymptotic approach without resorting to toolboxes for SDEs and ODEs. Further, we design two accelerated variants, improving the convergence to  $1/T^2$  for the ODE-based sampler and  $1/T$  for the DDPM-type sampler, which might be of independent theoretical and empirical interest.

**Keywords:** diffusion models, score-based generative modeling, non-asymptotic theory, reverse SDE, probability flow ODE, denoising diffusion probabilistic model

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preliminaries</b>	<b>3</b>
2.1	Diffusion generative models . . . . .	4
2.2	Deterministic vs. stochastic samplers: a continuous-time interpretation . . . . .	4
<b>3</b>	<b>Algorithms and main results</b>	<b>6</b>
3.1	Assumptions and learning rates . . . . .	6
3.2	Deterministic samplers . . . . .	7
3.3	Stochastic samplers . . . . .	9
<b>4</b>	<b>Other related works</b>	<b>10</b>
<b>5</b>	<b>Discussion</b>	<b>11</b>

---

\*Department of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA.

<sup>†</sup>Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

<b>A Analysis for the deterministic samplers (Theorems 1 and 2)</b>	<b>15</b>
A.1 Preliminary facts . . . . .	15
A.2 Proof of Theorem 1 . . . . .	15
A.3 Proof of Theorem 2 . . . . .	18
A.4 Proof of properties (41) . . . . .	19
A.5 Proof of Lemma 1 . . . . .	20
A.6 Proof of Lemma 2 . . . . .	22
A.7 Proof of Lemma 3 . . . . .	26
<b>B Analysis for the stochastic samplers (Theorems 3 and 4)</b>	<b>28</b>
B.1 Step 1: approximation of $p_{X_{t-1} X_t}(x_{t-1} x_t)$ . . . . .	28
B.2 Step 2: uniform control of the density ratios . . . . .	29
B.3 Step 3: computing the KL divergence . . . . .	31
B.4 In summary . . . . .	32
B.5 Proof of Claim (94) . . . . .	32

# 1 Introduction

Diffusion models have emerged as a cornerstone in contemporary generative modeling, a task that learns to generate new data instances (e.g., images, text, audio) that look similar in distribution to the training data (Ho et al., 2020; Sohl-Dickstein et al., 2015; Song and Ermon, 2019; Dhariwal and Nichol, 2021; Jolicœur-Martineau et al., 2021; Chen et al., 2021; Kong et al., 2021; Austin et al., 2021). Originally proposed by Sohl-Dickstein et al. (2015) and later popularized by Song and Ermon (2019); Ho et al. (2020), the mainstream diffusion generative models — e.g., denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) and denoising diffusion implicit models (DDIMs) (Song et al., 2020a) — have underpinned major successes in content generators like DALL·E (Ramesh et al., 2022), Stable Diffusion (Rombach et al., 2022) and Imagen (Saharia et al., 2022), claiming state-of-the-art performance in the now broad field of generative artificial intelligence (AI). See Yang et al. (2022); Croitoru et al. (2023) for overviews of recent development.

In a nutshell, a diffusion generative model is based upon two stochastic processes in  $\mathbb{R}^d$ :

- 1) a forward process

$$X_0 \rightarrow X_1 \rightarrow \cdots \rightarrow X_T \quad \text{eq:forward-process-informal} \quad (1)$$

that starts from a sample drawn from the target data distribution (e.g., of natural images) and gradually diffuses it into a noise-like distribution (e.g., standard Gaussian);

- 2) a reverse process

$$Y_T \rightarrow Y_{T-1} \rightarrow \cdots \rightarrow Y_0 \quad \text{eq:reverse-process-informal} \quad (2)$$

that starts from pure noise (e.g., standard Gaussian) and successively converts it into new samples sharing similar distributions as the target data distribution.

Transforming data into noise in the forward process is straightforward, often hand-crafted by increasingly injecting more noise into the data at hand. What is challenging is the construction of the reverse process: how to generate the desired information out of pure noise? To do so, a diffusion model learns to build a reverse process (2) that imitates the dynamics of the forward process (1) in a time-reverse fashion; more precisely, the design goal is to ascertain distributional proximity<sup>1</sup>

$$Y_t \overset{d}{\approx} X_t, \quad t = T, \dots, 1 \quad (3)$$

through proper learning based on how the training data propagate in the forward process. Encouragingly, there often exist feasible strategies to achieve this goal as long as faithful estimates about the (Stein) score functions — the gradients of the log marginal density of the forward process — are available, an intriguing fact that can be illuminated by the existence and construction of reverse-time stochastic differential equations

---

<sup>1</sup>Two random vectors  $X$  and  $Y$  are said to obey  $X \overset{d}{=} Y$  (resp.  $X \overset{d}{\approx} Y$ ) if they are equivalent (resp. close) in distribution.

(SDEs) (Anderson, 1982) (see Section 2.2 for more precise discussions). Viewed in this light, a diverse array of diffusion models are frequently referred to as *score-based generative modeling* (SGM). The popularity of SGM was initially motivated by, and has since further inspired, numerous recent studies on the problem of learning score functions, a subroutine that also goes by the name of score matching (e.g., Hyvärinen (2005, 2007); Vincent (2011); Song et al. (2020b)).

Nonetheless, despite the mind-blowing empirical advances, a mathematical theory for diffusion generative models is still in its infancy. Given the complexity of developing a full-fledged end-to-end theory, a divide-and-conquer approach has been advertised, decoupling the score learning phase (i.e., how to estimate score functions reliably from training data) and the generative sampling phase (i.e., how to generate new data instances given the score estimates). In particular, the past two years have witnessed growing interest and remarkable progress from the theoretical community towards understanding the generative sampling phase (Block et al., 2020; De Bortoli et al., 2021; Liu et al., 2022; De Bortoli, 2022; Lee et al., 2023; Pidstrigach, 2022; Chen et al., 2022b,a, 2023b). For instance, polynomial-time convergence guarantees have been established for stochastic samplers (e.g., Chen et al. (2022b,a)) and deterministic samplers (e.g., Chen et al. (2023b)), both of which accommodated a fairly general family of data distributions.

**This paper.** The present paper contributes to this growing list of theoretical endeavors by developing a new suite of non-asymptotic theory for several score-based generative modeling algorithms. We concentrate on two types of samplers (Song et al., 2021b) in discrete time: (i) a deterministic sampler based on a sort of ordinary differential equations (ODEs) called probability flow ODEs (which is closely related to DDIM); and (ii) a DDPM-type stochastic sampler motivated by reverse-time SDEs. In comparisons to past works, our main contributions are three-fold.

- *Non-asymptotic convergence guarantees.* For a popular deterministic sampler, we demonstrate that the number of steps needed to yield  $\varepsilon$ -accuracy — meaning that the total variance (TV) distance between the distribution of  $X_1$  and that of  $Y_1$  is no larger than  $\varepsilon$  — is proportional to  $1/\varepsilon$  (in addition to other polynomial dimension dependency), which improves upon prior results Chen et al. (2023b). For another DDPM-type stochastic sampler, we establish an iteration complexity proportional to  $1/\varepsilon^2$ , matching existing theory Chen et al. (2022b,a) in terms of the  $\varepsilon$ -dependency.
- *Accelerating data generation processes.* In order to further speed up the sampling processes, we develop an accelerated variant for each of the above two samplers, taking advantage of estimates of a small number of additional quantities. As it turns out, these variants achieve more rapid convergence, with the deterministic (resp. stochastic) variant exhibiting a  $1/\sqrt{\varepsilon}$  (resp.  $1/\varepsilon$ ) scaling in the accuracy level  $\varepsilon$ .
- *An elementary non-asymptotic analysis framework.* From the technical point of view, the analysis framework laid out in this paper is fully non-asymptotic in nature. In contrast to prior theoretical analyses that take a detour to study the continuum limits and then control the discretization error, our approach tackles the discrete-time processes directly using elementary analysis strategies. No knowledge of SDEs or ODEs is required for establishing our theory, thereby resulting in a more versatile framework and sometimes lowering the technical barrier towards understanding diffusion models.

**Notation.** Before proceeding, we introduce a couple of notation to be used throughout. For any two functions  $f(d, T)$  and  $g(d, T)$ , we adopt the notation  $f(d, T) \lesssim g(d, T)$  or  $f(d, T) = O(g(d, T))$  (resp.  $f(d, T) \gtrsim g(d, T)$ ) to mean that there exists some universal constant  $C_1 > 0$  such that  $f(d, T) \leq C_1 g(d, T)$  (resp.  $f(d, T) \geq C_1 g(d, T)$ ) for all  $d$  and  $T$ ; moreover, the notation  $f(d, T) \asymp g(d, T)$  indicates that  $f(d, T) \lesssim g(d, T)$  and  $f(d, T) \gtrsim g(d, T)$  hold at once. The notation  $\tilde{O}(\cdot)$  is defined similar to  $O(\cdot)$  except that it hides the logarithmic dependency. We shall often use capital letters to denote random variables/vectors/processes, and lowercase letters for deterministic variables. For any two probability measures  $P$  and  $Q$ , the total variation (TV) distance between them is defined to be  $\text{TV}(P, Q) := \frac{1}{2} \int |dP - dQ|$ .

## 2 Preliminaries

sec:preliminaries

In this section, we introduce the basics of diffusion generative models. The ultimate goal of a generative model can be concisely stated: given data samples drawn from an unknown distribution of interest  $p_{\text{data}}$  in

$\mathbb{R}^d$ , we wish to generate new samples whose distributions closely resemble  $p_{\text{data}}$ .

## 2.1 Diffusion generative models

Towards achieving the above goal, a diffusion generative model typically encompasses two Markov processes: a forward process and a reverse process, as described below.

**The forward process.** In the forward chain, one progressively injects noise into the data samples to diffuse and obscure the data. The distributions of the injected noise are often hand-picked, with the standard Gaussian distribution receiving widespread adoption. More specifically, the forward Markov process produces a sequence of  $d$ -dimensional random vectors  $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_T$  as follows: eq:forward-process

$$X_0 \sim p_{\text{data}}, \tag{4a}$$

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} W_t, \quad 1 \leq t \leq T, \tag{4b}$$

where  $\{W_t\}_{1 \leq t \leq T}$  indicates a sequence of independent noise vectors drawn from  $W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ . The hyper-parameters  $\{\beta_t \in (0, 1)\}$  represent prescribed learning rate schedules that control the variance of the noise injected in each step. If we define

$$\alpha_t := 1 - \beta_t, \quad \bar{\alpha}_t := \prod_{k=1}^t \alpha_k, \quad 1 \leq t \leq T, \tag{5}$$

then it can be straightforwardly verified that for every  $1 \leq t \leq T$ ,

$$X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \bar{W}_t \quad \text{for some } \bar{W}_t \sim \mathcal{N}(0, I_d). \tag{eqn:Xt-X0} \tag{6}$$

Clearly, if the covariance of  $X_0$  is also equal to  $I_d$ , then the covariance of  $X_t$  is preserved throughout the forward process; for this reason, this forward process (4) is sometimes referred to as *variance-preserving* (Song et al., 2021b). Throughout this paper, we employ the notation

$$q_t := \text{law}(X_t) \tag{eq:defn-qt} \tag{7}$$

to denote the distribution of  $X_t$ . As long as  $\bar{\alpha}_T$  is vanishingly small, one has the following property for a general family of data distributions:

$$q_T \approx \mathcal{N}(0, I_d). \tag{8}$$

**The reverse process.** The reverse chain  $Y_T \rightarrow Y_{T-1} \rightarrow \dots \rightarrow Y_1$  is designed to (approximately) revert the forward process, allowing one to transform pure noise into new samples with matching distributions as the original data. To be more precise, by initializing it as eq:goal-reverse-process

$$Y_T \sim \mathcal{N}(0, I_d), \tag{9a}$$

we seek to design a reverse-time Markov process with nearly identical marginals as the forward process, namely,

$$(\text{goal}) \quad Y_t \stackrel{d}{\approx} X_t, \quad t = T, T-1, \dots, 1. \tag{9b}$$

Throughout the paper, we shall often employ the following notation to indicate the distribution of  $Y_t$ :

$$p_t := \text{law}(Y_t). \tag{eq:defn-pt} \tag{10}$$

## 2.2 Deterministic vs. stochastic samplers: a continuous-time interpretation

sec:det-stochastic-samplers

Evidently, the most crucial step of the diffusion model lies in effective design of the reverse process. Two mainstream approaches stand out:

- *Deterministic samplers.* Starting from  $Y_T \sim \mathcal{N}(0, I_d)$ , this approach selects a set of functions  $\{\Phi_t(\cdot)\}_{1 \leq t \leq T}$  and computes:

$$Y_{t-1} = \Phi_t(Y_t), \quad t = T, \dots, 1. \quad \text{eq:deterministic-sampler} \quad (11)$$

Clearly, the sampling process is fully deterministic except for the initialization  $Y_T$ .

- *Stochastic samplers.* Initialized again at  $Y_T \sim \mathcal{N}(0, I_d)$ , this approach computes another collection of functions  $\{\Psi_t(\cdot, \cdot)\}_{1 \leq t \leq T}$  and performs the updates:

$$Y_{t-1} = \Psi_t(Y_t, Z_t), \quad t = T, \dots, 1, \quad \text{eq:stochastic-sampler} \quad (12)$$

where the  $Z_t$ 's are independent noise vectors obeying  $Z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ .

In order to elucidate the feasibility of the above two approaches, we find it helpful to look at the continuum limit through the lens of SDEs and ODEs. It is worth emphasizing, however, that the development of our main theory does *not* rely on any knowledge of SDEs and ODEs.

- *The forward process.* A continuous-time analog of the forward diffusion process can be modeled as

$$dX_t = f(X_t, t)dt + g(t)dW_t \quad (0 \leq t \leq T), \quad X_0 \sim p_{\text{data}} \quad \text{eq:forward-SDE-general} \quad (13)$$

for some functions  $f(\cdot, \cdot)$  and  $g(\cdot)$  (denoting respectively the drift and diffusion coefficient), where  $W_t$  denotes a  $d$ -dimensional standard Brownian motion. As a special example, the continuum limit of (4) takes the following form (Song et al., 2021b)

$$dX_t = -\frac{1}{2}\beta(t)X_tdt + \sqrt{\beta(t)}dW_t \quad (0 \leq t \leq T), \quad X_0 \sim p_{\text{data}} \quad \text{eq:forward-SDE} \quad (14)$$

for some function  $\beta(t)$ . As before, we denote by  $q_t$  the distribution of  $X_t$  in (13).

- *The reverse process.* As it turns out, the following two reverse processes are both capable of reconstructing the distribution of the forward process, motivating the design of two distinctive samplers. Here and throughout, we use  $\nabla \log q_t(X)$  to abbreviate  $\nabla_X \log q_t(X)$  for notational simplicity.

- One feasible approach is to resort to the so-called *probability flow ODE* (Song et al., 2021b)

$$dY_t^{\text{ode}} = \left( f(Y_t^{\text{ode}}, t) - \frac{1}{2}g(t)^2 \nabla \log q_t(Y_t^{\text{ode}}) \right) dt \quad (0 \leq t \leq T), \quad Y_0^{\text{ode}} \stackrel{\text{eq:prob-flow-ODE}}{\sim} q_T, \quad (15)$$

which exhibits matching distributions as follows:

$$Y_{T-t}^{\text{ode}} \stackrel{d}{=} X_t, \quad 0 \leq t \leq T.$$

The deterministic nature of this approach often enables faster sampling. It has been shown that this family of deterministic samplers is closely related to the DDIM sampler (Karras et al., 2022; Song et al., 2021b).

- In view of the seminal result by Anderson (1982), one can also construct a “reverse-time” SDE

$$dY_t^{\text{sde}} = \left( f(Y_t^{\text{sde}}, t) - g(t)^2 \nabla \log q_t(Y_t^{\text{sde}}) \right) dt + g(t)dZ_t^{\text{sde}} \quad (0 \leq t \leq T), \quad Y_0^{\text{sde}} \stackrel{\text{eq:reverse-SDE}}{\sim} q_T \quad (16)$$

with  $Z_t^{\text{sde}}$  being a standard Brownian motion, which also satisfies

$$Y_{T-t}^{\text{sde}} \stackrel{d}{=} X_t, \quad 0 \leq t \leq T.$$

The popular DDPM sampler (Ho et al., 2020; Nichol and Dhariwal, 2021) falls under this category.

Interestingly, in addition to the functions  $f$  and  $g$  that define the forward process, construction of both (15) and (16) relies only upon the knowledge of the gradient of the log density  $\nabla \log q_t(\cdot)$  of the intermediate steps of the forward diffusion process — often referred to as the (Stein) score function. Consequently, a key enabler of the above paradigms lies in reliable learning of the score function, and hence the name *score-based generative modeling*.

### 3 Algorithms and main results

sec:main-results

In this section, we analyze a couple of diffusion generative models, including both deterministic and stochastic samplers. While the proofs for our main theory are all postponed to the appendix, it is worth emphasizing upfront that our analysis framework directly tackles the discrete-time reverse process without resorting to any toolbox of SDEs and ODEs tailored to the continuous-time limits. This elementary approach might potentially be versatile for analyzing a broad class of variations of these samplers.

#### 3.1 Assumptions and learning rates

Before proceeding, we impose some assumptions on the score estimates and the target data distributions, and specify the hyper-parameters  $\{\alpha_t\}$ , which shall be adopted throughout all cases.

