

Accelerating Convergence of Score-Based Diffusion Models, Provably

Gen Li^{*†}
CUHK

Yu Huang^{*‡}
UPenn

Timofey Efimov[§]
CMU

Yuting Wei[‡]
UPenn

Yuejie Chi[§]
CMU

Yuxin Chen[‡]
UPenn

March 6, 2024

Abstract

Score-based diffusion models, while achieving remarkable empirical performance, often suffer from low sampling speed, due to extensive function evaluations needed during the sampling phase. Despite a flurry of recent activities towards speeding up diffusion generative modeling in practice, theoretical underpinnings for acceleration techniques remain severely limited. In this paper, we design novel training-free algorithms to accelerate popular deterministic (i.e., DDIM) and stochastic (i.e., DDPM) samplers. Our accelerated deterministic sampler converges at a rate $O(\frac{1}{T^2})$ with T the number of steps, improving upon the $O(\frac{1}{T})$ rate for the DDIM sampler; and our accelerated stochastic sampler converges at a rate $O(\frac{1}{T})$, outperforming the rate $O(\frac{1}{\sqrt{T}})$ for the DDPM sampler. The design of our algorithms leverages insights from higher-order approximation, and shares similar intuitions as popular high-order ODE solvers like the DPM-Solver-2. Our theory accommodates ℓ_2 -accurate score estimates, and does not require log-concavity or smoothness on the target distribution.

Keywords: diffusion models, training-free samplers, DDPM, DDIM, probability flow ODE, higher-order ODE

Contents

1	Introduction	2
1.1	Score-based diffusion models	2
1.2	Non-asymptotic convergence theory and acceleration	3
1.3	Our contributions	4
1.4	Other related works	4
1.5	Notation	5
2	Problem settings	5
2.1	Model and sampling process	5
2.2	Assumptions	7
3	Algorithm and main theory	7
3.1	Accelerated ODE-based sampler	8
3.2	Accelerated SDE-based sampler	10

*The first two authors contributed equally.

†Department of Statistics, The Chinese University of Hong Kong, Hong Kong.

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

§Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

4 Experiments	11
4.1 Practical implementation	11
4.2 Experimental results	12
5 Discussion	12
A Preliminaries	16
A.1 Basic facts	16
A.2 Proof of Lemma 4	18
B Analysis for the accelerated ODE sampler (proof of Theorem 1)	19
B.1 Main steps of the proof	19
B.2 Proof of Lemma 6	25
B.2.1 Proof of property (49)	25
B.2.2 Proof of property (50a)	27
B.2.3 Proof of property (50b)	29
B.2.4 Proof of property (51)	31
B.2.5 Proof of additional lemmas	32
B.3 Proof of Lemma 7	35
C Analysis for the accelerated DDPM sampler (proof of Theorem 2)	36
C.1 Main steps of the proof	36
C.2 Proof of Lemma 11	38
C.3 Proof of Lemma 13	39
C.4 Proof of Lemma 14	40
C.5 Proof of Lemma 15	41

1 Introduction

Initially introduced by [Sohl-Dickstein et al. \(2015\)](#) and subsequently gaining momentum through the works [Ho et al. \(2020\)](#); [Song et al. \(2021\)](#), diffusion models have risen to the forefront of generative modeling. Remarkably, score-based diffusion models have demonstrated superior performance across various domains like computer vision, natural language processing, medical imaging, and bioinformatics ([Croitoru et al., 2023](#); [Yang et al., 2023](#); [Kazerouni et al., 2023](#); [Guo et al., 2023](#)), outperforming earlier generative methods such as GANs ([Goodfellow et al., 2020](#)) and VAEs ([Kingma and Welling, 2014](#)) on multiple fronts ([Dhariwal and Nichol, 2021](#)).

1.1 Score-based diffusion models

On a high level, diffusion-based generative modeling begins by considering a forward Markov diffusion process that progressively diffuses a data distribution into noise:

$$X_0 \xrightarrow{\text{add noise}} X_1 \xrightarrow{\text{add noise}} X_2 \xrightarrow{\text{add noise}} \dots \xrightarrow{\text{add noise}} X_T, \quad (1)$$

where $X_0 \sim p_{\text{data}}$ is drawn from the target data distribution in \mathbb{R}^d , and X_T resembles pure noise (e.g., with a distribution close to $\mathcal{N}(0, I_d)$). The pivotal step then lies in learning to construct a reverse Markov process

$$Y_0 \xleftarrow{\text{use scores}} Y_1 \xleftarrow{\text{use scores}} Y_2 \xleftarrow{\text{use scores}} \dots \xleftarrow{\text{use scores}} Y_T, \quad (2)$$

which starts from pure noise $Y_T \sim \mathcal{N}(0, I_d)$ and maintains distributional proximity throughout in the sense that $Y_t \stackrel{d}{\approx} X_t$ ($t \leq T$). To accomplish this goal, Y_{t-1} in each step is typically obtained from Y_t with the aid of (Stein) score functions — namely, $\nabla_X \log p_{X_t}(X)$, with p_{X_t} denoting the distribution of X_t — where the score functions are pre-trained by means of score matching techniques (e.g., [Hyvärinen \(2005\)](#); [Ho et al. \(2020\)](#); [Hyvärinen \(2007\)](#); [Vincent \(2011\)](#); [Song and Ermon \(2019\)](#); [Pang et al. \(2020\)](#)).