**Score estimates.** Given that the score functions are an essential component in score-based generative modeling, we assume access to faithful estimates of the score functions  $\nabla \log q_t(\cdot)$  across all intermediate steps  $t$ , thus disentangling the score learning phase and the data generation phase. This is made precise in the following assumption.

assumption:score-estimate

**Assumption 1.** Suppose that we have access to the score estimates  $s_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$  ( $1 \leq t \leq T$ ) as follows:

$$s_t := \arg \min_{s: \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{X \sim q_t} \left[ \|s(X) - \nabla \log q_t(X)\|_2^2 \right], \quad 1 \leq t \leq T. \quad \text{eqn:training-score} \quad (17)$$

As has been pointed out by previous works concerning score matching (e.g., Hyvärinen (2005); Vincent (2011); Chen et al. (2022b)), this score estimate admits an alternative form as follows (owing to properties of Gaussian distributions):

$$s_t := \arg \min_{s: \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{W \sim \mathcal{N}(0, I_d), X_0 \sim p_{\text{data}}} \left[ \left\| s(\sqrt{\alpha_t} X_0 + \sqrt{1 - \alpha_t} W) + \frac{\nabla \log q_t(X_0)}{\sqrt{1 - \alpha_t}} \right\|_2^2 \right], \quad \text{eqn:training-score-equiv} \quad (18)$$

which is often more amenable to training.

**Target data distributions.** Our goal is to uncover the effectiveness of diffusion models in generating a broad family of data distributions. Throughout this paper, the only assumption we need to impose on the target data distribution  $p_{\text{data}}$  is the following:

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_3} \mid X_0 \sim p_{\text{data}}) = 1 \quad \text{eq:assumption-data-bounded} \quad (19)$$

for some arbitrarily large constant  $c_3 > 0$ . This assumption allows the radius of the support of  $p_{\text{data}}$  to be exceedingly large (given that the exponent  $c_3$  can be arbitrarily large).

**Learning rate schedule.** Let us also take a moment to specify the learning rates to be used for our theory and analyses. For some large enough numerical constants  $c_0, c_1 > 0$ , we set | eqn:alpha-t

$$\beta_1 = 1 - \alpha_1 = \frac{1}{T^{c_0}}; \quad (20a)$$

$$\beta_t = 1 - \alpha_t = \frac{c_1 \log T}{T} \min \left\{ \beta_1 \left( 1 + \frac{c_1 \log T}{T} \right)^t, 1 \right\}. \quad (20b)$$

We immediately single out several properties about the choice (20) that prove useful throughout: | eqn:properties-alpha

$$\frac{1 - \alpha_t}{1 - \alpha_{t-1}} \leq \frac{4c_1 \log T}{T}, \quad 1 \leq t \leq T \quad (21a)$$

$$\bar{\alpha}_T \leq \frac{1}{T^{c_2}} \quad (21b)$$

for some large enough numerical constant  $c_2 > 0$ . The proof of these properties can be found in Appendix (A.4).

## 3.2 Deterministic samplers

We begin by analyzing a deterministic sampler: a discrete-time version of the probability flow ODE.

### 3.2.1 An ODE-based deterministic sampler

Armed with the score estimates in Assumption 1, a discrete-time version of the probability flow ODE approach (cf. (15)) adopts the following update rule: eqn:ode-sampling

$$Y_T \sim \mathcal{N}(0, I_d), \quad Y_{t-1} = \Phi_t(Y_t) \quad \text{for } t = T, \dots, 1, \quad \text{eqn:ode-sampling-Y} \quad (22a)$$

where  $\Phi_t(\cdot)$  is taken to be

$$\Phi_t(x) := \frac{1}{\sqrt{\alpha_t}} \left( x + \frac{1 - \alpha_t}{2} s_t(x) \right). \quad \text{eqn:phi-func} \quad (22b)$$

This approach, based on the probability flow ODE (15), often achieves faster sampling compared to the stochastic counterpart (Song et al., 2021b). Despite the empirical advances, however, the theoretical understanding of this type of deterministic samplers remained far from mature.

We first derive non-asymptotic convergence guarantees — measured by the total variance distance between the forward and the reverse processes — for the above deterministic sampler (22).

**Theorem 1.** Suppose that (19) holds true. Equipped with the score estimates in Assumption 1 and the learning rate schedule (20), the sampling process (22) satisfies thm:main-ODE

$$\text{TV}(q_1, p_1) \leq C_1 \frac{d^2 \log^4 T}{T} + C_1 \frac{d^6 \log^6 T}{T^2} \quad \text{eq:ratio-ODE} \quad (23)$$

for some universal constants  $C_1 > 0$ , where we recall that  $p_1$  (resp.  $q_1$ ) represents the distribution of  $Y_1$  (resp.  $X_1$ ).

In other words, in order to achieve  $\text{TV}(q_1, p_1) \leq \varepsilon$ , the number of steps  $T$  only needs to exceed

$$\tilde{O}\left(\frac{d^2}{\varepsilon} + \frac{d^3}{\sqrt{\varepsilon}}\right). \quad \text{eq:iteration-complexity-ODE} \quad (24)$$

To the best of our knowledge, the only non-asymptotic analysis for such deterministic samplers in prior literature was derived by a very recent work Chen et al. (2023b), which established the first polynomial-time convergence guarantees (see, e.g., Chen et al. (2023b, Theorem 4.1)). However, it fell short of providing concrete polynomial dependency. In contrast, our result in Theorem 1 uncovers a concrete  $d^2/\varepsilon$  scaling (ignoring lower-order and logarithmic terms), which was previously unavailable.

[Yuxin: might need a discussion with the new paper]

### 3.2.2 An accelerated deterministic sampler

Thus far, we have demonstrated that the iteration complexity of the deterministic sampler (22) is proportional to  $1/\varepsilon$  (for small enough  $\varepsilon$ ). A natural question is whether this convergence rate can be further improved.

As it turns out, if we have access to reliable estimates of a few additional scalar quantities, then a modified version of the sampler (22) is able to achieve much improved convergence guarantees. These estimates are made precise in the following assumption.

**Assumption 2.** Suppose that we have access to the estimates  $A_t, \dots, E_t : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $F_t : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  ( $1 \leq t \leq T$ ) defined as follows: assumption:ODE-estimate  
eqn:training-ODE

$$A_t := \arg \min_{A: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \|W\|_2^2 - A(\sqrt{\alpha_t} X_0 + \sqrt{1 - \alpha_t} W) \right\}^2 \right], \quad (25a)$$

$$B_t := \arg \min_{B: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \sqrt{1 - \alpha_t} W^\top s_t(\sqrt{\alpha_t} X_0 + \sqrt{1 - \alpha_t} W) + B(\sqrt{\alpha_t} X_0 + \sqrt{1 - \alpha_t} W) \right\}^2 \right], \quad (25b)$$



$$C_t := \arg \min_{C: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \|W\|_2^4 - C(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W) \right\}^2 \right], \quad (25c)$$

$$D_t := \arg \min_{D: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \left( \sqrt{1-\alpha_t}W^\top s_t(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W) \right)^2 - D(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W) \right\}^2 \right], \quad (25d)$$

$$E_t := \arg \min_{E: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \sqrt{1-\alpha_t} \|W\|_2^2 W^\top s_t(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W) + E(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W) \right\}^2 \right], \quad (25e)$$

$$F_t := \arg \min_{F: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} \left[ \left\{ \langle M, M' \rangle - F(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W, \sqrt{\alpha_t}X'_0 + \sqrt{1-\alpha_t}W') \right\}^2 \right], \quad (25f)$$

where  $s_t(\cdot)$  is defined in Assumption 1, and  $\text{eq: defn-M-Mprime}$

$$M := (1 - \alpha_t)s_t(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W)s_t(\sqrt{\alpha_t}X_0 + \sqrt{1-\alpha_t}W)^\top - WW^\top; \quad (26a)$$

$$M' := (1 - \alpha_t)s_t(\sqrt{\alpha_t}X'_0 + \sqrt{1-\alpha_t}W')s_t(\sqrt{\alpha_t}X'_0 + \sqrt{1-\alpha_t}W')^\top - W'W'^\top. \quad (26b)$$

Here, the expectation is with respect to  $W, W' \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$  and  $X_0, X'_0 \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}$ .

Armed with the score estimate in Assumption 1 and the additional scalar estimates in Assumption 2, we are ready to introduce an accelerated variant of (22) as follows:  $\text{eqn: ode-sampling-R}$

$$Y_T \sim \mathcal{N}(0, I_d), \quad Y_{t-1} = \Phi_t(Y_t) \quad \text{for } t = T, \dots, 1, \quad \text{eqn: ode-sampling-R-Y} \quad (27a)$$

where the mapping  $\Phi_t(\cdot)$  is chosen to be

$$\Phi_t(x) = \frac{1}{\sqrt{\alpha_t}} \left( x + \frac{\lambda_t(x)(1 - \alpha_t)}{2} s_t(x) \right), \quad (27b)$$

Here,  $\lambda_t: \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$\lambda_t(x) = 1 + \frac{(1 - \alpha_t)w_t(x)}{4(\alpha_t - \bar{\alpha}_t)}, \quad (27c)$$

where  $w_t: \mathbb{R}^d \mapsto \mathbb{R}$  is chosen to be a solution to the following equation:

$$(1 - \bar{\alpha}_t)s_t^\top \nabla w + (d - A_t)w + (d + B_t - A_t)^2 - 2d^2 + 7d - 6A_t + 7B_t - C_t - D_t + 2E_t + \tilde{F}_t = 0 \quad (27d)$$

with  $\tilde{F}_t(x) := F_t(x, x)$ . Notably, this new variant (27) is closely related to the original sampler (22); in fact, they both move along the direction specified by the score estimate  $s_t$ , except that the stepsize of the accelerated variant (mainly the quantity  $\lambda_t$ ) is adaptively chosen based on additional information.

Encouragingly, our non-asymptotic analysis framework can be extended to derive enhanced convergence guarantees for the sampler (27), as stated below.

**Theorem 2.** Suppose that (19) holds true. Equipped with the estimates in Assumptions 1-2 and the learning rate schedule (20), the sampling process (27) obeys  $\text{thm: main-ODE-fast}$

$$\text{TV}(q_1, p_1) \leq C_1 \frac{d^6 \log^6 T}{T^2} \quad \text{eq: ratio-ODE-R} \quad (28)$$

for some universal constants  $C_1 > 0$ , where  $p_1$  (resp.  $q_1$ ) is the distribution of  $Y_1$  (resp.  $X_1$ ).

Theorem 2 reveals that: in order to achieve  $\text{TV}(q_1, p_1) \leq \varepsilon$ , the accelerated deterministic sampler (27) only requires the number of steps  $T$  to be on the order of

$$\tilde{O}\left(\frac{d^3}{\sqrt{\varepsilon}}\right), \quad \text{eq: iteration-complexity-ODE-fast} \quad (29)$$

thus improving the dependency on  $\varepsilon$  from  $\tilde{O}(1/\varepsilon)$  (cf. (24)) to  $\tilde{O}(1/\sqrt{\varepsilon})$  for small enough  $\varepsilon$ . Consequently, the improved convergence result underscores the crucial role of stepsize selection even when the search direction has been determined.



### 3.3 Stochastic samplers

#### 3.3.1 A DDPM-type stochastic sampler

Armed with the score estimates  $\{s_t\}$  in Assumption 1, we can readily introduce the following stochastic sampler that operates in discrete time, motivated by the reverse-time SDE (16): eqn:sde-sampling

$$Y_T \sim \mathcal{N}(0, I_d), \quad Y_{t-1} = \Psi_t(Y_t, Z_t) \quad \text{for } t = T, \dots, 1 \quad (30a)$$

where  $Z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ , and

$$\Psi_t(y, z) = \frac{1}{\sqrt{\alpha_t}} \left( y + (1 - \alpha_t) s_t(y) \right) + \sigma_t z \quad \text{with } \sigma_t^2 = \frac{1}{\alpha_t} - 1. \quad (30b)$$

The key difference between this sampler and the deterministic sampler (22) is that: (i) there exists an additional pre-factor of  $1/2$  on  $s_t$  in the deterministic sampler; and (ii) the stochastic sampler injects additional noise  $Z_t$  in each step.

In contrast to deterministic samplers, the stochastic samplers have received more theoretical attention, with the state-of-the-art results established by Chen et al. (2022b,a). The elementary approach developed in the current paper is also applicable towards understanding this type of samplers, leading to the following non-asymptotic theory.

thm:main-SDE

**Theorem 3.** Suppose that (19) holds true. Equipped with the estimates in Assumption 1 and the learning rate schedule (20), the stochastic sampler (30) achieves

$$\text{TV}(q_1, p_1) \leq C_1 \frac{d^2 \log^3 T}{\sqrt{T}} \quad \text{eq:ratio-SDE} \quad (31)$$

for some universal constants  $C_1 > 0$ .

Theorem 3 establishes non-asymptotic convergence guarantees for the stochastic sampler (30). As asserted by the theorem, the number of steps needed to attain  $\varepsilon$ -accuracy (measured by the TV distance between  $p_1$  and  $q_1$ ) is proportional to  $1/\varepsilon^2$ , matching the state-of-the-art  $\varepsilon$ -dependency derived in Chen et al. (2022a), albeit exhibiting a worse dimensional dependency. Our analysis follows a completely different path compared with the SDE-based approach in Chen et al. (2022a), thus offering complementary interpretations for this important sampler. In order to further illustrate the versatility of our analysis approach, we shall demonstrate how it can be applied to study an accelerated version in the next subsection.

#### 3.3.2 An accelerated stochastic sampler

In this subsection, we come up with a potential strategy to speed up the stochastic sampler (30), assuming access to reliable estimates of additional objects as described below.

assumption:variance-estimate

**Assumption 3.** Suppose that we have access to the estimates  $v_t : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  ( $1 \leq t \leq T$ ) as follows:

$$v_t := \arg \min_{v: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E} \left[ \left\| WW^\top Z - v(\sqrt{\alpha_t} X + \sqrt{1 - \alpha_t} W, Z) \right\|_2^2 \right], \quad 1 \leq t \leq T, \quad \text{eqn:training-variance} \quad (32)$$

where  $X, W, Z$  are independently generated obeying  $X \sim p_{\text{data}}$ ,  $W \sim \mathcal{N}(0, I_d)$ , and  $Z \sim \mathcal{N}(0, I_d)$ .

With the estimates in Assumption 1 and Assumption 3 in place, we are positioned to introduce the proposed accelerated sampler as follows: eqn:sde-sampling-R

$$Y_T \sim \mathcal{N}(0, I_d), \quad Y_{t-1} = \Psi_t(Y_t, Z_t) \quad \text{for } t = T, \dots, 1, \quad \text{eqn:sde-sampling-R-Y} \quad (33a)$$

where we choose the mapping  $\Psi_t(\cdot, \cdot)$  as follows

$$\Psi_t(y, z) = \frac{1}{\sqrt{\alpha_t}} \left( y + (1 - \alpha_t) s_t(y) \right) + \sigma_t \left\{ z - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} \left[ z + (1 - \bar{\alpha}_t) s_t(y) (s_t(y))^\top z - v_t(y, z) \right] \right\} \quad \text{eqn:sde-sampling-R-Psi} \quad (33b)$$

with

$$\sigma_t^2 = \frac{1}{\alpha_t} - 1. \quad (33c)$$

Clearly, the modified update mapping (33b) is still mainly a linear combination of the score estimate  $s_{t-1}$  and the additive noise  $Z_t$ , except that a correction term  $v_t$  (learned by solving (32)) needs to be included for acceleration purposes.

We now apply our analysis strategy to establish performance guarantees for the above stochastic sampler.

thm:main-SDE-R

**Theorem 4.** *Suppose that (19) holds true. Equipped with the estimates in Assumption 1, 3 and the learning rate schedule (20), the sampling process (33) satisfies*

$$\text{TV}(q_1, p_1) \leq C_1 \frac{d^3 \log^{4.5} T}{T} \quad \text{eq:ratio-SDE-R} \quad (34)$$

for some universal constants  $C_1 > 0$ .

In comparison to the stochastic sampler (30), Theorem 4 asserts that the iteration complexity of the sampler (33) is at most

$$\tilde{O}\left(\frac{d^3}{\varepsilon}\right), \quad \text{eq:iteration-complexity-SDE-fast} \quad (35)$$

thus significantly reducing the scaling  $\tilde{O}(1/\varepsilon^2)$  for the original sampler (30) to  $\tilde{O}(1/\varepsilon)$  regarding the  $\varepsilon$ -dependency. All in all, our theory reveals that having information about a small number of additional objects might substantially speed up the data generation process.

## 4 Other related works

sec:related-works

**Theory for SGMs.** Early theoretical efforts in understanding the convergence of score-based stochastic samplers suffered from being either not quantitative (De Bortoli et al., 2021; Liu et al., 2022; Pidstrigach, 2022), or the curse of dimensionality (e.g., exponential dependencies in the convergence guarantees) (Block et al., 2020; De Bortoli, 2022). The recent work Lee et al. (2022) provided the first polynomial convergence guarantee in the presence of  $L_2$ -accurate score estimates, for any smooth distribution satisfying the log-Sobolev inequality, effectively only allowing unimodal distributions though. Chen et al. (2022b); Lee et al. (2023); Chen et al. (2022a) subsequently lifted such a stringent data distribution assumption. More concretely, Chen et al. (2022b) accommodated a broad family of data distributions under the premise that the score functions over the entire trajectory of the forward process are Lipschitz; Lee et al. (2023) only required certain smoothness assumptions but came with worse dependence on the problem parameters; and more recent results in Chen et al. (2022a) applied to literally any data distribution with bounded second-order moment. In addition, Wibisono and Yang (2022) also established a convergence theory for score-based generative models, assuming that the error of the score estimator has a bounded moment generating function and that the data distribution satisfies the log-Sobolev inequality. Turning attention to deterministic samplers, Chen et al. (2023b) derived the first non-asymptotic bounds for the probability flow ODE. Additionally, theoretical justifications for DDPM in the context of image in-painting have been developed by Rout et al. (2023).