The mainstream approaches for constructing the reverse-time process (2) can roughly be divided into two categories, as described below.

- *Stochastic (or SDE-based) samplers.* A widely adopted strategy involves exploiting both the score function and some injected random noise when generating each Y_{t-1} ; that is, Y_{t-1} is taken to be a function of Y_t and some independent noise Z_t . A prominent example of this kind is the Denoising Diffusion Probabilistic Model (DDPM) (Ho et al., 2020), to be detailed in Section 2. Notably, this approach has intimate connections with certain stochastic differential equations (SDEs), which can be elucidated via celebrated SDE results concerning the existence of reverse-time diffusion processes (Anderson, 1982; Haussmann and Pardoux, 1986).
- *Deterministic (or ODE-based) samplers.* In contrast, another approach is purely deterministic (except for the generation of Y_T), constructing Y_{t-1} as a function of the previously computed steps (e.g., Y_t) without injecting any additional noise. This approach was introduced by Song et al. (2021), as inspired by the existence of ordinary differential equations (ODEs) — termed *probability flow ODEs* — exhibiting the same marginal distributions as the above-mentioned reverse-time diffusion process. A notable example in this category is often referred to as the Denoising Diffusion Implicit Model (DDIM) (Song et al., 2020).

In practice, it is often observed that DDIM converges more rapidly than DDPM, although the final data instances produced by DDPM (given sufficient runtime) enjoy better diversity compared to the output of DDIM.

1.2 Non-asymptotic convergence theory and acceleration

Despite the astounding empirical success, theoretical analysis for diffusion-based generative modeling is still in its early stages of development. Treating the score matching step as a blackbox and exploiting only (crude) information about the score estimation error, a recent strand of works have explored the convergence rate of the data generating process (i.e., the reverse Markov process) in a non-asymptotic fashion, in an attempt to uncover how fast sampling can be performed (e.g., Lee et al. (2022, 2023); Chen et al. (2022, 2023a,c,b); Li et al. (2023); Benton et al. (2023b,a); Liang et al. (2024)). In what follows, let us give a brief overview of the state-of-the-art results in this direction. Here and throughout, the iteration complexity of a sampler refers to the number of steps T needed to attain ε accuracy in the sense that $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$, where $\text{TV}(\cdot, \cdot)$ represents the total-variation (TV) distance between two distributions, and p_{X_1} (resp. p_{Y_1}) stands for the distribution of X_1 (resp. Y_1).

- *Convergence rate of stochastic samplers.* Assuming Lipschitz continuity (or smoothness) of the score functions across all steps, Chen et al. (2022) proved that the iteration complexity of the DDPM sampler is proportional to $1/\varepsilon^2$. The Lipschitz assumption is then relaxed by Chen et al. (2023a); Benton et al. (2023a); Li et al. (2023), revealing that the scaling $1/\varepsilon^2$ is achievable for a fairly general family of data distributions.
- *Convergence rate of deterministic samplers.* As alluded to previously, deterministic samplers often exhibit faster convergence in both practice and theory. For instance, Chen et al. (2023c) provided the first polynomial convergence guarantees for the probability flow ODE sampler under exact scores, whereas Li et al. (2023) demonstrated that its iteration complexity scales proportionally to $1/\varepsilon$ allowing score estimation error. Additionally, it is noteworthy that an iteration complexity proportional to $1/\varepsilon$ has also been established by Chen et al. (2023b) for a variant of the probability flow ODE sampler, although the sampler studied therein incorporates a stochastic corrector step in each iteration.

Acceleration? While the theoretical studies outlined above have offered non-asymptotic convergence guarantees for both the stochastic and deterministic samplers, one might naturally wonder whether there is potential for achieving faster rates. In practice, the evaluation of Stein scores in each step often entails computing the output of a large neural network, thereby calling for new solutions to reduce the number of score evaluations without compromising sampling fidelity. Indeed, this has inspired a large strand of recent works focused on speeding up diffusion generative modeling. Towards this end, one prominent approach is *distillation*, which attempts to distill a pre-trained diffusion model into another model (e.g., progressive distillation, consistency model) that can be executed in significantly fewer steps (Luhman and Luhman, 2021; Salimans and Ho, 2021; Meng et al., 2023; Song et al., 2023). However, while distillation-based techniques have

achieved outstanding empirical performance, they often necessitate additional training processes, imposing high computational burdens beyond score matching. In contrast, an alternative route towards acceleration is “training-free,” which directly invokes the pre-trained diffusion model (particularly the pre-trained score functions) for sampling without requiring additional training processes. Examples of training-free accelerated samplers include DPM-Solver (Lu et al., 2022a), DPM-Solver++ (Lu et al., 2022b), DEIS (Zhang and Chen, 2022), UniPC (Zhao et al., 2023), the SA-Solver (Xue et al., 2023), among others, which leverage faster solvers for ODE and SDE using only the pre-trained score functions. Nevertheless, non-asymptotic convergence analyses for these methods remain largely absent, making it challenging to rigorize the degrees of acceleration compared to the non-accelerated results (Lee et al., 2023; Chen et al., 2022, 2023a; Li et al., 2023; Benton et al., 2023a). All of this leads to the following question that we aim to explore in this work:

Can we design a training-free deterministic (resp. stochastic) sampler that converges provably faster than the DDIM (resp. DDPM)?

1.3 Our contributions

In this paper, we answer the above question in the affirmative. Our main contributions can be summarized as follows.

- In the deterministic setting, we demonstrate how to speed up the ODE-based sampler (i.e., the DDIM-type sampler). The proposed sampler, which exploits some sort of momentum term to adjust the update rule, leverages insights from higher-order ODE approximation in discrete time and shares similar intuitions with the fast ODE-based sampler DPM-Solver-2 (Lu et al., 2022a). We establish non-asymptotic convergence guarantees for the accelerated DDIM-type sampler, showing that its iteration complexity scales proportionally to $1/\sqrt{\varepsilon}$ (up to log factor). This substantially improves upon the prior convergence theory for the original DDIM sampler (Li et al., 2023) (which has an iteration complexity proportional to $1/\varepsilon$).
- In the stochastic setting, we propose a novel sampling procedure to accelerate the SDE-based sampler (i.e., the DDPM-type sampler). For this new sampler, we establish an iteration complexity bound proportional to $1/\varepsilon$ (modulo some log factor), thus unveiling the superiority of the proposed sampler compared to the original DDPM sampler (recall that the original DDPM sampler has an iteration complexity proportional to $1/\varepsilon^2$ (Li et al., 2023; Chen et al., 2023a, 2022)).

In addition, two aspects of our theory are worth emphasizing: (i) our theory accommodates ℓ_2 -accurate score estimates, rather than requiring ℓ_∞ score estimation accuracy; (ii) our theory covers a fairly general family of target data distributions, without imposing stringent assumptions like log-concavity and smoothness on the target distributions.

1.4 Other related works

We now briefly discuss additional related works in the prior art.

Convergence of score-based generative models (SGMs). For stochastic samplers of SGMs, the convergence guarantees were initially provided by early works including but not limited to De Bortoli et al. (2021); Liu et al. (2022b); Pidstrigach (2022); Block et al. (2020); De Bortoli (2022); Wibisono and Yang (2022); Gao et al. (2023), which often faced issues of either being not quantitative or suffering from the curse of dimensionality. More recent research has advanced this field by relaxing the assumptions on the score function and achieving polynomial convergence rates (Lee et al., 2022, 2023; Chen et al., 2022, 2023a,b; Li et al., 2023; Benton et al., 2023a; Liang et al., 2024; Tang and Zhao, 2024b). Furthermore, theoretical insights into probability flow-based ODE samplers, though less abundant, have been explored in recent works (Chen et al., 2023c; Li et al., 2023; Chen et al., 2023b; Benton et al., 2023b; Gao and Zhu, 2024). Additionally, Tang and Zhao (2024a) provided a continuous-time sampling error guarantee for a novel class of contraction diffusion models. Gao and Zhu (2024) studies the convergence properties for general probability flow ODEs w.r.t. Wasserstein distances. Most recently, Chen and Ying (2024) makes a step towards the convergence analysis of discrete state space diffusion model. Note that this body of research primarily aims

to quantify the proximity between distributions generated by SGMs and the ground truth distributions, assuming availability of an accurate score estimation oracle. Interestingly, a very recent research by Li et al. (2024b) reveals that even SGMs with empirically optimized score functions might underperform due to strong memorization effects. Moreover, some works delve into other aspects of the theoretical understanding of diffusion models. Furthermore, Wu et al. (2024) investigated how diffusion guidance combined with DDPM and DDIM samplers influences the conditional sampling quality.

Fast sampling in diffusion models. A recent strand of works to achieve few-step sampling — or even one-step sampling — falls under the category of training-based samplers, primarily focused on knowledge distillation (Meng et al., 2023; Salimans and Ho, 2021; Song et al., 2023). This method aims to distill a pre-trained diffusion model into another model that can be executed in significantly fewer steps. The recent work Li et al. (2024a) provided a first attempt towards theoretically understanding the sampling efficiency of consistency models. Another line of works aims to design training-free samplers (Lu et al., 2022a,b; Zhao et al., 2023; Zhang and Chen, 2022; Liu et al., 2022a; Zhang et al., 2022), which addresses the efficiency issue by developing faster solvers for the reverse-time SDE or ODE without requiring other information beyond the pre-trained SGMs. In addition, Li et al. (2023); Liang et al. (2024) introduced accelerated samplers that require additional training pertaining to estimating Hessian information at each step. Furthermore, combining GAN with diffusion has shown to be an effective strategy to speed up the sampling process (Wang et al., 2022; Xiao et al., 2021).

1.5 Notation

Before continuing, we find it helpful to introduce some notational conventions to be used throughout this paper. Capital letters are often used to represent random variables/vectors/processes, while lowercase letters denote deterministic variables. When considering two probability measures P and Q , we define their total-variation (TV) distance as $\text{TV}(P, Q) := \frac{1}{2} \int |dP - dQ|$, and the Kullback-Leibler (KL) divergence as $\text{KL}(P \parallel Q) := \int (\log \frac{dP}{dQ}) dP$. We use $p_X(\cdot)$ and $p_{X|Y}(\cdot|\cdot)$ to denote the probability density function of a random vector X , and the conditional probability of X given Y , respectively. For matrices, $\|A\|$ and $\|A\|_F$ refer to the spectral norm and Frobenius norm of a matrix A , respectively. For vector-valued functions f , we use J_f or $\frac{\partial f}{\partial x}$ to represent the Jacobian matrix of f . Given two functions $f(d, T)$ and $g(d, T)$, we employ 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 indicate the existence of a universal constant $C_1 > 0$ such that for all d and T , $f(d, T) \leq C_1 g(d, T)$ (resp. $f(d, T) \geq C_1 g(d, T)$). The notation $f(d, T) \asymp g(d, T)$ indicates that both $f(d, T) \lesssim g(d, T)$ and $f(d, T) \gtrsim g(d, T)$ hold at once.