**Score matching.** Hyvärinen (2005) showed that the score function can be estimated via integration by parts, a result that was further extended in Hyvärinen (2007). Song et al. (2020b) proposed sliced score matching to tame the computational complexity in high dimension. The consistency of the score matching estimator was studied in Hyvärinen (2005), with asymptotic normality established in Forbes and Lauritzen (2015). Optimizing the score matching loss has been shown to be intimately connected to minimizing upper bounds on the Kullback-Leibler divergence (Song et al., 2021a) and Wasserstein distance (Kwon et al., 2022) between the generated distribution and the target data distribution. From a non-asymptotic perspective, Koehler et al. (2022) studied the statistical efficiency of score matching by connecting it with the isoperimetric properties of the distribution.

**Other theory for diffusion models.** [Oko et al. \(2023\)](#) studied the approximation and generalization capabilities of diffusion modeling for distribution estimation. Assuming that the data are supported on a low-dimensional linear subspace, [Chen et al. \(2023a\)](#) developed a sample complexity bound for diffusion models. Moreover, [Ghimire et al. \(2023\)](#) adopted a geometric perspective and showed that the forward and backward processes of diffusion models are essentially Wasserstein gradient flows operating in the space of probability measures. Recently, the idea of stochastic localization, which is closely related to diffusion models, is adopted to sample from posterior distributions ([Montanari and Wu, 2023](#); [El Alaoui et al., 2022](#)), which has been implemented using the approximate message passing algorithm ([Donoho et al. \(2009\)](#); [Li and Wei \(2022\)](#)).

## 5 Discussion

sec:discussion

In this paper, we have developed a new suite of non-asymptotic theory for establishing the convergence and faithfulness of diffusion generative modeling, assuming access to reliable estimates of the (Stein) score functions. Our analysis framework seeks to track the dynamics of the reverse process directly using elementary tools, which eliminates the need to look at the continuous-time limit and invoke the SDE and ODE toolboxes. In addition to demonstrating the non-asymptotic iteration complexities of two mainstream discrete-time samplers — a deterministic sampler based on the probability flow ODE, and a DDPM-type stochastic sampler — we have discovered potential strategies to further accelerate the sampling process, taking advantage of estimates of a small number of additional objects. The analysis framework laid out in the current paper might shed light on how to analyze other variants of score-based generative models as well.

Moving forward, there are plenty of questions that require in-depth theoretical understanding. For instance, the dimension dependency in our convergence results remains sub-optimal; can we further refine our theory in order to reveal tight dependency in this regard? To what extent can we further accelerate the sampling process, without requiring much more information than the score functions? It would also be of paramount interest to establish end-to-end performance guarantees that take into account both the score learning phase and the sampling phase.

## Acknowledgements

Y. Wei is supported in part by the the NSF grants DMS-2147546/2015447, CAREER award DMS-2143215, CCF-2106778, and the Google Research Scholar Award. Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR grant FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009 and CCF-1907661. Y. Chi are supported in part by the grants ONR N00014-19-1-2404, NSF CCF-2106778, DMS-2134080 and CNS-2148212.

## References

- Anderson, B. D. (1982). Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326.
- Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and van den Berg, R. (2021). Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993.
- Block, A., Mroueh, Y., and Rakhlin, A. (2020). Generative modeling with denoising auto-encoders and Langevin sampling. *arXiv preprint arXiv:2002.00107*.
- Chen, H., Lee, H., and Lu, J. (2022a). Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. *arXiv preprint arXiv:2211.01916*.
- Chen, M., Huang, K., Zhao, T., and Wang, M. (2023a). Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. *arXiv preprint arXiv:2302.07194*.

- Chen, N., Zhang, Y., Zen, H., Weiss, R. J., Norouzi, M., and Chan, W. (2021). WaveGrad: Estimating gradients for waveform generation. In *International Conference on Learning Representations*.
- Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., and Zhang, A. R. (2022b). Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint arXiv:2209.11215*.
- Chen, S., Daras, G., and Dimakis, A. G. (2023b). Restoration-degradation beyond linear diffusions: A non-asymptotic analysis for DDIM-type samplers. *arXiv preprint arXiv:2303.03384*.
- Croitoru, F.-A., Hondru, V., Ionescu, R. T., and Shah, M. (2023). Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- De Bortoli, V. (2022). Convergence of denoising diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2208.05314*.
- De Bortoli, V., Thornton, J., Heng, J., and Doucet, A. (2021). Diffusion Schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709.
- Dhariwal, P. and Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794.
- Donoho, D. L., Maleki, A., and Montanari, A. (2009). Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919.
- El Alaoui, A., Montanari, A., and Sellke, M. (2022). Sampling from the sherrington-kirkpatrick gibbs measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 323–334.
- Forbes, P. G. and Lauritzen, S. (2015). Linear estimating equations for exponential families with application to Gaussian linear concentration models. *Linear Algebra and its Applications*, 473:261–283.
- Ghimire, S., Liu, J., Comas, A., Hill, D., Masoomi, A., Camps, O., and Dy, J. (2023). Geometry of score based generative models. *arXiv preprint arXiv:2302.04411*.
- Hajek, B. (2015). *Random processes for engineers*. Cambridge university press.
- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851.
- Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4).
- Hyvärinen, A. (2007). Some extensions of score matching. *Computational statistics & data analysis*, 51(5):2499–2512.
- Jolicœur-Martineau, A., Piché-Taillefer, R., Mitliagkas, I., and des Combes, R. T. (2021). Adversarial score matching and improved sampling for image generation. In *International Conference on Learning Representations*.
- Karras, T., Aittala, M., Aila, T., and Laine, S. (2022). Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems*, volume 35, pages 26565–26577.
- Koehler, F., Hekett, A., and Risteski, A. (2022). Statistical efficiency of score matching: The view from isoperimetry. *arXiv preprint arXiv:2210.00726*.
- Kong, Z., Ping, W., Huang, J., Zhao, K., and Catanzaro, B. (2021). DiffWave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*.
- Kwon, D., Fan, Y., and Lee, K. (2022). Score-based generative modeling secretly minimizes the wasserstein distance. In *Advances in Neural Information Processing Systems*.

- Lee, H., Lu, J., and Tan, Y. (2022). Convergence for score-based generative modeling with polynomial complexity. In *Advances in Neural Information Processing Systems*.
- Lee, H., Lu, J., and Tan, Y. (2023). Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985.
- Li, G. and Wei, Y. (2022). A non-asymptotic framework for approximate message passing in spiked models. *arXiv preprint arXiv:2208.03313*.
- Liu, X., Wu, L., Ye, M., and Liu, Q. (2022). Let us build bridges: Understanding and extending diffusion generative models. *arXiv preprint arXiv:2208.14699*.
- Montanari, A. and Wu, Y. (2023). Posterior sampling from the spiked models via diffusion processes. *arXiv preprint arXiv:2304.11449*.
- Nichol, A. Q. and Dhariwal, P. (2021). Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171.
- Oko, K., Akiyama, S., and Suzuki, T. (2023). Diffusion models are minimax optimal distribution estimators. *arXiv preprint arXiv:2303.01861*.
- Pidstrigach, J. (2022). Score-based generative models detect manifolds. *arXiv preprint arXiv:2206.01018*.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. (2022). Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695.
- Rout, L., Parulekar, A., Caramanis, C., and Shakkottai, S. (2023). A theoretical justification for image inpainting using denoising diffusion probabilistic models. *arXiv preprint arXiv:2302.01217*.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al. (2022). Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265.
- Song, J., Meng, C., and Ermon, S. (2020a). Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*.
- Song, Y., Durkan, C., Murray, I., and Ermon, S. (2021a). Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–1428.
- Song, Y. and Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.
- Song, Y., Garg, S., Shi, J., and Ermon, S. (2020b). Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. (2021b). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*.
- Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674.
- Wibisono, A. and Yang, K. Y. (2022). Convergence in KL divergence of the inexact Langevin algorithm with application to score-based generative models. *arXiv preprint arXiv:2211.01512*.

Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Shao, Y., Zhang, W., Cui, B., and Yang, M.-H. (2022). Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*.

## A Analysis for the deterministic samplers (Theorems 1 and 2)

In this section, we present the proofs for our main results tailored to the two deterministic samplers.

### A.1 Preliminary facts

Before proceeding, we gather a couple of useful facts for the proof. First of all, in view of the alternative expression (18) for the score estimate  $s_t$  and the property of the minimum mean square error estimator (e.g., Hajek (2015, Section 3.3.1)), we know that  $s_t$  is given by the conditional expectation

$$\begin{aligned} s_t(x) &= \mathbb{E} \left[ -\frac{1}{\sqrt{1-\bar{\alpha}_t}} W \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} W = x \right] = \frac{1}{1-\bar{\alpha}_t} \mathbb{E} [\sqrt{\bar{\alpha}_t} X_0 - x \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} W = x] \\ &= \frac{1}{1-\bar{\alpha}_t} \int_{x_0} (\sqrt{\bar{\alpha}_t} x_0 - x) p_{X_0|X_t}(x_0 | x) dx_0. \end{aligned} \quad \text{eq:st-MMSE-expression} \quad (36)$$

Next, given that  $X_t \stackrel{d}{=} \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} W$  with  $W \sim \mathcal{N}(0, I_d)$ , we have the following tail bound for the random vector  $X_0$  conditional on  $X_t$ . lem:x0

**Lemma 1.** *Suppose that there exists some numerical constant  $c_R > 0$  obeying*

$$\mathbb{P}(\|X_0\|_2 \leq R) = 1 \quad \text{and} \quad R = T^{c_R}. \quad \text{eq:cR-defn-lem} \quad (37)$$

Consider any  $y \in \mathbb{R}$  obeying

$$-\log p_{X_t}(y) \leq c_6 d \log T \quad \text{eqn:choice-y} \quad (38)$$

for some large enough constant  $c_6 > 0$ , and let  $X$  be a random vector whose distribution  $p_X(\cdot)$  obeys

$$p_X(x) = p_{X_0|X_t}(x | y). \quad \text{eq:pX-properties} \quad (39)$$

Then for any  $c_5 \geq \sqrt{2 + c_0/2 + c_R}$  (with  $c_0$  defined in (20)), one has

$$\|\sqrt{\bar{\alpha}_t} X - y\|_2 \leq 5c_5 \sqrt{d(1-\bar{\alpha}_t) \log T} \quad (40)$$

with probability at least  $1 - \exp(-c_5^2 d \log T)$ .

*Proof.* See Appendix A.5. □

Additionally, we isolate a few useful properties about the learning rates specified by  $\{\alpha_t\}$  in (20): eqn:properties-alpha

$$\frac{1-\alpha_t}{1-\bar{\alpha}_{t-1}} \leq \frac{4c_1 \log T}{T}, \quad 1 \leq t \leq T \quad \text{eqn:properties-alpha-proof-1} \quad (41a)$$

$$\bar{\alpha}_T \leq \frac{1}{T^{c_2}} \quad \text{eqn:properties-alpha-proof-alphaT} \quad (41b)$$

for some large enough numerical constant  $c_2 > 0$ . In addition, if  $\frac{d(1-\alpha_t)}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ , then one has

$$\left( \frac{1-\bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} = 1 + \frac{d(1-\alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{d(d-2)(1-\alpha_t)^2}{8(\alpha_t - \bar{\alpha}_t)^2} + O\left(d^3 \left( \frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3\right). \quad \text{eq:expansion-ratio-1-alpha} \quad (41c)$$

The proof of these properties is postponed to Appendix A.4.

### A.2 Proof of Theorem 1

Before proceeding, we find it convenient to introduce a function

$$\phi_t(x) = x + \frac{1-\alpha_t}{2} s_t(x) = x - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} \int_{x_0} (x - \sqrt{\bar{\alpha}_t} x_0) p_{X_0|X_t}(x_0 | x) dx_0, \quad \text{defn:phi-t-x} \quad (42)$$

where the second identity follows from (36). This allows us to express the update rule (22) as follows:

$$Y_{t-1} = \Phi_t(Y_t) = \frac{1}{\sqrt{\alpha_t}} \phi_t(Y_t). \quad \text{eq:Yt-phi-ODE} \quad (43)$$

Our proof consists of three steps below.



**Step 1: bounding the density ratios of interest.** To begin with, we note that for any vectors  $y_{t-1}$  and  $y_t$ , elementary properties about transformation of probability distributions give

$$\begin{aligned} \frac{p_{Y_{t-1}}(y_{t-1})}{p_{X_{t-1}}(y_{t-1})} &= \frac{p_{\sqrt{\alpha_t}Y_{t-1}}(\sqrt{\alpha_t}y_{t-1})}{p_{\sqrt{\alpha_t}X_{t-1}}(\sqrt{\alpha_t}y_{t-1})} \\ &= \frac{p_{\sqrt{\alpha_t}Y_{t-1}}(\sqrt{\alpha_t}y_{t-1})}{p_{Y_t}(y_t)} \cdot \left( \frac{p_{\sqrt{\alpha_t}X_{t-1}}(\sqrt{\alpha_t}y_{t-1})}{p_{X_t}(y_t)} \right)^{-1} \cdot \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)}, \end{aligned} \quad \text{eq:recursion (44)}$$

thus converting the density ratio of interest into the product of three other density ratios. Motivated by this expression, we develop a lemma related to some of these density ratios, which plays a crucial role in establishing Theorem 1. The proof of this result is postponed to Appendix A.6.

**Lemma 2.** Suppose  $\frac{d^2(1-\alpha_t)\log T}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ . For every  $x \in \mathbb{R}$  obeying  $-\log p_{X_t}(x) \leq c_6 d \log T$  for some large enough constant  $c_6 > 0$ , it holds that lem:main-ODE

$$\begin{aligned} \frac{p_{\sqrt{\alpha_t}X_{t-1}}(\phi_t(x))}{p_{X_t}(x)} &= 1 + \frac{d(1-\alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + O\left(d^2 \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^2 \log^2 T\right) \\ &+ \frac{(1-\alpha_t) \left( \left\| \int_{x_0} (x - \sqrt{\alpha_t}x_0) p_{X_0|X_t}(x_0|x) dx_0 \right\|_2^2 - \int_{x_0} \|x - \sqrt{\alpha_t}x_0\|_2^2 p_{X_0|X_t}(x_0|x) dx_0 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)}. \end{aligned} \quad \text{eq:xt (45a)}$$

Moreover, for any random vector  $X$ , one has

$$\begin{aligned} \frac{p_{\phi_t(X)}(\phi_t(x))}{p_X(x)} &= 1 + \frac{d(1-\alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + O\left(d^2 \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^2 \log^2 T + d^6 \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3 \log^3 T\right) \\ &+ \frac{(1-\alpha_t) \left( \left\| \int_{x_0} (x - \sqrt{\alpha_t}x_0) p_{X_0|X_t}(x_0|x) dx_0 \right\|_2^2 - \int_{x_0} \|x - \sqrt{\alpha_t}x_0\|_2^2 p_{X_0|X_t}(x_0|x) dx_0 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)}. \end{aligned} \quad \text{eq:yt (45b)}$$

**Step 2: bounding the TV distance with the aid of a high-probability event.** In order to bound the TV distance of interest, we find it helpful to single out the following set:

$$\mathcal{E}_t := \left\{ y : \left| \frac{q_t(y)}{p_t(y)} - 1 \right| \leq c_5 \left( \frac{d^2 \log^4 T}{T^2} + \frac{d^6 \log^6 T}{T^3} \right) (T - t + 1) \right\}$$

for some large enough constant  $c_5 > 0$ . Informally speaking, this set  $\mathcal{E}_t$  contains all points such that  $q_t(y) \approx p_t(y)$ . We now claim that the event  $Y_t \in \mathcal{E}_t$  occurs with high probability, in the sense that

$$\mathbb{P}(Y_t \in \mathcal{E}_t) \geq 1 - (T - t + 1) \exp(-c_3 d \log T), \quad t \geq 1 \quad \text{eq:claim-Py (46)}$$

for some constant  $c_3 > 0$ . Suppose that this claim is valid for the moment. Then taking  $t = 1$  leads to

$$\mathbb{P}(Y_1 \in \mathcal{E}_1) = \mathbb{P}\left( \left| \frac{q_1(Y_1)}{p_1(Y_1)} - 1 \right| \leq c_5 \left( \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2} \right) \right) \geq 1 - \exp(-c_3 d \log T), \quad (47)$$

which also reveals that

$$\begin{aligned} \int_{y \notin \mathcal{E}_1} |p_1(y) - q_1(y)| dy &\leq c_5 \left( \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2} \right) \int_{y \notin \mathcal{E}_1} p_1(y) dy = c_5 \left( \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2} \right) \mathbb{P}(Y_1 \notin \mathcal{E}_1) \\ &\leq c_5 \left( \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2} \right) \exp(-c_3 d \log T) \leq \exp(-c_3 d \log T). \end{aligned}$$

The above results taken together with the definition of the total variation distance gives

$$\text{TV}(q_1, p_1) = \frac{1}{2} \int_x |q_1(x) - p_1(x)| dx = \frac{1}{2} \mathbb{E}_{Y_1 \sim p_1} \left[ \left| \frac{q_1(Y_1)}{p_1(Y_1)} - 1 \right| \cdot \mathbf{1}\{Y_1 \in \mathcal{E}_1\} \right] + \frac{1}{2} \int_{y \notin \mathcal{E}_1} |p_1(y) - q_1(y)| dy$$

$$\begin{aligned} &\leq c_5 \left( \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2} \right) + \frac{1}{2} \exp(-c_3 d \log T) \\ &\asymp \frac{d^2 \log^4 T}{T} + \frac{d^6 \log^6 T}{T^2}. \end{aligned}$$

This establishes the advertised result in Theorem 1, provided that Claim (46) can be verified.