2 Problem settings

In this section, we formulate the problem, and introduce a couple of key assumptions.

2.1 Model and sampling process

Forward process. Consider the forward Markov process (1) in discrete time that starts from the target data distribution $X_0 \sim p_{\text{data}}$ in \mathbb{R}^d and proceeds as follows:

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} W_t, \quad t = 1, \dots, T, \quad (3)$$

where the W_t 's are independently drawn from $\mathcal{N}(0, I_d)$. This process is said to be “variance-preserving,” in the sense that the covariance $\text{Cov}(X_t) = I_d$ holds throughout if $\text{Cov}(X_0) = I_d$. Taking

$$\bar{\alpha}_t := \prod_{k=1}^t \alpha_k \quad \text{with } \alpha_t := 1 - \beta_t \quad (4)$$

for every $1 \leq t \leq T$, one can write

$$X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \bar{W}_t \quad \text{for } \bar{W}_t \sim \mathcal{N}(0, I_d). \quad (5)$$

Throughout the paper, we shall use $q_t(\cdot)$ or $p_{X_t}(\cdot)$ interchangeably to denote the probability density function (PDF) of X_t . While we shall concentrate on the discrete-time process in the current paper, we shall note that the forward process has also been commonly studied in the continuous-time limit through the following diffusion process

$$dX_t = -\frac{1}{2}\beta(t)X_t dt + \sqrt{\beta(t)}dW_t \quad (0 \leq t \leq T), \quad X_0 \sim p_{\text{data}} \quad (6)$$

for some function $\beta(t)$ related to the learning rate, where W_t is the standard Brownian motion.

Score functions and score estimates. A key ingredient that plays a pivotal role in the sampling process is the (Stein) score function, defined as the log marginal density of the forward process.

Definition 1 (Score function). The score function, denoted by $s_t^* : \mathbb{R}^d \rightarrow \mathbb{R}^d (1 \leq t \leq T)$, is defined as

$$s_t^*(X) := \nabla \log q_t(X) = -\frac{1}{1 - \bar{\alpha}_t} \int_{x_0} p_{X_0 | X_t}(x_0 | x)(x - \sqrt{\bar{\alpha}_t}x_0) dx_0. \quad (7)$$

Here, the last identity follows from standard properties about Gaussians; see, e.g., [Chen et al. \(2022\)](#). In most applications, we have no access to perfect score functions; instead, what we have available are certain estimates for the score functions, to be denoted by $\{s_t(\cdot)\}_{1 \leq t \leq T}$ throughout.

Data generation process. The sampling process is performed via careful construction of the reverse process (2) to ensure distributional proximity. Working backward from $t = T, \dots, 1$, we assume throughout that $Y_T \sim \mathcal{N}(0, I_d)$.

- *Deterministic sampler.* A deterministic sampler typically chooses Y_{t-1} for each t to be a function of $\{Y_t, \dots, Y_T\}$. For instance, the following construction

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \frac{1 - \alpha_t}{2} s_t(Y_t) \right), \quad t = T, \dots, 1 \quad (8)$$

can be viewed as a DDIM-type sampler in discrete time. Note that the DDIM sampler is intimately connected with the following ODE — called the probability flow ODE or the diffusion ODE — in the continuous-time limit:

$$d\tilde{Y}_t = -\frac{1}{2}\beta(t) \left(\tilde{Y}_t + \nabla \log q_t(\tilde{Y}_t) \right) dt, \quad \tilde{Y}_T \sim q_T, \quad (9)$$

which enjoys matching marginal distributions as the forward diffusion process (6) in the sense that $\tilde{Y}_t \stackrel{d}{=} X_t$ for all $0 \leq t \leq T$ ([Song et al., 2021](#)).

- *Stochastic sampler.* In contrast to the deterministic case, each Y_{t-1} is a function of not only $\{Y_t, \dots, Y_T\}$ but also an additional independent noise $Z_t \sim \mathcal{N}(0, I_d)$. One example is the following sampler:

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + (1 - \alpha_t)s_t(Y_t) + \sqrt{1 - \alpha_t}Z_t \right), \quad t = T, \dots, 1 \quad (10)$$

which is closely related to the DDPM sampler in discrete time. The design of DDPM draws inspiration from a well-renowned result in the SDE literature ([Anderson, 1982](#); [Haussmann and Pardoux, 1986](#)); namely, there exists a reverse-time SDE

$$d\hat{Y}_t = -\frac{1}{2}\beta(t) \left(\hat{Y}_t + 2\nabla \log q_t(\hat{Y}_t) \right) dt + \sqrt{\beta(t)}d\hat{Z}_t \quad (11)$$

with $\hat{Y}_T \sim q_T$ that exhibits the same marginals — $\hat{Y}_t \stackrel{d}{=} X_t$ for all t — as the forward diffusion process (6). Here, \hat{Z}_t indicates a backward standard Brownian motion.

2.2 Assumptions

Before moving on to our algorithms and theory, let us introduce several assumptions that shall be used multiple times in this paper. To begin with, we impose the following assumption on the target data distribution.

Assumption 1. *Suppose that X_0 is a continuous random vector, and obeys*

$$\mathbb{P}(\|X_0\|_2 \leq R = T^{c_R} \mid X_0 \sim p_{\text{data}}) = 1 \quad (12)$$

for some arbitrarily large constant $c_R > 0$.

In words, the size of X_0 is allowed to grow polynomially (with arbitrarily large constant degree) in the number of steps, which suffices to accommodate the vast majority of practical applications.

Next, we specify the learning rates $\{\beta_t\}$ (or $\{\alpha_t\}$) employed in the forward process (3). Throughout this paper, we select the same learning rate schedule as in Li et al. (2023), namely,

$$\beta_1 = 1 - \alpha_1 = \frac{1}{T^{c_0}}, \quad (13a)$$

$$\beta_t = 1 - \alpha_t = \frac{c_1 \log T}{T} \min \left\{ \beta_1 \left(1 + \frac{c_1 \log T}{T} \right)^t, 1 \right\}, \quad t > 1 \quad (13b)$$

for some large enough numerical constants $c_0, c_1 > 0$. In short, there are two phases here: at first β_t grows exponentially fast, and then stays unchanged after surpassing some threshold. This also resembles the learning rate choices recommended by Benton et al. (2023a).

Moreover, let us also introduce two assumptions regarding the accuracy of the score estimates $\{s_t\}$, which are adopted in Li et al. (2023). Here and throughout, we denote by

$$J_{s_t^*} = \frac{\partial s_t^*}{\partial x} \quad \text{and} \quad J_{s_t} = \frac{\partial s_t}{\partial x} \quad (14)$$

the Jacobian matrices of $s_t^*(\cdot)$ and $s_t(\cdot)$, respectively.

Assumption 2. *Suppose that the mean squared estimation error of the score estimates $\{s_t\}_{1 \leq t \leq T}$ obeys*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} \left[\|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2.$$

Assumption 3. *Suppose that $s_t(\cdot)$ is continuously differentiable for each $1 \leq t \leq T$, and that the Jacobian matrices associated with the score estimates $\{s_t\}_{1 \leq t \leq T}$ satisfy*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} \left[\|J_{s_t}(X) - J_{s_t^*}(X)\| \right] \leq \varepsilon_{\text{Jacobi}}.$$

In short, Assumption 2 is concerned with the ℓ_2 score estimation error averaged across all steps, whereas Assumption 3 is about the average discrepancy in the associated Jacobian matrices. It is worth noting that Assumption 3 will only be imposed when analyzing the convergence of deterministic samplers, and is completely unnecessary for the stochastic counterpart.

3 Algorithm and main theory

In this section, we put forward two accelerated samplers — an ODE-based algorithm and an SDE-based algorithm — and present convergence theory to confirm the acceleration compared with prior DDIM and DDPM approaches.

3.1 Accelerated ODE-based sampler

The first algorithm we propose is an accelerated variant of the ODE-based deterministic sampler. Specifically, starting from $Y_T \sim \mathcal{N}(0, I_d)$, the proposed discrete-time sampler adopts the following update rule:

$$Y_t^- = \Phi_t(Y_t), \quad Y_{t-1} = \Psi_t(Y_t, Y_t^-) \quad \text{for } t = T, \dots, 1 \quad (15a)$$

where the mappings $\Phi_t(\cdot)$ and $\Psi_t(\cdot, \cdot)$ are chosen to be

$$\Phi_t(x) = \sqrt{\alpha_{t+1}} \left(x - \frac{1 - \alpha_{t+1}}{2} s_t(x) \right), \quad (15b)$$

$$\Psi_t(x, y) = \frac{1}{\sqrt{\alpha_t}} \left(x + \frac{1 - \alpha_t}{2} s_t(x) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} (s_t(x) - \sqrt{\alpha_{t+1}} s_{t+1}(y)) \right), \quad (15c)$$

and we remind the reader that s_t is the score estimate. In contrast to the original DDIM-type solver (8), the proposed accelerated sampler enjoys two distinguishing features:

- In each iteration t , the proposed sampler computes a mid-point $Y_t^- = \Phi_t(Y_t)$ (cf. (15b)). As it turns out, this mid-point is designed as a prediction of the probability flow ODE at time $t + 1$ using Y_t .
- In contrast to (8), the proposed update rule $Y_{t-1} = \Psi_t(Y_t, Y_t^-)$ (see (15c)) includes an additional term that is a properly scaled version of $s_t(Y_t) - \sqrt{\alpha_{t+1}} s_{t+1}(Y_t^-)$. In some sense, this term can be roughly viewed as exploiting “momentum” in adjusting the original sampling rule.

Theoretical guarantees. Let us proceed to present our convergence theory and its implications for the proposed deterministic sampler.