**Step 3: justifying the claim (46).** We would like to prove this claim by induction, for which we start with the base case with  $t = T$ . Recall that  $X_T \stackrel{d}{=} \sqrt{\alpha_T} X_0 + \sqrt{1 - \alpha_T} B$  and  $Y_T \stackrel{d}{=} B$  with  $B \sim \mathcal{N}(0, I_d)$  independent of  $X_0$ , and that  $\|X_0\|_2 \leq R$  with  $R = T^{c_R}$  for some constant  $c_R > 0$ . For large enough  $T$ , it immediately follows from (41b) that

$$\mathbb{P}(Y_T \in \mathcal{E}_T) \geq 1 - \exp(-c_3 d \log T) \quad \text{eq:claim-Py-1} \quad (48)$$

for some constant  $c_3 > 0$  large enough.

Suppose now that the claim (46) holds for some  $t \geq 2$ , and we wish to prove the claim for  $t - 1$ . We would first like to claim that with probability at least  $1 - (T - t) \exp(-c_3 d \log T)$  for some constant  $c_3 > 0$ , one has

$$q_t(Y_t) \geq \exp(-c_6 d \log T). \quad \text{eqn:x-paganini} \quad (49)$$

With (49) in place, one sees that Lemma 2 is applicable to  $x = Y_t$  with high probability.

Next, observe that for any  $y$  obeying  $-\log p_{X_t}(y) \leq c_6 \log T$ , applying relations (45a) and (45b) in Lemma 2 leads to

$$\begin{aligned} \frac{p_{\sqrt{\alpha_t} Y_{t-1}}(\phi_t(y))}{p_{Y_t}(y)} \left( \frac{p_{\sqrt{\alpha_t} X_{t-1}}(\phi_t(y))}{p_{X_t}(y)} \right)^{-1} &= \frac{p_{\phi_t(Y_t)}(\phi_t(y))}{p_{Y_t}(y)} \left( \frac{p_{\sqrt{\alpha_t} X_{t-1}}(\phi_t(y))}{p_{X_t}(y)} \right)^{-1} \\ &= 1 + O \left( d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T + d^6 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3 \log^3 T \right). \end{aligned}$$

Replacing  $y$  with  $Y_t$  in the above display, using the fact that  $Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \phi_t(Y_t)$ , and invoking the relation (44), we immediately arrive at

$$\frac{p_{t-1}(Y_{t-1})}{q_{t-1}(Y_{t-1})} = \left\{ 1 + O \left( d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T + d^6 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3 \log^3 T \right) \right\} \cdot \frac{p_t(Y_t)}{q_t(Y_t)}$$

with probability exceeding  $1 - (T - t) \exp(-c_3 d \log T)$ . This concludes the proof of Claim (46) via standard induction arguments.

*Proof of property (49).* First, for any  $y$  obeying  $y \in \mathcal{E}_t$  and  $p_t(y) \geq \exp(-c_6 d \log T)$ , one has

$$q_t(y) \geq \left( 1 - \left| \frac{q_t(y)}{p_t(y)} - 1 \right| \right) p_t(y) \geq \frac{1}{2} p_t(y) \geq \exp(-c_6 d \log T)$$

and, as a result,

$$\left\{ y \mid p_t(y) \geq 2 \exp(-c_6 d \log T) \right\} \cap \mathcal{E}_t \subseteq \left\{ y \mid q_t(y) \geq \exp(-c_6 d \log T) \right\}. \quad \text{eq:y-set-pt-cap-E} \quad (50)$$

We can then deduce that

$$\begin{aligned} \mathbb{P}(q_t(Y_t) \geq \exp(-c_6 d \log T)) &\geq \mathbb{P}(p_t(Y_t) \geq 2 \exp(-c_6 d \log T) \text{ and } Y_t \in \mathcal{E}_t) \\ &\geq \mathbb{P}(Y_t \in \{y : p_t(y) \geq 2 \exp(-c_6 d \log T)\} \cap \mathcal{E}_t \cap \{y : \|y\|_2 \leq T^{c_y}\}) \\ &\geq 1 - \mathbb{P}\{Y_t \notin \mathcal{E}_t\} - \mathbb{P}(Y_t \in \{y : p_t(y) < 2 \exp(-c_6 d \log T)\} \cap \{y : \|y\|_2 \leq T^{c_y}\}) - \mathbb{P}\{\|Y_t\|_2 > T^{c_y}\} \end{aligned}$$

$$\begin{aligned}
&\geq 1 - (T - t - 1) \exp(-c_3 d \log T) - \int_{y: \|y\|_2 \leq T^{c_y}} 2 \exp(-c_6 d \log T) dy - \mathbb{P}\{\|Y_T\|_2 > T^{c_y/2}\} \\
&\geq 1 - (T - t) \exp(-c_3 d \log T).
\end{aligned}$$

Here, the last line makes use of the fact  $y_T \sim \mathcal{N}(0, I_d)$  and holds as long as  $c_6$  is large enough; regarding the penultimate line, it suffices to recognize that (cf. (42) and (21))

$$\|\phi_t(x)\|_2 \leq \left(1 - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)}\right) \|x\|_2 + \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} R$$

and hence  $\|y_t\|_2 \lesssim \text{poly}(T)$  unless  $\|y_T\|_2 \geq T^{c_y/2}$  (assuming that  $c_y$  is large).  $\square$

### A.3 Proof of Theorem 2

Before continuing, we introduce the following additional notation

$$\varphi_t(x) := x + \frac{\lambda_t(1 - \alpha_t)}{2} s_t(x) = x - \frac{\lambda_t(1 - \alpha_t)}{2(1 - \bar{\alpha}_t)} \int_{x_0} (x - \sqrt{\bar{\alpha}_t} x_0) p_{X_0|X_t}(x_0|x) dx_0, \quad (51)$$

where the second identity follows from (36). Here, it is worth noting that  $|\lambda_t - 1| \lesssim d^{3/2} \frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}$ . With this notation in place, we can write

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \varphi_t(Y_t). \quad \text{eq:Yt-recursion-accelerate-proof (52)}$$

As it turns out, Theorem 2 can be established in a very similar way as in the proof of Theorem 1. The only step that needs to be changed is to replace (45) with (53) in the following lemma. The proof of this lemma can be found in Appendix A.7.

**Lemma 3.** *Consider every  $x \in \mathbb{R}$  such that  $-\log p_{X_t}(x) \lesssim d \log T$ . Suppose  $\frac{d^2 \log T(1 - \alpha_t)}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ , it satisfies that* lem:main-ODE-fast

$$\begin{aligned}
\frac{p_{\sqrt{\alpha_t} X_{t-1}}(\varphi_t(x))}{p_{X_t}(x)} &= 1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{d(d - 2)(1 - \alpha_t)^2}{8(\alpha_t - \bar{\alpha}_t)^2} + O\left(d^3 \log^3 T \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right) \\
&\quad - \frac{(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} \left(A_t - \lambda_t B_t - \frac{1 - \alpha_t}{4(\alpha_t - \bar{\alpha}_t)} (-B_t + C_t + D_t - 2E_t)\right) \quad \text{eq:xt-higher (53a)}
\end{aligned}$$

and

$$\begin{aligned}
\frac{p_{\varphi_t(X)}(\varphi_t(x))}{p_X(x)} &= 1 + \frac{\lambda_t d(1 - \alpha_t)}{2(1 - \bar{\alpha}_t)} + \frac{\lambda_t(1 - \alpha_t)(B_t - A_t)}{2(1 - \bar{\alpha}_t)} + O\left(d^6 \log^3 T \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right) \\
&\quad + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)^2} \left[\left(d + B_t - A_t\right)^2 + d + 2(B_t - A_t) + F_t\right] - \frac{(1 - \alpha_t)}{2} s_t^\top \nabla \lambda_t, \quad \text{eq:yt-higher (53b)}
\end{aligned}$$

where  $A_t, \dots, F_t$  are defined in (25).

With Lemma 3 in place, we are able to prove, in the same way as in the proof of Theorem 1, that

$$\mathbb{P}(Y_t \in \mathcal{E}_t) \geq 1 - (T - t + 1) \exp(-\Omega(d \log T)) \quad \text{eq:claim (54)}$$

for every  $t \geq 1$ , where the set  $\mathcal{E}_t$  is defined as

$$\mathcal{E}_t := \left\{y : \left| \frac{p_{Y_t}(y)}{p_{X_t}(y)} - 1 \right| \lesssim (T - t + 1) \frac{d^6 \log^6 T}{T^3} \right\}.$$

This implies the desired result in Theorem 2.

## A.4 Proof of properties (41)

sec:proof-properties-alpha

**Proof of property (41a).** We start by proving (41a). Let  $\tau$  be an integer obeying

$$\beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^\tau \leq 1 < \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{\tau+1}, \quad \text{eq:defn-tau-proof-alpha} \quad (55)$$

and we divide into two cases based on  $\tau$ .

- Consider any  $t$  satisfying  $t \leq \tau$ . In this case, it suffices to prove that

$$1 - \bar{\alpha}_{t-1} \geq \frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^t. \quad \text{eq:induction-proof-alpha} \quad (56)$$

Clearly, if (56) is valid, then any  $t \leq \tau$  obeys

$$\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}} = \frac{\beta_t}{1 - \bar{\alpha}_{t-1}} \leq \frac{\frac{c_1 \log T}{T} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^t}{\frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^t} = \frac{3c_1 \log T}{T}$$

as claimed. Towards proving (56), first note that the base case with  $t = 2$  holds true trivially since  $1 - \bar{\alpha}_1 = 1 - \alpha_1 = \beta_1 \geq \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^2 / 3$ . Next, let  $t_0 > 2$  be the first time that Condition (56) fails to hold and suppose that  $t_0 \leq \tau$ . It then follows that

$$1 - \bar{\alpha}_{t_0-2} = 1 - \frac{\bar{\alpha}_{t_0-1}}{\alpha_{t_0-1}} \leq 1 - \bar{\alpha}_{t_0-1} < \frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0} \leq \frac{1}{2} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1} < \frac{1}{2}, \quad (57)$$

where the last inequality result from (55) and the assumption  $t_0 \leq \tau$ . This taken together with the assumptions (56) and  $t_0 \leq \tau$  implies that

$$\frac{(1 - \alpha_{t_0-1})\bar{\alpha}_{t_0-1}}{1 - \bar{\alpha}_{t_0-2}} \geq \frac{\frac{c_1 \log T}{T} \beta_1 \min \left\{ \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1}, 1 \right\} \cdot \left(1 - \frac{1}{2}\right)}{\frac{1}{2} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1}} = \frac{\frac{c_1 \log T}{T} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1}}{\beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1}} = \frac{c_1 \log T}{T}.$$

As a result, we can further derive

$$\begin{aligned} 1 - \bar{\alpha}_{t_0-1} &= 1 - \alpha_{t_0-1} \bar{\alpha}_{t_0-2} = 1 - \bar{\alpha}_{t_0-2} + (1 - \alpha_{t_0-1})\bar{\alpha}_{t_0-2} \\ &= \left(1 + \frac{(1 - \alpha_{t_0-1})\bar{\alpha}_{t_0-2}}{1 - \bar{\alpha}_{t_0-2}}\right)(1 - \bar{\alpha}_{t_0-2}) \\ &\geq \left(1 + \frac{c_1 \log T}{T}\right)(1 - \bar{\alpha}_{t_0-2}) \geq \left(1 + \frac{c_1 \log T}{T}\right) \cdot \left\{ \frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0-1} \right\} \\ &= \frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{t_0}, \end{aligned}$$

where the penultimate line holds since (56) is first violated at  $t = t_0$ ; this, however, contradicts with the definition of  $t_0$ . Consequently, one must have  $t_0 > \tau$ , meaning that (56) holds for all  $t \leq \tau$ .

- We then turn attention to those  $t$  obeying  $t > \tau$ . In this case, it suffices to make the observation that

$$1 - \bar{\alpha}_{t-1} \geq 1 - \bar{\alpha}_{\tau-1} \geq \frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^\tau = \frac{\frac{1}{3} \beta_1 \left(1 + \frac{c_1 \log T}{T}\right)^{\tau+1}}{1 + \frac{c_1 \log T}{T}} \geq \frac{1}{4}, \quad (58)$$

where the second and the third inequalities come from (56). Therefore, one obtains

$$\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}} \leq \frac{\frac{c_1 \log T}{T}}{1/4} \leq \frac{4c_1 \log T}{T}.$$

**Proof of property (41b).** Turning attention to the second claim (41b), we note that for any  $t$  obeying  $t \geq \frac{T}{2} \gtrsim \frac{T}{\log T}$ , one has

$$1 - \alpha_t = \frac{c_1 \log T}{T} \min \left\{ \beta_1 \left( 1 + \frac{c_1 \log T}{T} \right)^t, 1 \right\} = \frac{c_1 \log T}{T}.$$

This in turn allows one to deduce that

$$\bar{\alpha}_T \leq \prod_{t: t \geq T/2} \alpha_t \leq \left( 1 - \frac{c_1 \log T}{T} \right)^{T/2} \leq \frac{1}{T^{c_2}}$$

for an arbitrarily large constant  $c_2 > 0$ .

**Proof of property (41c).** Finally, it is easily seen from the Taylor expansion that the learning rates  $\{\alpha_t\}$  satisfy

$$\begin{aligned} \left( \frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} &= \left( 1 + \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} \\ &= 1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{d(d-2)(1 - \alpha_t)^2}{8(\alpha_t - \bar{\alpha}_t)^2} + O\left(d^3 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3\right), \end{aligned}$$

provided that  $\frac{d(1 - \alpha_t)}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ .

## A.5 Proof of Lemma 1

To establish this lemma, we first make the following claim, whose proof is deferred to the end of this subsection. sec:proof-lem:x0  
eq:claim-123

**Claim 1.** Consider any  $c_5 \geq 5$ . There exist some  $x_0 \in \mathbb{R}$  such that eqn:lemma-x0-claim

$$\begin{aligned} \|\sqrt{\bar{\alpha}_t}x_0 - y\|_2 &\leq c_5 \sqrt{d(1 - \bar{\alpha}_t) \log T} \quad \text{and} \quad \text{eqn:lemma-x0-claim-1} \\ \mathbb{P}(\|X_0 - x_0\|_2 \leq \epsilon) &\geq \left( \frac{\epsilon}{2T^{c_6 + c_R}} \right)^d \quad \text{with} \quad \epsilon = \frac{1}{T^{c_0/2}} \quad \text{eqn:lemma-x0-claim-2} \end{aligned} \tag{59a}$$

hold simultaneously, where  $c_0$  is defined in (20) and  $c_R$  is defined in (37).

With the above claim in place, we are ready to prove Lemma 1. Let us look at a set:

$$\mathcal{E} := \left\{ x : \sqrt{\bar{\alpha}_t} \|x - x_0\|_2 > 4c_5 \sqrt{d(1 - \bar{\alpha}_t) \log T} \right\},$$

where  $c_5 \geq 5$  (see Claim 1). Combining this with (59a) results in

$$\mathbb{P}(\|\sqrt{\bar{\alpha}_t}X - y\|_2 > 5c_5 \sqrt{d(1 - \bar{\alpha}_t) \log T}) \leq \mathbb{P}(X \in \mathcal{E}). \quad \text{eq:UB-P-set-E} \tag{60}$$

Consequently, everything boils down to bounding  $\mathbb{P}(X \in \mathcal{E})$ . Towards this, we first invoke the Bayes rule  $p_{X_0|X_t}(x|y) \propto p_{X_0}(x)p_{X_t|X_0}(y|x)$  to derive

$$\begin{aligned} \mathbb{P}(X_0 \in \mathcal{E} | X_t = y) &= \frac{\int_{x \in \mathcal{E}} p_{X_0}(x) p_{X_t|X_0}(y|x) dx}{\int_x p_{X_0}(x) p_{X_t|X_0}(y|x) dx} \\ &\leq \frac{\int_{x \in \mathcal{E}} p_{X_0}(x) p_{X_t|X_0}(y|x) dx}{\int_{x: \|x - x_0\|_2 \leq \epsilon} p_{X_0}(x) p_{X_t|X_0}(y|x) dx} \\ &\leq \frac{\sup_{x \in \mathcal{E}} p_{X_t|X_0}(y|x)}{\inf_{x: \|x - x_0\|_2 \leq \epsilon} p_{X_t|X_0}(y|x)} \cdot \frac{\mathbb{P}(X_0 \in \mathcal{E})}{\mathbb{P}(\|X_0 - x_0\|_2 \leq \epsilon)}. \quad \text{eq:UB1-P-X0-Xt} \end{aligned} \tag{61}$$

To further bound this quantity, note that: in view of the definition of  $\mathcal{E}$  and expression (59a), one has

$$\begin{aligned} \sup_{x \in \mathcal{E}} p_{X_t | X_0}(y | x) &= \sup_{x: \|\sqrt{\bar{\alpha}_t}x - \sqrt{\bar{\alpha}_t}x_0\|_2 > 4c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_t | X_0}(y | x) \\ &\leq \sup_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 > 3c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_t | X_0}(y | x) \\ &\leq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \exp\left(-\frac{9c_5^2 d \log T}{2}\right) \end{aligned}$$

and

$$\begin{aligned} \inf_{x: \|x-x_0\|_2 \leq \epsilon} p_{X_t | X_0}(y | x) &\geq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \inf_{x: \|x-x_0\|_2 \leq \epsilon} \exp\left(-\frac{\|y - \sqrt{\bar{\alpha}_t}x\|_2^2}{2(1-\bar{\alpha}_t)}\right) \\ &\geq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \inf_{x: \|x-x_0\|_2 \leq \epsilon} \exp\left(-\frac{\|y - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{1-\bar{\alpha}_t} - \frac{\|\sqrt{\bar{\alpha}_t}x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{1-\bar{\alpha}_t}\right) \\ &\geq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|y - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{1-\bar{\alpha}_t} - \frac{\epsilon^2}{1-\bar{\alpha}_t}\right) \\ &\geq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \exp\left(-c_5^2 d \log T - \frac{1}{T^{c_0}} \frac{1}{1-\bar{\alpha}_t}\right) \\ &\geq \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \exp(-2c_5^2 d \log T), \end{aligned}$$

where the penultimate line relies on (59a), and the last line holds true since  $1-\bar{\alpha}_t \geq 1-\alpha_1 = 1/T^{c_0}$  (see (20)). Substitution of the above two displays into (61), we arrive at

$$\begin{aligned} \mathbb{P}(X_0 \in \mathcal{E} | X_t = y) &\leq \exp(-2.5c_5^2 d \log T) \cdot \frac{1}{\mathbb{P}(\|X_0 - x_0\|_2 \leq \epsilon)} \\ &\leq \exp(-2.5c_5^2 d \log T) \cdot \left(2T^{c_6+c_0/2+c_R}\right)^d \\ &\leq \exp\left(-(2c_5^2 - c_6 - c_0/2 - c_R)d \log T\right). \end{aligned} \tag{62}$$

Substituting this into (60) and recalling the distribution (39) of  $X$ , we arrive at

$$\mathbb{P}\left(\|\sqrt{\bar{\alpha}_t}X - y\|_2 > 5c_5 \sqrt{d \log T(1-\bar{\alpha}_t)}\right) \leq \exp\left(-(2c_5^2 - c_6 - c_0/2 - c_R)d \log T\right).$$

This concludes the proof of Lemma 1 for sufficiently large  $c_5$ , as long as Claim 1 can be justified.

*Proof of Claim 1.* We prove this claim by contradiction. Specifically, suppose instead that: for every  $x$  obeying  $\|\sqrt{\bar{\alpha}_t}x - y\|_2 \leq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}$ , we have

$$\mathbb{P}(\|X_0 - x\|_2 \leq \epsilon) \leq \left(\frac{\epsilon}{2T^{c_6}R}\right)^d \quad \text{with } \epsilon = \frac{1}{T^{c_0/2}}. \quad \text{eq:contradiction-x-y} \tag{63}$$

Clearly, the choice of  $\epsilon$  ensures that  $\epsilon < \frac{1}{2} \sqrt{d(1-\bar{\alpha}_t) \log T}$ . In the following, we would like to show that this assumption leads to contradiction.

First of all, let us look at  $p_{X_t}$ , which obeys

$$\begin{aligned} p_{X_t}(y) &= \int_x p_{X_0}(x) p_{X_t | X_0}(y | x) dx \\ &= \int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 \geq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_0}(x) p_{X_t | X_0}(y | x) dx \end{aligned}$$

$$\begin{aligned}
& + \int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 < c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_0}(x) p_{X_t | X_0}(y | x) dx \\
& \leq \int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 \geq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_t | X_0}(y | x) dx + \int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 < c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_0}(x) dx. \tag{64}
\end{aligned}$$

To further control (64), we make two observations:

- 1) The first term on the right-hand side of (64) can be bounded by

$$\begin{aligned}
& \int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 \geq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_t | X_0}(y | x) dx \\
& = \int_{z: \|z\|_2 \geq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} \frac{1}{(2\pi(1-\bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|z\|_2^2}{2(1-\bar{\alpha}_t)}\right) dz \\
& \leq \frac{1}{2} \exp(-c_6 d \log T), \tag{65}
\end{aligned}$$

for some constant  $c_6 > 0$ , provided that  $c_5$  is sufficiently large. Here, we have used  $X_t \stackrel{(i)}{=} \sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}W$  with  $W \sim \mathcal{N}(0, I_d)$  as well as standard properties about Gaussian distributions.

- 2) Regarding the second term on the right-hand side of (64), let us construct an epsilon-net  $\mathcal{N}_\epsilon = \{z_i\}$  for the set

$$\{x : \|\sqrt{\bar{\alpha}_t}x - y\|_2 \leq c_5 \sqrt{d(1-\bar{\alpha}_t) \log T} \text{ and } \|x\|_2 \leq R\},$$

so that for each  $x$  in this set, one can find a vector  $z_i \in \mathcal{N}_\epsilon$  such that  $\|x - z_i\|_2 \leq \epsilon$ . Define  $\mathcal{B}_i := \{x \mid \|x - z_i\|_2 \leq \epsilon\}$  for each  $z_i \in \mathcal{N}_\epsilon$ . Armed with these sets, we can derive

$$\begin{aligned}
\int_{x: \|\sqrt{\bar{\alpha}_t}x - y\|_2 < c_5 \sqrt{d(1-\bar{\alpha}_t) \log T}} p_{X_0}(x) dx & \leq \sum_{i=1}^{|\mathcal{N}_\epsilon|} \mathbb{P}(X_0 \in \mathcal{B}_i) \\
& \leq \left(\frac{\epsilon}{2T^{c_6} R}\right)^d \left(\frac{R}{\epsilon}\right)^d \\
& \leq \frac{1}{2} \exp(-c_6 d \log T),
\end{aligned}$$

where the penultimate step comes from the assumption (63).

The above results taken collectively lead to

$$p_{X_t}(y) \leq \exp(-c_6 d \log T), \tag{66}$$

thus contradicting the assumption (38). This in turn validates this claim.  $\square$

## A.6 Proof of Lemma 2

sec:proof-lem:main-ODE

### A.6.1 Proof of relation (45a)

Recall the definition of  $\phi_t$  in (42), and introduce the following vector:

$$u := x - \phi_t(x) = \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} \int_{x_0} (x - \sqrt{\bar{\alpha}_t}x_0) p_{X_0 | X_t}(x_0 | x) dx_0. \tag{67}$$

The proof consists of the following steps.



**Step 1: decomposing**  $p_{\sqrt{\alpha_t}X_{t-1}}(\phi_t(x))/p_{X_t}(x)$ . Recognizing that

$$X_t \stackrel{d}{=} \sqrt{\alpha_t}X_0 + \sqrt{1 - \alpha_t}W \quad \text{with } W \sim \mathcal{N}(0, I_d) \quad \text{eq:Xt-dist-proof1 (68)}$$

and making use of the Bayes rule, we can express the conditional distribution  $p_{X_0|X_t}(\phi_t(x))$  as

$$p_{X_0|X_t}(x_0|x) = \frac{p_{X_0}(x_0)}{p_{X_t}(x)} p_{X_t|X_0}(x|x_0) = \frac{p_{X_0}(x_0)}{p_{X_t}(x)} \cdot \frac{1}{(2\pi(1 - \alpha_t))^{d/2}} \exp\left(-\frac{\|x - \sqrt{\alpha_t}x_0\|_2^2}{2(1 - \alpha_t)}\right). \quad \text{eqn: bayes (69)}$$

In turn, this taken together with (68) allows one to rewrite  $p_{\sqrt{\alpha_t}X_{t-1}}$  such that: for any set  $\mathcal{E}$ ,

$$\begin{aligned} \frac{p_{\sqrt{\alpha_t}X_{t-1}}(\phi_t(x))}{p_{X_t}(x)} &\stackrel{(i)}{=} \frac{1}{p_{X_t}(x)} \int_{x_0} p_{X_0}(x_0) \frac{1}{(2\pi(\alpha_t - \bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|\phi_t(x) - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)}\right) dx_0 \\ &\stackrel{(ii)}{=} \frac{1}{p_{X_t}(x)} \int_{x_0} p_{X_0}(x_0) \frac{1}{(2\pi(\alpha_t - \bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|x - \sqrt{\alpha_t}x_0\|_2^2}{2(1 - \bar{\alpha}_t)}\right) \\ &\quad \cdot \exp\left(-\frac{(1 - \alpha_t)\|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)}\right) dx_0 \\ &\stackrel{(iii)}{=} \left(\frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t}\right)^{d/2} \cdot \int_{x_0} p_{X_0|X_t}(x_0|x) \cdot \\ &\quad \exp\left(-\frac{(1 - \alpha_t)\|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)}\right) dx_0 \\ &\stackrel{(iv)}{=} \left\{1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + O\left(d^2\left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^2\right)\right\} \cdot \\ &\quad \int_{x_0} p_{X_0|X_t}(x_0|x) \exp\left(-\frac{(1 - \alpha_t)\|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)}\right) dx_0. \end{aligned} \quad \text{eqn: fei (70)}$$

Here, identity (i) holds true since  $\sqrt{\alpha_t}X_{t-1} \stackrel{d}{=} \sqrt{\alpha_t}(\sqrt{\bar{\alpha}_{t-1}}X_0 + \sqrt{1 - \bar{\alpha}_{t-1}}W) = \sqrt{\alpha_t}X_0 + \sqrt{\alpha_t - \bar{\alpha}_t}W$  and hence

$$p_{\sqrt{\alpha_t}X_{t-1}}(x) = \int_{x_0} p_{X_0}(x_0) p_{\sqrt{\alpha_t - \bar{\alpha}_t}W}(x - \sqrt{\alpha_t}x_0) dx_0;$$

identity (ii) follows from elementary algebra; relation (iii) is a consequence of the Bayes rule (69); and relation (iv) results from (41c).

**Step 2: controlling the integral in the decomposition (70).** In order to further control the right-hand side of expression (70), we need to evaluate the integral in (70). To this end, we make a few observations.

- To begin with, Lemma 1 tells us that <sup>eqn: BBB</sup>

$$\mathbb{P}\left(\|\sqrt{\alpha_t}X_0 - x\|_2 > 5c_5\sqrt{d(1 - \bar{\alpha}_t)\log T} \mid X_t = x\right) \leq \exp(-c_5^2 d \log T) \quad \text{eqn: brahms (71a)}$$

for any large enough  $c_5$ , provided that  $x$  satisfies  $-\log p_{X_t}(x) \leq 2d \log T$ .

- A little algebra based on this relation allows one to bound  $u$  (cf. (67)) as follows:

$$\begin{aligned} \|u\|_2 &\leq \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} \int_{x_0: \|x - \sqrt{\alpha_t}x_0\|_2 \leq c_5\sqrt{d(1 - \bar{\alpha}_t)\log T}} p_{X_0|X_t}(x_0|x) \|x - \sqrt{\alpha_t}x_0\|_2 dx_0 \\ &\quad + \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} \int_{x_0: \|x - \sqrt{\alpha_t}x_0\|_2 > c_5\sqrt{d(1 - \bar{\alpha}_t)\log T}} p_{X_0|X_t}(x_0|x) \|x - \sqrt{\alpha_t}x_0\|_2 dx_0 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} \cdot c_5 \sqrt{d(1-\bar{\alpha}_t) \log T} + \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} \cdot \int_{c_5}^{\infty} \left( d(1-\bar{\alpha}_t) \log T \right) \tau \exp \left( -\frac{1}{25} \tau^2 d \log T \right) d\tau \\
&\leq \frac{2c_5(1-\alpha_t) \sqrt{d \log T}}{3\sqrt{1-\bar{\alpha}_t}} \leq \frac{2c_5}{3} \sqrt{d(1-\alpha_t) \log T},
\end{aligned} \tag{71b}$$

with the proviso that  $-\log p_{X_t}(x) \leq 2d \log T$ .

Equipped with the above properties, let us define

$$\mathcal{E} := \left\{ x : \|x - \sqrt{\bar{\alpha}_t} x_0\|_2 \leq 5c_5 \sqrt{d(1-\bar{\alpha}_t) \log T} \right\}. \tag{72}$$

For any  $x \in \mathcal{E}$ , the Taylor expansion  $e^{-x} = 1 - x + O(x^2)$  (for all  $|x| < 1$ ) gives

$$\begin{aligned}
&\exp \left( -\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) \\
&= 1 - \frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2}{2(\alpha_t - \bar{\alpha}_t)} + \frac{u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{\alpha_t - \bar{\alpha}_t} + O \left( d^2 \left( \frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T \right) \\
&= 1 - \frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} + \frac{u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{\alpha_t - \bar{\alpha}_t} + O \left( d^2 \left( \frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T \right),
\end{aligned} \tag{73}$$

where the penultimate line invokes (71), and the last line holds true since, according to (71b),

$$\frac{\|u\|_2^2}{|\alpha_t - \bar{\alpha}_t|} \leq \frac{1}{\alpha_t - \bar{\alpha}_t} \cdot \frac{9c_5^2(1-\alpha_t)^2 d \log T}{(1-\bar{\alpha}_t)} \leq \frac{9c_5^2(1-\alpha_t)^2 d \log T}{(\alpha_t - \bar{\alpha}_t)^2}.$$

In contrast, for any  $x \notin \mathcal{E}$ , we invoke the crude bound

$$\begin{aligned}
&-\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2}{2(\alpha_t - \bar{\alpha}_t)} + \frac{2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \\
&\leq -\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2}{2(\alpha_t - \bar{\alpha}_t)} + \frac{\frac{1-\alpha_t}{1-\bar{\alpha}_t} \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 + \frac{1-\bar{\alpha}_t}{1-\alpha_t} \|u\|_2^2}{2(\alpha_t - \bar{\alpha}_t)} \\
&\leq \frac{\left( \frac{1-\bar{\alpha}_t}{1-\alpha_t} - 1 \right) \|u\|_2^2}{2(\alpha_t - \bar{\alpha}_t)} = \frac{\|u\|_2^2}{2(1-\alpha_t)},
\end{aligned}$$

and as a result,

$$\exp \left( -\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) \leq \exp \left( \frac{\|u\|_2^2}{2(1-\alpha_t)} \right) \leq \exp \left( \frac{2c_5^2 d \log T}{9} \right). \tag{74}$$

Combine (73) and (74) to show that: if  $\frac{d(1-\alpha_t) \log T}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ , then one has

$$\begin{aligned}
&\int_{x_0} p_{X_0 | X_t}(x_0 | x) \exp \left( -\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) dx_0 \\
&= \left( \int_{x_0 \in \mathcal{E}} + \int_{x_0 \notin \mathcal{E}} \right) p_{X_0 | X_t}(x_0 | x) \exp \left( -\frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) dx_0 \\
&= \left( 1 + O \left( d^2 \left( \frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T \right) \right) \int_{x_0} p_{X_0 | X_t}(x_0 | x) \left( 1 - \frac{(1-\alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1-\bar{\alpha}_t)} + \frac{u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{\alpha_t - \bar{\alpha}_t} \right) dx_0 \\
&\quad + O \left( \exp \left( \frac{2c_5^2 d \log T}{9} \right) \int_{x_0 \notin \mathcal{E}} p_{X_0 | X_t}(x_0 | x) dx_0 \right)
\end{aligned}$$

$$\begin{aligned}
&= 1 - \frac{(1 - \alpha_t) \left( \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\alpha_t} x_0\|_2^2 dx_0 - \left\| \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\alpha_t} x_0) dx_0 \right\|_2^2 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} \\
&\quad + O\left(d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T\right) + O\left(\exp\left(-\frac{7c_5^2 d \log T}{9}\right)\right) \\
&= 1 - \frac{(1 - \alpha_t) \left( \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\alpha_t} x_0\|_2^2 dx_0 - \left\| \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\alpha_t} x_0) dx_0 \right\|_2^2 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} \\
&\quad + O\left(d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T\right), \tag{eq:exp-UB-135702}
\end{aligned}$$

where the last relation makes use of the definition (67) of  $u$ .

**Step 3: putting everything together.** Substitution of (75) into (70) yields

$$\begin{aligned}
\frac{p_{\sqrt{\alpha_t} X_{t-1}}(\phi_t(x))}{p_{X_t}(x)} &= 1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + O\left(d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T\right) - \\
&\quad \frac{(1 - \alpha_t) \left( \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\alpha_t} x_0\|_2^2 dx_0 - \left\| \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\alpha_t} x_0) dx_0 \right\|_2^2 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)}
\end{aligned}$$

as claimed.

#### A.6.2 Proof of relation (45b)

To understand the density ratio  $p_{\phi_t(X)}(\phi_t(x))/p_X(x)$ , we make note of the transformation

$$p_{\phi_t(X)}(\phi_t(x)) = \det\left(\frac{\partial \phi_t(x)}{\partial x}\right)^{-1} p_X(x), \tag{76}$$

where  $\frac{\partial \phi_t(x)}{\partial x}$  denotes the Jacobian matrix. It thus suffices to control the quantity  $\det\left(\frac{\partial \phi_t(x)}{\partial x}\right)^{-1}$ .