Theorem 1. *Suppose that Assumptions 1, 2 and 3 hold. Then the proposed sampler (15) with the learning rate schedule (13b) satisfies*

$$\text{TV}(q_1, p_1) \leq C_1 \frac{d^6 \log^6 T}{T^2} + C_1 \sqrt{d \log^3 T} \varepsilon_{\text{score}} + C_1 (d \log T) \varepsilon_{\text{Jacobi}} \quad (16)$$

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

We now take a moment to discuss the implications about this theorem.

- *Iteration complexity.* When the target accuracy level ε is small enough, the number of iterations needed to yield $\text{TV}(q_1, p_1) \leq \varepsilon$ is no larger than

$$(\text{iteration complexity}) \quad \frac{\text{poly}(d)}{\sqrt{\varepsilon}}, \quad (17)$$

ignoring any logarithmic factor in $1/\varepsilon$. Clearly, the dependency on $1/\varepsilon$ substantially improves upon the vanilla DDIM sampler, the latter of which has an iteration complexity proportional to $1/\varepsilon$ (Li et al., 2023).

- *Stability vis-a-vis score errors.* The discrepancy between the distribution of Y_1 and the target distribution of X_1 is proportional to the ℓ_2 score estimation error $\varepsilon_{\text{score}}$ defined in Assumption 2, as well as the Jacobian error $\varepsilon_{\text{Jacobi}}$ defined in Assumption 3. It is worth noting, however, that the same result might not hold if we remove Assumption 3. More specifically, when only score estimation accuracy is assumed, the deterministic sampler is not guaranteed to achieve small TV error; see Li et al. (2023) for an illustrative example.

Interpretation via second-order ODE. In order to help elucidate the rationale of the proposed sampler, we make note of an intimate connection between (15) and high-order ODE, the latter of which has facilitated the design of fast deterministic samplers (e.g., DPM-Solver (Lu et al., 2022a)).

In view of the relation (5), for any $0 < \gamma < 1$, let us first abuse the notation and introduce

$$X(\gamma) \stackrel{d}{=} \sqrt{\gamma}X_0 + \sqrt{1-\gamma}Z, \quad Z \sim \mathcal{N}(0, I_d) \quad (18a)$$

$$s_\gamma^*(X) := \nabla_X \log p_{X(\gamma)}(X). \quad (18b)$$

We further consider the following continuous-time analog $\bar{\alpha}(t)$ of the discrete learning rate $\bar{\alpha}_t$ (cf. (4)):

$$\frac{d\bar{\alpha}(t)}{dt} = -\beta(t)\bar{\alpha}(t), \quad \bar{\alpha}(T) = \bar{\alpha}_T. \quad (18c)$$

Given that the probability flow ODE (9) yields identical marginal distributions as the forward process X_t (cf. (6)) for every t , invoking (18c), we can easily see that $X(\bar{\alpha}(t)) \stackrel{d}{=} X_t$ can be generated as follows:

$$\frac{dX(\bar{\alpha}(t))}{d\bar{\alpha}(t)} = \frac{1}{2\bar{\alpha}(t)} \left(X(\bar{\alpha}(t)) + s_{\bar{\alpha}(t)}^*(X(\bar{\alpha}(t))) \right), \quad X(\bar{\alpha}(T)) \sim q_T, \quad (19)$$

By taking $f(\gamma) = \frac{1}{\sqrt{\gamma}}X(\gamma)$, we can apply (19) to derive

$$\frac{df(\gamma)}{d\gamma} = -\frac{1}{2\sqrt{\gamma^3}}X(\gamma) + \frac{1}{\sqrt{\gamma}} \frac{dX(\gamma)}{d\gamma} = \frac{1}{2\sqrt{\gamma^3}}s_\gamma^*(X(\gamma)).$$

This taken together with $\bar{\alpha}_t = \bar{\alpha}_{t-1}\alpha_t$ (cf. (4)) immediately implies that

$$\begin{aligned} \frac{1}{\sqrt{\bar{\alpha}_{t-1}}}X(\bar{\alpha}_{t-1}) &= \frac{1}{\sqrt{\bar{\alpha}_t}}X(\bar{\alpha}_t) + \frac{1}{2} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}}s_\gamma^*(X(\gamma))d\gamma, \\ \implies X(\bar{\alpha}_{t-1}) &= \frac{1}{\sqrt{\alpha_t}}X(\bar{\alpha}_t) + \frac{\sqrt{\bar{\alpha}_{t-1}}}{2} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}}s_\gamma^*(X(\gamma))d\gamma. \end{aligned} \quad (20)$$

With this relation in mind, we are ready to discuss the following approximation in discrete time:

- *Scheme 1:* If we approximate $s_\gamma^*(X(\gamma))$ for $\gamma \in [\bar{\alpha}_t, \bar{\alpha}_{t-1}]$ by $s_\gamma^*(X(\gamma)) \approx s_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t)) \approx s_t(X_t)$, then we arrive at

$$\begin{aligned} X(\bar{\alpha}_{t-1}) &\approx \frac{1}{\sqrt{\alpha_t}}X(\bar{\alpha}_t) + \left(\frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} - 1 \right) s_t(X_t) \\ &\approx \frac{1}{\sqrt{\alpha_t}} \left\{ X(\bar{\alpha}_t) + \frac{1-\alpha_t}{2} s_t(X_t) \right\}, \end{aligned}$$

where we use the facts that $\bar{\alpha}_t/\bar{\alpha}_{t-1} = \alpha_t$ and $\alpha_t \approx 1$. This coincides with the deterministic sampler (8).