To begin with, elementary calculations of the derivative — in conjunction with the fact that  $X_t \stackrel{d}{=} \sqrt{\alpha_t} X_0 + \sqrt{1 - \alpha_t} W$  with  $W \sim \mathcal{N}(0, I_d)$  — yield that: for any  $1 \leq i, j \leq d$ ,

$$\begin{aligned}
&\frac{\partial \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x_i - \sqrt{\alpha_t} x_{0,i}) dx_0}{\partial x_j} \\
&= \mathbb{1}\{i = j\} + \frac{1}{1 - \bar{\alpha}_t} \left\{ \left( \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x_i - \sqrt{\alpha_t} x_{0,i}) dx_0 \right) \left( \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x_j - \sqrt{\alpha_t} x_{0,j}) dx_0 \right) \right. \\
&\quad \left. - \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x_i - \sqrt{\alpha_t} x_{0,i}) (x_j - \sqrt{\alpha_t} x_{0,j}) dx_0 \right\} \\
&= \mathbb{1}\{i = j\} + \frac{1}{1 - \bar{\alpha}_t} \left\{ \mathbb{E}[g_i(X_0) | X_t = x] \mathbb{E}[g_j(X_0) | X_t = x] - \mathbb{E}[g_i(X_0) g_j(X_0) | X_t = x] \right\}, \tag{eqn:derivative}
\end{aligned}$$

where  $g_i(x_0) := x_i - \sqrt{\alpha_t} x_{0,i}$ . Combining (77) with the definition of  $\phi_t$  (cf. (42)), one can readily see that

$$\begin{aligned}
\text{Tr}\left(I - \frac{\partial \phi_t(x)}{\partial x}\right) &= \frac{d(1 - \alpha_t)}{2(1 - \bar{\alpha}_t)} + \\
&\quad \frac{(1 - \alpha_t) \left( \left\| \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\alpha_t} x_0) dx_0 \right\|_2^2 - \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\alpha_t} x_0\|_2^2 dx_0 \right)}{2(1 - \bar{\alpha}_t)^2}. \tag{78a}
\end{aligned}$$

Moreover, define a matrix  $\mathbf{B} = [B_{i,j}]_{1 \leq i, j \leq d}$  with

$$B_{i,j} := \mathbb{E}[g_i(X_0) g_j(X_0) | X_t = x] - \mathbb{E}[g_i(X_0) | X_t = x] \mathbb{E}[g_j(X_0) | X_t = x],$$

which clearly satisfies (using Jensen's inequality)

$$\begin{aligned}
\|\mathbf{B}\|_{\mathbb{F}} &= \left\| \mathbb{E} \left[ \begin{bmatrix} g_1(X_0) \\ \vdots \\ g_d(X_0) \end{bmatrix} \begin{bmatrix} g_1(X_0) & \cdots & g_d(X_0) \end{bmatrix} \mid X_t = x \right] \right\|_{\mathbb{F}} \\
&\leq \mathbb{E} \left[ \left\| \begin{bmatrix} g_1(X_0) \\ \vdots \\ g_d(X_0) \end{bmatrix} \begin{bmatrix} g_1(X_0) & \cdots & g_d(X_0) \end{bmatrix} \mid X_t = x \right\|_{\mathbb{F}}^2 \right] = \mathbb{E} \left[ \sum_{i=1}^d (g_i(X_0))^2 \mid X_t = x \right] \\
&= \int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0.
\end{aligned}$$

Taking (77) and (42) together and using the above bound on  $\|\mathbf{B}\|_{\mathbb{F}}$  also reveal that

$$\begin{aligned}
\left\| \frac{\partial \phi_t(x)}{\partial x} - I \right\| &\leq \left\| \frac{\partial \phi_t(x)}{\partial x} - I \right\|_{\mathbb{F}} \lesssim \frac{1 - \alpha_t}{1 - \bar{\alpha}_t} \left( \sqrt{d} + \frac{1}{1 - \bar{\alpha}_t} \|\mathbf{B}\|_{\mathbb{F}} \right) \\
&\lesssim \frac{1 - \alpha_t}{1 - \bar{\alpha}_t} \left( \sqrt{d} + \frac{\int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0}{1 - \bar{\alpha}_t} \right) \\
&\lesssim \frac{d(1 - \alpha_t) \log T}{1 - \bar{\alpha}_t},
\end{aligned} \tag{78b}$$

where the last line invokes the relation (71a). Additionally, the Taylor expansion guarantees that

$$\det(I + A) = 1 + \text{Tr}(A) + O((\text{Tr}(A))^2 + \|A\|_{\mathbb{F}}^2 + d^3 \|A\|^3)$$

as long as  $d\|A\| \lesssim 1$ . The above properties taken collectively allow us to demonstrate that

$$\begin{aligned}
\frac{p_{\phi_t(X)}(\phi_t(x))}{p_X(x)} &= \det \left( \frac{\partial \phi_t(x)}{\partial x} \right)^{-1} \\
&= 1 - \text{Tr} \left( \frac{\partial \phi_t(x)}{\partial x} - I \right) + O \left( d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T + d^6 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3 \log^3 T \right) \\
&= 1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{(1 - \alpha_t) \left( \left\| \int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x) (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right\|_2^2 - \int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0 \right)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} \\
&\quad + O \left( d^2 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \log^2 T + d^6 \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3 \log^3 T \right),
\end{aligned} \tag{79}$$

with the proviso that  $\frac{d^2(1 - \alpha_t) \log T}{\alpha_t - \bar{\alpha}_t} \lesssim 1$ .

## A.7 Proof of Lemma 3

sec:proof-lem:main-ODE-fast

The proof of Lemma 3 is derived in the same way as the proof of Lemma 2.

### A.7.1 Proof of relation (53a)

First let us recall that

$$\begin{aligned}
p_{\sqrt{\bar{\alpha}_t} X_{t-1}}(\varphi_t(x)) &= p_{X_t}(x) \left( \frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} \\
&\quad \cdot \int_{x_0} dx_0 p_{X_0 \mid X_t}(x_0 \mid x) \exp \left( - \frac{(1 - \alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top (x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right),
\end{aligned} \tag{80}$$

where  $u = x - \varphi_t(x) = \frac{\lambda_t(x)(1-\alpha_t)}{2(1-\bar{\alpha}_t)} \int_{x_0} p_{X_0|X_t}(x_0|x)(x - \sqrt{\bar{\alpha}_t}x_0)dx_0$ . Moreover, according to Lemma 1, we have  $\|x - \sqrt{\bar{\alpha}_t}x_0\| \lesssim \sqrt{d \log T(1 - \bar{\alpha}_t)}$ , and then

$$\begin{aligned} & \int_{x_0} p_{X_0|X_t}(x_0|x) \exp\left(-\frac{(1-\alpha_t)\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)}\right) dx_0 \\ &= \left(1 + O\left(d^3 \log^3 T \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right)\right) \int_{x_0} p_{X_0|X_t}(x_0|x) \left(1 - \frac{(1-\alpha_t)\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)}\right. \\ &\quad \left. + \frac{1}{2} \left(\frac{(1-\alpha_t)\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{u^\top(x - \sqrt{\bar{\alpha}_t}x_0)}{\alpha_t - \bar{\alpha}_t}\right)^2\right) dx_0 \\ &= 1 - \frac{(1-\alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} \left(A_t - \lambda_t B_t - \frac{1-\alpha_t}{4(\alpha_t - \bar{\alpha}_t)} [-B_t + C_t + D_t - 2E_t]\right) + O\left(d^3 \log^3 T \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right), \end{aligned} \quad (81)$$

where in the last step, we plug in the definitions in Assumption 2. Therefore, the relation (53a) can be easily seen by recognizing the fact that

$$\left(\frac{1-\bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t}\right)^{d/2} = 1 + \frac{d(1-\alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{d(d-2)(1-\alpha_t)^2}{4(\alpha_t - \bar{\alpha}_t)^2} + O\left(d^3 \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right).$$

#### A.7.2 Proof of relation (53b)

Recall expression (77) that

$$\begin{aligned} & \frac{\partial \int_{x_0} p_{X_0|X_t}(x_0|x)(x_i - \sqrt{\bar{\alpha}_t}x_{0,i})dx_0}{\partial x_j} \\ &= \mathbb{1}\{i=j\} + \frac{1}{1-\bar{\alpha}_t} \left\{ \left( \int_{x_0} p_{X_0|X_t}(x_0|x)(x_i - \sqrt{\bar{\alpha}_t}x_{0,i})dx_0 \right) \left( \int_{x_0} p_{X_0|X_t}(x_0|x)(x_j - \sqrt{\bar{\alpha}_t}x_{0,j})dx_0 \right) \right. \\ &\quad \left. - \int_{x_0} p_{X_0|X_t}(x_0|x)(x_i - \sqrt{\bar{\alpha}_t}x_{0,i})(x_j - \sqrt{\bar{\alpha}_t}x_{0,j})dx_0 \right\}, \end{aligned}$$

which implies that

$$\text{Tr}\left(\frac{\partial \phi_t(x)}{\partial x} - I\right) = -\frac{\lambda_t(1-\alpha_t)(d+B_t-A_t)}{2(1-\bar{\alpha}_t)}, \quad (82a)$$

$$\left\| \frac{\partial \phi_t(x)}{\partial x} - I \right\|_F^2 = \frac{\lambda_t^2(1-\alpha_t)^2}{4(1-\bar{\alpha}_t)^2} \left( d + 2(B_t - A_t) + \sum F_{ij}^2 \right), \quad (82b)$$

$$\left\| \frac{\partial \phi_t(x)}{\partial x} - I \right\|_{\text{op}} \lesssim \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_t}. \quad (82c)$$

Here, we denote

$$\begin{aligned} F_{ij} &:= \frac{1}{1-\bar{\alpha}_t} \left( \int_{x_0} p_{X_0|X_t}(x_0|x)(x_i - \sqrt{\bar{\alpha}_t}x_{0,i})dx_0 \right) \left( \int_{x_0} p_{X_0|X_t}(x_0|x)(x_j - \sqrt{\bar{\alpha}_t}x_{0,j})dx_0 \right) \\ &\quad - \frac{1}{1-\bar{\alpha}_t} \int_{x_0} p_{X_0|X_t}(x_0|x)(x_i - \sqrt{\bar{\alpha}_t}x_{0,i})(x_j - \sqrt{\bar{\alpha}_t}x_{0,j})dx_0. \end{aligned} \quad (83)$$

Now in view of the basic relation

$$\det(I+A)^{-1} = 1 - \text{Tr}(A) + \frac{1}{2} [\text{Tr}(A)^2 + \|A\|_F^2] + O(d^3 \|A\|_{\text{op}}^3), \quad (84)$$

for any symmetric matrix  $A$ , we can derive

$$p_{\phi_t(X)}(\phi_t(x)) = \det\left(\frac{\partial \phi_t(x)}{\partial x}\right)^{-1} p_X(x)$$

$$\begin{aligned}
&= \left(1 - \text{Tr}\left(\frac{\partial\phi_t(x)}{\partial x} - I\right) + \frac{1}{2}\left[\text{Tr}\left(\frac{\partial\phi_t(x)}{\partial x} - I\right)^2 + \left\|\frac{\partial\phi_t(x)}{\partial x} - I\right\|_F^2\right]\right. \\
&\quad \left.+ \frac{(1-\alpha_t)}{2(1-\bar{\alpha}_t)} \nabla\lambda_t(x)^\top \int_{x_0} p_{X_0|X_t}(x_0|x)(x - \sqrt{\bar{\alpha}_t}x_0)dx_0 + O\left(d^6 \log^3 T \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right)\right) p_X(x) \\
&= \left(1 + \frac{\lambda_t d(1-\alpha_t)}{2(1-\bar{\alpha}_t)} + \frac{\lambda_t(1-\alpha_t)(B_t - A_t)}{2(1-\bar{\alpha}_t)} + \frac{(1-\alpha_t)^2}{8(1-\bar{\alpha}_t)^2} \left[(d + B_t - A_t)^2 + d + 2(B_t - A_t) + \sum F_{ij}^2\right]\right. \\
&\quad \left.+ \frac{(1-\alpha_t)}{2(1-\bar{\alpha}_t)} \nabla\lambda_t(x)^\top \int_{x_0} p_{X_0|X_t}(x_0|x)(x - \sqrt{\bar{\alpha}_t}x_0)dx_0 + O\left(d^6 \log^3 T \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right)\right) p_X(x),
\end{aligned} \tag{85}$$

provided that  $d^2 \log T \left(\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}\right) \lesssim 1$ . We thus complete the proof of Lemma 3.

## B Analysis for the stochastic samplers (Theorems 3 and 4)

sec:analysis-stochastic-samplers

The proofs of Theorem 3 and Theorem 4 share the same structure, therefore in the following, we shall treat them simultaneously. We point out their differences as we move along.

In order to compute the KL divergence between  $p_{X_1}$  and  $p_{Y_1}$ . Let us recall the basic inequality that

$$\begin{aligned}
\text{KL}(p_{X_1} \parallel p_{Y_1}) &\leq \text{KL}(p_{X_1, \dots, X_T} \parallel p_{Y_1, \dots, Y_T}) \\
&= \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{t=2}^T \mathbb{E}_{X_t} [\text{KL}(p_{X_{t-1}|X_t} \parallel p_{Y_{t-1}|Y_t})].
\end{aligned} \tag{eqn:kl-decomp} \tag{86}$$

Recall that in the forward process (4), one has

$$X_T = \sqrt{\bar{\alpha}_T} X_0 + \sqrt{1 - \bar{\alpha}_T} \bar{W}_T, \quad \bar{W}_t \sim \mathcal{N}(0, I_d).$$

By direct calculations, we find

$$\begin{aligned}
\text{KL}(p_{X_T} \parallel p_{Y_T}) &= \int p_{X_T}(x) \log \frac{p_{X_T}(x)}{p_{Y_T}(x)} dx \\
&\stackrel{(i)}{=} \int p_{X_T}(x) \log \frac{\int_{y+z=x, |y| \leq 1/T^c} p_{\sqrt{\bar{\alpha}_T} X_0}(y) p_{\sqrt{1-\bar{\alpha}_T} \bar{W}_T}(z) dy dz}{p_{Y_T}(x)} dx \\
&\leq \int p_{X_T}(x) \log \frac{\max_{y+z=x, |y| \leq 1/T^c} p_{\sqrt{1-\bar{\alpha}_T} \bar{W}_T}(z)}{p_{Y_T}(x)} dx \\
&\leq \int p_{X_T}(x) \left( -d/2 \log(1 - \bar{\alpha}_T) - \frac{\bar{\alpha}_T \|x\|_2^2 + \|x\|_2/T^c + 1/T^{2c}}{2(1 - \bar{\alpha}_T)} \right) dx \\
&= O(\bar{\alpha}_T d).
\end{aligned} \tag{eqn:KL-T} \tag{87}$$

where in (i), we recall that  $\|X_0\|_2 \lesssim T^{c_3}$  while  $\bar{\alpha}_T \lesssim T^{-c_2}$ . It thus suffices for us to control  $\text{KL}(p_{X_{t-1}|X_t} \parallel p_{Y_{t-1}|Y_t})$  for each  $t$ .

### B.1 Step 1: approximation of $p_{X_{t-1}|X_t}(x_{t-1}|x_t)$

Given any constant  $\gamma \in [0, 1]$ , let us write

$$x_t(\gamma) = \gamma x_{t-1} + (1 - \gamma) \hat{x}_t \quad \text{and} \quad \hat{x}_t = x_t / \sqrt{\alpha_t}. \tag{88}$$

With this piece of notation, we proceed to computing  $p_{X_{t-1}|X_t}(x_{t-1}|x_t)$ . Some direct calculations regarding the forward process (4) give

$$\begin{aligned}
&p_{X_{t-1}|X_t}(x_{t-1}|x_t) \\
&\propto p_{X_{t-1}, X_t}(x_{t-1}, x_t) \propto \exp(\log p_{X_{t-1}}(x_{t-1}) + \log p_{X_t|X_{t-1}}(x_t|x_{t-1}))
\end{aligned}$$

$$\begin{aligned}
&= \exp \left( \log p_{X_{t-1}}(\hat{x}_t) + \int_0^1 \nabla \log p_{X_{t-1}}(x_t(\gamma))^\top (x_{t-1} - \hat{x}_t) d\gamma + \log p_{X_t | X_{t-1}}(x_t | x_{t-1}) \right) \\
&\propto \exp \left( (x_{t-1} - \hat{x}_t)^\top \int_0^1 d\gamma \int_{x_0} \frac{\nabla p_{X_{t-1} | X_0}(x_t(\gamma) | x_0) p_{X_0}(x_0)}{p_{X_{t-1}}(x_t(\gamma))} dx_0 + \log p_{X_t | X_{t-1}}(x_t | x_{t-1}) \right). \tag{eqn:sde-cond} \tag{89}
\end{aligned}$$

To continue, in view of the mean value theorem, there exists some  $\tilde{x}_t = x_t(\tilde{\gamma})$  such that

$$\begin{aligned}
&(x_{t-1} - \hat{x}_t)^\top \int_0^1 d\gamma \int_{x_0} \frac{\nabla p_{X_{t-1} | X_0}(x_t(\gamma) | x_0) p_{X_0}(x_0)}{p_{X_{t-1}}(x_t(\gamma))} dx_0 \\
&= (x_{t-1} - \hat{x}_t)^\top \int_0^1 d\gamma \int_{x_0} \frac{\nabla p_{X_{t-1} | X_0}(x_t(\gamma) | x_0)}{p_{X_{t-1} | X_0}(x_t(\gamma) | x_0)} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) dx_0 \\
&= -(x_{t-1} - \hat{x}_t)^\top \int_0^1 d\gamma \int_{x_0} \frac{x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0}{1 - \bar{\alpha}_{t-1}} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) dx_0 \\
&=: -(x_{t-1} - \hat{x}_t)^\top \int_0^1 F(\gamma) d\gamma \tag{90} \\
&= - \frac{(x_{t-1} - \hat{x}_t)^\top \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0 + \frac{1}{2} (x_{t-1} - \hat{x}_t)^\top J_{t-1}(\tilde{x}_t) (x_{t-1} - \hat{x}_t)}{1 - \bar{\alpha}_{t-1}}. \tag{eqn:schubert} \tag{91}
\end{aligned}$$

Here, we compute the Jacobian matrix as in (77) where

$$\begin{aligned}
J_{t-1,ij} &:= \frac{\partial F(\gamma)}{\partial x_j} = \frac{\partial \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x) (x_i - \sqrt{\bar{\alpha}_t} x_{0,i}) dx_0}{\partial x_j} \\
&= \delta_{ij} + \frac{1}{1 - \bar{\alpha}_t} \left( \left( \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x) (x_i - \sqrt{\bar{\alpha}_t} x_{0,i}) dx_0 \right) \left( \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x) (x_j - \sqrt{\bar{\alpha}_t} x_{0,j}) dx_0 \right) \right. \\
&\quad \left. - \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x) (x_i - \sqrt{\bar{\alpha}_t} x_{0,i}) (x_j - \sqrt{\bar{\alpha}_t} x_{0,j}) dx_0 \right) \\
&= \delta_{ij} - \frac{1}{1 - \bar{\alpha}_t} \left( \mathbb{E}_{X_0}[X_0 X_0^\top | X_{t-1} = x] - \mathbb{E}_{X_0}[X_0 | X_{t-1} = x] \cdot \mathbb{E}_{X_0}[X_0 | X_{t-1} = x]^\top \right). \tag{eqn:derivative-2} \tag{92}
\end{aligned}$$

We prove at the end of this section that if we define set

$$\mathcal{E} := \left\{ (x_t, x_{t-1}) \mid -\log p_{X_t}(x_t) \lesssim d \log T, \|x_{t-1} - x_t / \sqrt{\alpha_t}\| \lesssim \sqrt{d \log T (1 - \alpha_t)} \right\}, \tag{eqn:eset} \tag{93}$$

for every  $(x_t, x_{t-1}) \in \mathcal{E}$ , one has <sup>eqn:unfinished</sup>

$$\begin{aligned}
p_{X_{t-1} | X_t}(x_{t-1} | x_t) &\propto \exp \left( -\frac{\alpha_t}{2(1 - \alpha_t)} \left\| \left( I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(\hat{x}_t) \right)^{-1} (x_{t-1} - \mu_t) \right\|^2 + O \left( d^3 \log^3 T \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^{3/2} \right) \right) \\
&\propto \exp \left( -\frac{\alpha_t}{2(1 - \alpha_t)} \|x_{t-1} - \mu_t\|^2 + O \left( d^2 \log^2 T \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right) \right) \right). \tag{eq:cond-dist} \tag{94a} \\
&\tag{eq:cond-dist-crude} \tag{94b}
\end{aligned}$$

Here,  $\mu_t$  is given by

$$\mu_t := \hat{x}_t - \frac{(1 - \alpha_t) \int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0}{\sqrt{\bar{\alpha}_t} (1 - \bar{\alpha}_t)}. \tag{eqn:mu-t-2} \tag{95}$$

We remind the readers that the original DDPM (i.e. Theorem 3) approximates the conditional distribution of  $X_{t-1} | X_t$  via the form (94b) while the accelerated DDPM (i.e. Theorem 4) approximates using the form as in (94a).

## B.2 Step 2: uniform control of the density ratios

Recalling the distribution of  $Y_{t-1}$  (cf. (30)) to obtain

$$p_{Y_{t-1} | Y_t}(x_{t-1} | x_t) \propto \exp \left( -\frac{\alpha_t}{2(1 - \alpha_t)} \left\| \left( I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(\hat{x}_t) \right)^{-1} (x_{t-1} - \mu_t) \right\|^2 \right), \tag{96}$$



with  $\mu$  and  $J$  defined as in expressions (95) and (92). In this section, we prove that

$$\log \frac{p_{X_{t-1} | X_t}(x_{t-1} | x_t)}{p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)} \lesssim \text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2) + \text{poly}(T). \quad \text{eq:SDE-ratio-crude} \quad (97)$$

According to Cauchy's inequality, one can easily check that  $J_t(\hat{x}_t) \preceq I$ , and

$$J_t(\hat{x}_t) \succ -\frac{1}{1 - \bar{\alpha}_{t-1}} \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|\hat{x}_t - \sqrt{\bar{\alpha}_t} x_0\|^2 dx_0 \succ -\frac{2(\|\hat{x}_t\|^2 + \text{poly}(T))}{1 - \bar{\alpha}_t} I, \quad (98)$$

where we recall that  $\|X_0\|_2 \leq \text{poly}(T)$ . If we write

$$Y_{t-1} | Y_t = x_t \sim \mathcal{N}(\mu_t, U_t M_t U_t^\top),$$

where  $\{M_{t,i}\}$  denote eigenvalues of covariance matrix  $\frac{1-\alpha_t}{\alpha_t}(I - (1-\alpha_t)/2(1-\bar{\alpha}_t)J_t(\hat{x}_t))$ , in view of the above properties of  $J_t$ , one has

$$M_{t,i} \gtrsim 1 - \alpha_t, \quad \text{and} \quad \log M_{t,i} \lesssim \log(\|\hat{x}_t\| + \text{poly}(T)).$$

Given these relations, we are ready to bound the density function  $p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)$ . Recall that the basic relation that for  $Z \sim \mathcal{N}(0, \Sigma)$

$$\log \frac{1}{p_Z(z)} = \frac{z^\top \Sigma^{-1} z}{2} + \frac{1}{2} \log \det(\Sigma) + \frac{d}{2} \log(2\pi),$$

which leads to

$$\begin{aligned} \log \frac{1}{p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)} &\lesssim \frac{\|x_{t-1} - \mu_t\|^2}{1 - \alpha_t} + d \log(\|\hat{x}_t\| + \text{poly}(T)) \\ &\lesssim \text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2) + \text{poly}(T), \end{aligned} \quad \text{eqn:y-condi} \quad (99)$$

Here, the last inequality uses

$$\begin{aligned} \|x_{t-1} - \mu_t\|^2 &\leq 2\|x_{t-1} - \hat{x}_t\|^2 + 2\|\hat{x}_t - \mu_t\|^2 \\ &\lesssim \|x_{t-1} - \hat{x}_t\|^2 + \frac{(1 - \alpha_t)}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \int_{x_0} p_{X_0 | X_t}(x_0 | x_t) \|\sqrt{\bar{\alpha}_t} \hat{x}_t - \sqrt{\bar{\alpha}_t} x_0\|^2 dx_0 \\ &\leq \text{poly}(T)(\|x_{t-1} - \hat{x}_t\|^2 + \|\hat{x}_t\|^2) + \text{poly}(T) \end{aligned}$$

for  $x_t = \sqrt{\bar{\alpha}_t} \hat{x}_t$  and  $\|X_0\|_2 \leq \text{poly}(T)$ .

In addition, since  $X_t | X_{t-1}$  follows a Gaussian distribution, we can write

$$\begin{aligned} \log p_{X_{t-1} | X_t}(x_{t-1} | x_t) &= \log \frac{p_{X_t | X_{t-1}}(x_t | x_{t-1}) p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} \\ &= \log \frac{p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} + \frac{\|x_t - \sqrt{\bar{\alpha}_t} x_{t-1}\|^2}{2(1 - \alpha_t)} - d/2 \log(2\pi(1 - \alpha_t)) \\ &\leq \log \frac{p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} + \text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2). \end{aligned}$$

Some direct calculations yield

$$\begin{aligned} \log \frac{p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} &= \log \frac{\int_{x_0} p_{X_0}(x_0) \exp\left(-\frac{\|x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2}{2(1 - \bar{\alpha}_{t-1})}\right) dx_0}{\int_{x_0} p_{X_0}(x_0) \exp\left(-\frac{\|x_t - \sqrt{\bar{\alpha}_t} x_0\|^2}{2(1 - \bar{\alpha}_t)}\right) dx_0} - \frac{d}{2} \log\left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\right) \\ &\leq \max_{x_0} \left\{ \|x_t - \sqrt{\bar{\alpha}_t} x_0\|^2 - \|x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2 \right\} - \frac{d}{2} \log\left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\right) \\ &\leq \bar{\alpha}_t \|\hat{x}_t - x_0\|^2 \leq \|x_{t-1}\|^2 + \text{poly}(T). \end{aligned}$$

Combining the above two relations together, we arrive at

$$\log p_{X_{t-1} | X_t}(x_{t-1} | x_t) \lesssim \text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2) + \text{poly}(T), \quad \text{eqn:x-condi} \quad (100)$$

Putting everything together, we complete the proof of (97).

### B.3 Step 3: computing the KL divergence

Based on the analyses laid out previously, we are ready to compute the KL divergence between  $p_{X_{t-1} | X_t}$  and  $p_{Y_{t-1} | Y_t}$ . In view of relations (94a), (97) and set (93), we decompose the integration as

$$\begin{aligned} & \mathbb{E}_{X_t} [\text{KL}(p_{X_{t-1} | X_t} \parallel p_{Y_{t-1} | Y_t})] \\ &= \left( \int_{\mathcal{E}} + \int_{\mathcal{E}^c} \right) p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) \log \frac{p_{X_{t-1} | X_t}(x_{t-1} | x_t)}{p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)} dx_{t-1} dx_t, \\ &\stackrel{(i)}{=} \int_{\mathcal{E}} p_{X_t}(x_t) \left( p_{X_{t-1} | X_t}(x_{t-1} | x_t) - p_{Y_{t-1} | Y_t}(x_{t-1} | x_t) + p_{X_{t-1} | X_t}(x_{t-1} | x_t) \cdot O\left(d^6 \log^6 T \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right) \right) dx_{t-1} dx_t \\ &\quad + \int_{\mathcal{E}^c} p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) (\text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2) + \text{poly}(T)) dx_{t-1} dx_t. \end{aligned} \tag{eqn:serenade} \tag{101}$$

Here, expression (i) makes use of the facts that if  $|\frac{p_Y(x)}{p_X(x)} - 1| < \frac{1}{2}$ ,

$$\begin{aligned} p_X(x) \log \frac{p_X(x)}{p_Y(x)} &= -p_X(x) \log \left( 1 + \frac{p_Y(x) - p_X(x)}{p_X(x)} \right) \\ &= p_X(x) - p_Y(x) + p_X(x) O\left(\left(\frac{p_Y(x)}{p_X(x)} - 1\right)^2\right). \end{aligned} \tag{102}$$

In particular, (101) holds true when we instantiate the above inequality to  $p_{X_{t-1} | X_t}(x_{t-1} | x_t)$  and  $p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)$ , and recall expression (94a) that

$$p_{X_{t-1} | X_t}(x_{t-1} | x_t) = p_{Y_{t-1} | Y_t}(x_{t-1} | x_t) \cdot \exp\left(O\left(d^3 \log^3 T \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^{3/2}\right)\right).$$

To continue, let us bound each term on the right hand side of (101) respectively. Given set  $\mathcal{E}$  (cf. 93), direct calculations yield

$$\begin{aligned} \mathbb{P}((x_{t-1}, x_t) \in \mathcal{E}^c) &= \int_{\mathcal{E}^c} p_{X_{t-1}}(x_{t-1}) p_{X_t | X_{t-1}}(x_t | x_{t-1}) dx_{t-1} dx_t \\ &= \int_{\mathcal{E}^c} p_{X_{t-1}}(x_{t-1}) \exp\left(-d/2 \log(2\pi(1 - \alpha_t)) - \frac{\|x_t - \sqrt{\alpha_t} x_{t-1}\|^2}{2(1 - \alpha_t)}\right) dx_{t-1} dx_t \\ &\leq \exp(-\Omega(d \log T)), \end{aligned}$$

and similarly,

$$\int_{\mathcal{E}^c} p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) (\text{poly}(T)(\|\hat{x}_t\|^2 + \|x_{t-1} - \hat{x}_t\|^2) + \text{poly}(T)) dx_{t-1} dx_t \leq \exp(-\Omega(d \log T)).$$

In addition, for every  $x_t$  obeying  $-\log p_{X_t}(x_t) \lesssim d \log T$ , it satisfies that

$$\int_{\|x_{t-1} - x_t / \sqrt{\alpha_t}\|^2 \gtrsim \sqrt{d \log T (1 - \alpha_t)}} p_{Y_{t-1} | Y_t}(x_{t-1} | x_t) dx_{t-1} \leq \exp(-\Omega(d \log T)),$$

by recalling that

$$p_{Y_{t-1} | Y_t}(x_{t-1} | x_t) \propto \exp\left(-\frac{\alpha_t}{2(1 - \alpha_t)} \left\| \left(I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(\hat{x}_t)\right)^{-1} (x_{t-1} - \mu_t) \right\|^2\right), \tag{103}$$

where  $\|J_t(\hat{x}_t)\|_{\text{op}} \lesssim d \log T$  and  $\|\mu_t - \hat{x}_t\|_2 \lesssim (1 - \alpha_t) \sqrt{d \log T / (1 - \bar{\alpha}_t)}$ .

Taking everything collectively, for each  $t \geq 2$ , the right hand side of (101) satisfies

$$\mathbb{E}_{X_t} [\text{KL}(p_{X_{t-1} | X_t} \parallel p_{Y_{t-1} | Y_t})] = O\left(d^6 \log^6 T \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t}\right)^3\right). \tag{eqn:each-kl} \tag{104}$$

## B.4 In summary

To summarize, let us recall the decomposition (86) and the control of each term (104) to obtain

$$\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{2 \leq t \leq T} d^6 \log^6 T \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3$$

Putting the above together with our choice of learning rates (cf. (21)) and the bound of  $\text{KL}(p_{X_T} \parallel p_{Y_T})$  as in (87), yields

$$\text{KL}(p_{X_1} \parallel p_{Y_1}) \lesssim \frac{d^6 \log^9 T}{T^2}.$$

which completes the claimed bound (34) by virtue of Pinsker's inequality.

Similarly, if we instead consider the approximation (94b) concerning the original DDPM, the first term in (101) will be replaced by  $O(d^2 \log^2 T (\frac{1-\alpha_t}{\alpha_t - \bar{\alpha}_t}))$ . As a consequence, we can derive

$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) &\lesssim \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{t \geq 2}^T d^4 \log^4 T \left( \frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^2 \\ &\lesssim \frac{d^4 \log^6 T}{T}, \end{aligned} \tag{105}$$

which again leads to (31).

## B.5 Proof of Claim (94)

In order to prove these relations, note that putting together (89) and (91) leads to

$$\begin{aligned} &p_{X_{t-1} | X_t}(x_{t-1} | x_t) \\ &\propto \exp \left( (x_{t-1} - \hat{x}_t)^\top \int_0^1 d\gamma \int_{x_0} \frac{\nabla p_{X_{t-1} | X_0}(\tilde{x}_t | x_0) p_{X_0}(x_0)}{p_{X_{t-1}}(\tilde{x}_t)} dx_0 + \log p_{X_t | X_{t-1}}(x_t | x_{t-1}) \right) \\ &\propto \exp \left( - \frac{(x_{t-1} - \hat{x}_t)^\top \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0 + \frac{1}{2} (x_{t-1} - \hat{x}_t)^\top J_{t-1}(\tilde{x}_t) (x_{t-1} - \hat{x}_t)}{1 - \bar{\alpha}_{t-1}} \right. \\ &\quad \left. - \frac{\|x_t - \sqrt{\alpha_t} x_{t-1}\|^2}{2(1 - \alpha_t)} \right). \end{aligned} \tag{106} \quad \text{eqn:allegro}$$

If we can show that for every  $\gamma \in [0, 1]$ ,  $\text{eq:Jacobi}$

$$\|J_{t-1}(x_t(\gamma)) - I\| \lesssim d \log T, \tag{107a} \quad \text{eq:Jacobi-a}$$

$$\|J_{t-1}(x_t(\gamma)) - J_{t-1}(\hat{x}_t)\| \lesssim d^2 \log^2 T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}, \tag{107b} \quad \text{eq:Jacobi-b}$$

and  $\text{eq:approx-t}$

$$\left\| \frac{\int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0}{1 - \bar{\alpha}_{t-1}} - \frac{\int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0}{1 - \bar{\alpha}_t} \right\| \lesssim (1 - \alpha_t) \left( \frac{d \log T}{\alpha_t - \bar{\alpha}_t} \right)^{3/2}, \tag{108a} \quad \text{eq:approx-t-a}$$

$$\left\| \frac{J_{t-1}(\hat{x}_t)}{1 - \bar{\alpha}_{t-1}} - \frac{J_t(\hat{x}_t)}{1 - \bar{\alpha}_t} \right\| \lesssim \frac{d^2 \log^2 T (1 - \alpha_t)}{(\alpha_t - \bar{\alpha}_t)^2}, \tag{108b} \quad \text{eq:approx-t-b}$$

the relation (94) shall follow immediately by recognizing that  $\|x_{t-1} - \hat{x}_t\|_2 \lesssim \sqrt{d \log T (1 - \alpha_t)}$ .