- *Scheme 2:* If we invoke a more refined approximation for $s_\gamma^*(X(\gamma))$ as

$$s_\gamma^*(X(\gamma)) \approx s_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t)) + \frac{ds_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t))}{d\gamma} (\gamma - \bar{\alpha}_t) \quad (21)$$

$$\begin{aligned} &\approx s_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t)) + \frac{\gamma - \bar{\alpha}_t}{\bar{\alpha}_t - \bar{\alpha}_{t+1}} \left(s_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t)) - s_{\bar{\alpha}_{t+1}}^*(X(\bar{\alpha}_{t+1})) \right) \\ &\approx s_t(X_t) + \frac{\gamma - \bar{\alpha}_t}{\bar{\alpha}_t - \bar{\alpha}_{t+1}} (s_t(X_t) - s_{t+1}(X_{t+1})), \end{aligned} \quad (22)$$

then (20) can be approximated by

$$X(\bar{\alpha}_{t-1})$$

$$\begin{aligned}
&\approx \frac{1}{\sqrt{\alpha_t}} X(\bar{\alpha}_t) + \frac{\sqrt{\bar{\alpha}_{t-1}} s_t(X_t)}{2} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} d\gamma + \frac{\sqrt{\bar{\alpha}_{t-1}} (s_t(X_t) - s_{t+1}(X_{t+1}))}{2(\bar{\alpha}_t - \bar{\alpha}_{t+1})} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{\gamma - \bar{\alpha}_t}{\sqrt{\alpha^3}} d\gamma \\
&\approx \frac{1}{\sqrt{\alpha_t}} \left\{ X(\bar{\alpha}_t) + \frac{1 - \alpha_t}{2} s_t(X_t) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} (s_t(X_t) - \sqrt{\alpha_{t+1}} s_{t+1}(X_{t+1})) \right\}, \tag{23}
\end{aligned}$$

which resembles the proposed sampler (15), and is computationally more appealing since it reuses the previous score function evaluation.

It is worth noting that similar approximation as in Scheme 2 has been invoked previously in Lu et al. (2022a, Eqn (3.6)) to construct high-order ODE solvers (e.g., the DPM-Solver-2, with 2 indicating second-order ODEs). Consequently, the acceleration achieved by our sampler is achieved through ideas akin to the second-order ODE; in turn, our convergence guarantees shed light on the effectiveness of high-order ODE solvers like the popular DPM-Solver.

3.2 Accelerated SDE-based sampler

Next, we turn to stochastic samplers, and propose a new stochastic sampling procedure that enjoys improved convergence guarantees compared to the DDPM-type sampler (10). To be precise, the proposed sampler begins by drawing $Y_T \sim \mathcal{N}(0, I_d)$ and adopts the following update rule:

$$Y_t^+ = \Phi_t(Y_t, Z_t), \quad Y_{t-1} = \Psi_t(Y_t^+, Z_t^+) \tag{24a}$$

for $t = T, \dots, 1$, where $Z_t, Z_t^+ \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$, and

$$\Phi_t(x, z) = x + \sqrt{\frac{1 - \alpha_t}{2}} z, \tag{24b}$$

$$\Psi_t(y, z) = \frac{1}{\sqrt{\alpha_t}} \left(y + (1 - \alpha_t) s_t(y) + \sqrt{\frac{1 - \alpha_t}{2}} z \right). \tag{24c}$$

The key difference between the proposed sampler and the original DDPM-type sampler lies in the additional operation $\Phi_t(\cdot, \cdot)$. In this step, a random noise Z_t is injected into the current sample Y_t to obtain an intermediate point Y_t^+ , which together with another random noise Z_t^+ is subsequently fed into $\Psi_t(\cdot, \cdot)$ — a mapping identical to (10).

Theoretical guarantees. Let us present the convergence guarantees of the proposed stochastic sampler and their implications, followed by some interpretation of the design rationale of the algorithm.

Theorem 2. *Suppose that Assumptions 1 and 2 hold. Then the proposed stochastic sampler (24) with the learning rate schedule (13b) achieves*

$$\text{TV}(q_1, p_1) \leq \sqrt{\frac{1}{2} \text{KL}(q_1 \parallel p_1)} \leq C_1 \frac{d^3 \log^{4.5} T}{T} + C_1 \sqrt{d} \varepsilon_{\text{score}} \log^{1.5} T \tag{25}$$

for some universal constant $C_1 > 0$.

Theorem 2 provides non-asymptotic characterizations for the data generation quality of the accelerated stochastic sampler. In comparison with the convergence theory for the DDPM-type sampler — which has a convergence rate proportional to $1/\sqrt{T}$ (Chen et al., 2022, 2023a; Li et al., 2023; Benton et al., 2023a) — Theorem 2 asserts that the proposed accelerated sampler achieves a faster convergence rate proportional to $1/T$. In contrast to Theorem 1 for the ODE-based sampler, the SDE-based sampler does not require continuity of the Jacobian matrix (i.e., Assumption 3). As before, the total-variation distance between X_1 and Y_1 is proportional to the ℓ_2 score estimation error when T is sufficiently large, which covers a broad range of target data distributions with no requirement on the smoothness or log-concavity of the data distribution.