### B.5.1 Preliminaries

For  $(x_t, x_{t-1}) \in \mathcal{E}$ , it holds that

$$-\log p_{X_t}(x_t) \lesssim d \log T, \quad \text{and} \quad \|x_{t-1} - \hat{x}_t\| \lesssim \sqrt{d \log T(1 - \alpha_t)}. \quad (109)$$

Based on this relation, we derive that

$$-\log p_{X_{t-1}}(x_t(\gamma)) \lesssim d \log T. \quad \text{eqn:river} \quad (110)$$

**Proof of relation (110).** For every  $x$  such that  $\|\hat{x}_t - x\| \lesssim \sqrt{d \log T(1 - \alpha_t)}$ , consider the density ratio

$$\begin{aligned} \frac{p_{X_{t-1}|X_0}(x|x_0)}{p_{X_t|X_0}(x_t|x_0)} &= \left( \frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}} \right)^{d/2} \exp \left( \frac{\|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{2(1 - \bar{\alpha}_t)} - \frac{\|x - \sqrt{\bar{\alpha}_{t-1}}x_0\|^2}{2(1 - \bar{\alpha}_{t-1})} \right) \\ &= \exp \left( O \left( \frac{d(1 - \alpha_t) + \|\hat{x}_t - x\|^2 + \|\hat{x}_t - x\| \|x_t - \sqrt{\bar{\alpha}_t}x_0\|}{1 - \bar{\alpha}_{t-1}} + \frac{(1 - \alpha_t) \|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{(1 - \bar{\alpha}_{t-1})^2} \right) \right). \end{aligned}$$

If we define set

$$\tilde{\mathcal{E}} := \{x_0 : \|x_t - \sqrt{\bar{\alpha}_t}x_0\| \lesssim \sqrt{d \log T(1 - \bar{\alpha}_t)}\},$$

for every  $x_0 \in \tilde{\mathcal{E}}$ , we find

$$\frac{p_{X_{t-1}|X_0}(x|x_0)}{p_{X_t|X_0}(x_t|x_0)} = 1 + O \left( d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}} \right). \quad \text{eqn:scriabin} \quad (111)$$

Here we remind the readers on our choice of step sizes in (21). Otherwise, for  $x_0 \notin \tilde{\mathcal{E}}$ ,

$$\frac{p_{X_{t-1}|X_0}(x|x_0)}{p_{X_t|X_0}(x_t|x_0)} = \exp \left( O \left( \frac{\log T \|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{T(1 - \bar{\alpha}_t)} \right) \right).$$

Based on the above bounds, some direct calculations yield

$$\begin{aligned} p_{X_{t-1}}(x) &= \int_{x_0} p_{X_0}(x_0) p_{X_{t-1}|X_0}(x|x_0) dx_0 \\ &= \int_{x_0 \in \tilde{\mathcal{E}}} p_{X_0}(x_0) p_{X_t|X_0}(x_t|x_0) \cdot \frac{p_{X_{t-1}|X_0}(x|x_0)}{p_{X_t|X_0}(x_t|x_0)} dx_0 \\ &\quad + \int_{x_0 \in \tilde{\mathcal{E}}^c} p_{X_0}(x_0) p_{X_t|X_0}(x_t|x_0) \cdot \frac{p_{X_{t-1}|X_0}(x|x_0)}{p_{X_t|X_0}(x_t|x_0)} dx_0 \\ &= p_{X_t}(x_t) \int_{x_0 \in \tilde{\mathcal{E}}} \left( 1 + O \left( d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}} \right) \right) p_{X_0|X_t}(x_0|x_t) dx_0 \\ &\quad + p_{X_t}(x_t) \int_{x_0 \in \tilde{\mathcal{E}}^c} \exp \left( O \left( \frac{\log T \|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{T(1 - \bar{\alpha}_t)} \right) \right) p_{X_0|X_t}(x_0|x_t) dx_0. \quad \text{eqn:rachmaninoff} \quad (112) \end{aligned}$$

Recall that in Lemma 1, we showed that if  $-\log p_{X_t}(x_t) \lesssim d \log T$ , For random variable  $X = X_0 | X_t = x_t$ , it holds with probability at least  $1 - \exp(-\Omega(d \log T))$  that

$$\|\sqrt{\bar{\alpha}_t}X - y\| \lesssim \sqrt{d \log T(1 - \bar{\alpha}_t)}. \quad (113)$$

As a result, the right hand side of (112) equals to

$$p_{X_{t-1}}(x) = \left( 1 + O \left( d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}} \right) \right) (1 + \exp(-\Omega(d \log T))) p_{X_t}(x_t) + \exp(-\Omega(d \log T)) p_{X_t}(x_t)$$

$$= \left(1 + O\left(d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}\right)\right) p_{X_t}(x_t). \quad \text{eqn:prokofiev} \quad (114)$$

Therefore, relation (110) follows immediately as  $-\log p_{X_t}(x_t) \lesssim d \log T$ .

Combining relation (110) with Lemma 1, with probability at least  $1 - \exp(-\Omega(d \log T))$ , it satisfies that

$$X = X_0 | X_{t-1} = x_t(\gamma), \quad \sqrt{\bar{\alpha}_{t-1}} X = \hat{x}_t + O(\sqrt{d \log T(1 - \bar{\alpha}_{t-1})}). \quad \text{eqn:condi-g} \quad (115)$$

### B.5.2 Main analysis

**Proof of relation (107a).** As a consequence of property (115), inequality (107a) holds true by noticing that

$$\|J_{t-1}(x_t(\gamma)) - I\|_{\text{op}} \leq \frac{2 \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) \|x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2 dx_0}{1 - \bar{\alpha}_{t-1}} \lesssim d \log T,$$

where similar to the proof of (112), the integration is computed on the set  $\|\sqrt{\bar{\alpha}_{t-1}} x_0 - \hat{x}_t\| \lesssim \sqrt{d \log T(1 - \bar{\alpha}_{t-1})}$  and its complement respectively.

**Proof of relation (107b).** Consider  $x_0$  such that  $\|\sqrt{\bar{\alpha}_{t-1}} x_0 - \hat{x}_t\| \lesssim \sqrt{d \log T(1 - \bar{\alpha}_{t-1})}$  and

$$\|x_t(\gamma) - \hat{x}_t\| \leq \|x_{t-1} - \hat{x}_t\| \lesssim \sqrt{d \log T(1 - \alpha_t)} \leq \sqrt{\frac{1 - \bar{\alpha}_t}{d \log T}}. \quad (116)$$

It guarantees that

$$\begin{aligned} \frac{p_{X_{t-1} | X_0}(x_t(\gamma) | x_0)}{p_{X_t(\gamma) | X_0}(\hat{x}_t | x_0)} &= \exp\left(\frac{\|\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{\|x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2}{2(1 - \bar{\alpha}_{t-1})}\right) \\ &= 1 + O\left(\frac{\|\hat{x}_t - x_t(\gamma)\|^2 + \|\hat{x}_t - x_t(\gamma)\| \sqrt{d \log T(1 - \bar{\alpha}_{t-1})}}{1 - \bar{\alpha}_{t-1}}\right). \end{aligned}$$

Similar to the proof of (112), we can also derive

$$\begin{aligned} \frac{p_{X_{t-1}}(x_t(\gamma))}{p_{X_{t-1}}(\hat{x}_t)} &= \frac{\int_{x_0} p_{X_0}(x_0) p_{X_{t-1} | X_0}(x_t(\gamma) | x_0) dx_0}{\int_{x_0} p_{X_0}(x_0) p_{X_{t-1} | X_0}(\hat{x}_t | x_0) dx_0} \quad \text{eqn:cond-xt-1} \quad (117) \\ &= 1 + O\left(\frac{\|\hat{x}_t - x_t(\gamma)\|^2 + \|\hat{x}_t - x_t(\gamma)\| \sqrt{d \log T(1 - \bar{\alpha}_{t-1})}}{1 - \bar{\alpha}_{t-1}}\right). \quad (118) \end{aligned}$$

Putting everything together leads to

$$\frac{p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma))}{p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t)} = \frac{p_{X_{t-1} | X_0}(x_t(\gamma) | x_0) / p_{X_{t-1}}(x_t(\gamma))}{p_{X_{t-1} | X_0}(\hat{x}_t | x_0) / p_{X_{t-1}}(\hat{x}_t)} = 1 + O\left(d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}\right).$$

Equipped with the above relation, we find

$$\begin{aligned} &\left\| \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) (x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0 - \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0 \right\| \\ &= O\left(d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}\right) \cdot \left\| \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) (x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0 \right\| \\ &\lesssim \sqrt{d^3 \log^3 T(1 - \alpha_t)}, \end{aligned}$$

where the last step invokes the property (115) again. Similarly, we also arrive at

$$\left\| \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) (x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0) (x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0)^\top dx_0 \right\|$$

$$\begin{aligned}
& - \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0)^\top dx_0 \Big\| \\
& = O\left(d \log T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}\right) \cdot \left\| \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t(\gamma)) (x_t(\gamma) - \sqrt{\bar{\alpha}_{t-1}} x_0) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0)^\top dx_0 \right\| \\
& \lesssim d^2 \log^2 T \sqrt{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}.
\end{aligned}$$

Taking the above two above relations together with the definition (92) gives

$$\|J_{t-1}(x_t(\gamma)) - J_{t-1}(\hat{x}_t)\| \lesssim d^2 \log^2 T \sqrt{\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}}.$$

**Proof of relation (108a)** To control the quantity on the left of (108a), by triangle's inequality, let us write

$$\begin{aligned}
& \left\| \frac{\int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0) dx_0}{1 - \bar{\alpha}_{t-1}} - \frac{\int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0}{1 - \bar{\alpha}_t} \right\| \\
& \leq \left\| \frac{\int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 - \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0}{\alpha_t - \bar{\alpha}_t} \right\| \\
& \quad + \left\| \left( \frac{1}{\alpha_t - \bar{\alpha}_t} - \frac{1}{1 - \bar{\alpha}_t} \right) \int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right\|. \tag{eqn:adagio} \tag{119}
\end{aligned}$$

As computed previously, given  $-\log p_{X_t}(x_t) \lesssim d \log T$ , Lemma 1 implies that

$$\begin{aligned}
\int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 &= \int_{\|x_t - \sqrt{\bar{\alpha}_t} x_0\| \leq \sqrt{d \log T(1 - \bar{\alpha}_t)}} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 \\
&\quad + \int_{\|x_t - \sqrt{\bar{\alpha}_t} x_0\| > \sqrt{d \log T(1 - \bar{\alpha}_t)}} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 \\
&\lesssim \sqrt{d \log T(1 - \bar{\alpha}_t)}. \tag{eqn:roundo} \tag{120}
\end{aligned}$$

Additionally, similar argument gives

$$\begin{aligned}
& \int_{x_0} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 \\
&= \int_{x_0 \in \|x_t - \sqrt{\bar{\alpha}_t} x_0\| \leq \sqrt{d \log T(1 - \bar{\alpha}_t)}} p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 + \exp(-\Omega(d \log T)) \\
&\stackrel{(i)}{=} \left(1 + O\left(\frac{d \log T(1 - \alpha_t)}{1 - \bar{\alpha}_{t-1}}\right)\right) \int_{x_0 \in \|x_t - \sqrt{\bar{\alpha}_t} x_0\| \leq \sqrt{d \log T(1 - \bar{\alpha}_t)}} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0 + \exp(-\Omega(d \log T)) \\
&= \left(1 + O\left(\frac{d \log T(1 - \alpha_t)}{1 - \bar{\alpha}_{t-1}}\right)\right) \int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0, \tag{eqn:schertzo} \tag{121}
\end{aligned}$$

where in step (i), we make use of the following relation that

$$p_{X_0 | X_{t-1}}(x_0 | \hat{x}_t) = p_{X_0}(x_0) \frac{p_{X_{t-1} | X_0}(\hat{x}_t | x_0)}{p_{X_{t-1}}(\hat{x}_t)} = \left(1 + O\left(\frac{d \log T(1 - \alpha_t)}{1 - \bar{\alpha}_{t-1}}\right)\right) p_{X_0}(x_0) \frac{p_{X_t | X_0}(x_t | x_0)}{p_{X_t}(x_t)}. \tag{eqn:1-to-t} \tag{122}$$

To see why (122) holds, we first observe that

$$\begin{aligned}
\frac{p_{X_{t-1} | X_0}(\hat{x}_t | x_0)}{p_{X_t | X_0}(x_t | x_0)} &= \left(\frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}}\right)^{d/2} \exp\left(\frac{\|x_t - \sqrt{\bar{\alpha}_t} x_0\|^2}{2(1 - \bar{\alpha}_t)} - \frac{\|\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}} x_0\|^2}{2(1 - \bar{\alpha}_{t-1})}\right) \\
&= \left(1 + O\left(\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}}\right)\right)^{d/2} \exp\left(-\frac{(1 - \alpha_t)\|x_t - \sqrt{\bar{\alpha}_t} x_0\|^2}{2(1 - \bar{\alpha}_t)(\alpha_t - \bar{\alpha}_t)}\right)
\end{aligned}$$

$$= \exp \left( O \left( \frac{d(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} + \frac{(1-\alpha_t)\|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{(1-\bar{\alpha}_{t-1})^2} \right) \right).$$

In addition, for  $x_0 \in \mathcal{E} = \{x_0 \mid \|x_t - \sqrt{\bar{\alpha}_t}x_0\| \lesssim \sqrt{d \log T(1-\bar{\alpha}_t)}\}$ , it satisfies

$$\frac{p_{X_{t-1} \mid X_0}(x \mid x_0)}{p_{X_t \mid X_0}(x_t \mid x_0)} = 1 + O \left( \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} \right);$$

on the other hand, for  $x_0 \notin \mathcal{E}$ , we have

$$\frac{p_{X_{t-1} \mid X_0}(x \mid x_0)}{p_{X_t \mid X_0}(x_t \mid x_0)} \lesssim \exp \left( O \left( \frac{\log T \|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{T(1-\bar{\alpha}_t)} \right) \right).$$

Similar to (112), integrate to obtain

$$\begin{aligned} p_{X_t}(x) &= \int_{x_0} p_{X_0}(x_0) p_{X_{t-1} \mid X_0}(x \mid x_0) dx_0 \\ &= \int_{x_0 \in \mathcal{E}} p_{X_0}(x_0) p_{X_t \mid X_0}(x_t \mid x_0) \cdot \frac{p_{X_{t-1} \mid X_0}(x \mid x_0)}{p_{X_t \mid X_0}(x_t \mid x_0)} dx_0 \\ &\quad + \int_{x_0 \in \mathcal{E}^c} p_{X_0}(x_0) p_{X_t \mid X_0}(x_t \mid x_0) \cdot \frac{p_{X_{t-1} \mid X_0}(x \mid x_0)}{p_{X_t \mid X_0}(x_t \mid x_0)} dx_0 \\ &= p_{X_t}(x_t) \int_{x_0 \in \mathcal{E}} \left( 1 + O \left( \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} \right) \right) p_{X_0 \mid X_t}(x_0 \mid x_t) dx_0 \\ &\quad + p_{X_t}(x_t) \int_{x_0 \in \mathcal{E}^c} \exp \left( O \left( \frac{\log T \|x_t - \sqrt{\bar{\alpha}_t}x_0\|^2}{T(1-\bar{\alpha}_t)} \right) \right) p_{X_0 \mid X_t}(x_0 \mid x_t) dx_0 \\ &= \left( 1 + O \left( \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} \right) \right) (1 + \exp(-\Omega(d \log T))) p_{X_t}(x_t) + \exp(-\Omega(d \log T)) p_{X_t}(x_t) \\ &= \left( 1 + O \left( \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} \right) \right) p_{X_t}(x_t). \end{aligned}$$

Plugging expressions (121) and (120) into (119) gives

$$\left\| \frac{\int_{x_0} p_{X_0 \mid X_{t-1}}(x_0 \mid \hat{x}_t) (\hat{x}_t - \sqrt{\bar{\alpha}_{t-1}}x_0) dx_0}{1-\bar{\alpha}_{t-1}} - \frac{\int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x_t) (x_t - \sqrt{\bar{\alpha}_t}x_0) dx_0}{1-\bar{\alpha}_t} \right\| \lesssim (1-\alpha_t) \left( \frac{d \log T}{\alpha_t - \bar{\alpha}_t} \right)^{3/2}.$$

**Proof of relation (108b).** First observe that

$$\left\| \frac{J_{t-1}(\hat{x}_t)}{1-\bar{\alpha}_{t-1}} - \frac{J_t(\hat{x}_t)}{1-\bar{\alpha}_t} \right\| \leq \left\| \frac{J_{t-1}(\hat{x}_t) - J_t(\hat{x}_t)}{1-\bar{\alpha}_{t-1}} \right\| + \left\| \left( \frac{1}{1-\bar{\alpha}_{t-1}} - \frac{1}{1-\bar{\alpha}_t} \right) J_t(\hat{x}_t) \right\|.$$

By virtue of relation (122), one can deduce

$$\begin{aligned} &\int_{x_0} p_{X_0 \mid X_{t-1}}(x_0 \mid \hat{x}_t) (x_t - \sqrt{\bar{\alpha}_t}x_0) (x_t - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0 \\ &= \left( 1 + O \left( \frac{d \log T(1-\alpha_t)}{1-\bar{\alpha}_{t-1}} \right) \right) \int_{x_0} p_{X_0 \mid X_t}(x_0 \mid x_t) (x_t - \sqrt{\bar{\alpha}_t}x_0) (x_t - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0. \end{aligned}$$

This property, combined with relation (121), implies that

$$\left\| \frac{J_{t-1}(\hat{x}_t)}{1-\bar{\alpha}_{t-1}} - \frac{J_t(\hat{x}_t)}{1-\bar{\alpha}_t} \right\| \leq \frac{d^2 \log^2 T(1-\alpha_t)}{(\alpha_t - \bar{\alpha}_t)^2}.$$