Interpretation via higher-order approximation. Now we provide some insights into the motivation of the proposed sampler. We start with the characterizations of conditional density $p_{X_{t-1}|X_t}$. Denoting $\mu_t^*(x_t) := \frac{1}{\sqrt{\alpha_t}}(x_t + (1 - \alpha_t)s_t^*(x_t))$ and $J_t(x_t) = -(1 - \bar{\alpha}_t)J_{s_t^*}(x_t)$, we can approximate $p_{X_{t-1}|X_t}$ by

$$p_{X_{t-1}|X_t}(x_{t-1} | x_t) \approx \exp\left(-\frac{\alpha_t}{2(1 - \alpha_t)} \cdot \left\| \left(I - \frac{1 - \alpha_t}{2(\alpha_t - \bar{\alpha}_t)} J_t(x_t)\right)^{-1} (x_{t-1} - \mu_t^*(x_t)) \right\|^2\right). \quad (26)$$

which is tighter than the one used in analysis of the original SDE-based sampler (Li et al., 2023) by adopting a higher-order expansion. This in turn motivates us to consider the following sequence

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\underbrace{Y_t + \sqrt{\frac{1 - \alpha_t}{2}} Z_t + \sqrt{\frac{1 - \alpha_t}{2}} Z_t^+}_{\Phi(Y_t, Z_t)} + \underbrace{(1 - \alpha_t)s_t^*(Y_t) - \frac{(1 - \alpha_t)^{3/2}}{\sqrt{2}(\alpha_t - \bar{\alpha}_t)} J_t(Y_t) Z_t}_{\approx s_t^*(\Phi(Y_t, Z_t))} \right)$$

with $Z_t, Z_t^+ \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$, and $p_{Y_{t-1}|Y_t}(x_{t-1} | x_t)$ follows

$$\mathcal{N}\left(\mu_t^*(x_t), \frac{1 - \alpha_t}{\alpha_t} \left(I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(x_t)\right) \left(I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(x_t)\right)^\top\right)$$

which aligns with (26). In addition, if we further utilize $(1 - \alpha_t)s_t^*(Y_t) - \frac{(1 - \alpha_t)^{3/2}}{\sqrt{2}(\alpha_t - \bar{\alpha}_t)} J_t(Y_t) Z_t$ as a first-order approximation of $s_t^*(Y_t + \sqrt{\frac{1 - \alpha_t}{2}} Z_t)$, we can then arrive at the update rule of the proposed sampler in (24).

4 Experiments

In this section, we illustrate the performance of the proposed accelerated samplers, focusing on highlighting the relative comparisons with respect to the original DDIM/DDPM ones using the same pre-trained score functions. As an initial step, we focus on reporting result for the deterministic samplers, leaving the stochastic samplers to future work.

4.1 Practical implementation

In practice, the pre-trained score functions are often available in the form of noise-prediction networks $\epsilon_t(\cdot)$, which are connected via the following relationship in view of (7):

$$s_t^*(X) := -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t^*(X), \quad (27)$$

and $\epsilon_t(\cdot)$ is the estimate of $\epsilon_t^*(\cdot)$. To better align with the empirical practice, it is judicious that the integration in (20) be approximated in terms of $\epsilon_t^*(X)$, leading to an equivalent rewrite as

$$X(\bar{\alpha}_{t-1}) = \frac{1}{\sqrt{\alpha_t}} X(\bar{\alpha}_t) - \frac{\sqrt{\bar{\alpha}_{t-1}}}{2} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3 \sqrt{1 - \gamma}}} \epsilon_\gamma^*(X(\gamma)) d\gamma.$$

Following similar discussions in Section 3.1, we discuss its first-order and second-order approximations in discrete time.

- *Scheme 1:* If we approximate $\epsilon_\gamma^*(X(\gamma))$ for $\gamma \in [\bar{\alpha}_t, \bar{\alpha}_{t-1}]$ by $\epsilon_\gamma^*(X(\gamma)) \approx \epsilon_{\bar{\alpha}_t}^*(X(\bar{\alpha}_t)) \approx \epsilon_t(X_t)$, then we arrive at

$$X(\bar{\alpha}_{t-1}) \approx \frac{1}{\sqrt{\alpha_t}} X(\bar{\alpha}_t) + \left(\sqrt{1 - \bar{\alpha}_{t-1}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} \right) \epsilon_t(X_t), \quad (28)$$

which matches exactly with the DDIM sampler in Song et al. (2020).

- *Scheme 2:* If we invoke the refined approximation (22) in terms of $\epsilon_\gamma^*(X(\gamma))$, we have

$$X(\bar{\alpha}_{t-1}) \approx \frac{1}{\sqrt{\alpha_t}} X(\bar{\alpha}_t) - \frac{\sqrt{\bar{\alpha}_{t-1}} \epsilon_t(X_t)}{2} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3(1-\gamma)}} d\gamma \\ - \frac{\sqrt{\bar{\alpha}_{t-1}} (\epsilon_t(X_t) - \epsilon_{t+1}(X_{t+1}))}{2(\bar{\alpha}_t - \bar{\alpha}_{t+1})} \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{(\gamma - \bar{\alpha}_t)}{\sqrt{\gamma^3(1-\gamma)}} d\gamma,$$

which after integration becomes:

$$X(\bar{\alpha}_{t-1}) \approx \frac{1}{\sqrt{\alpha_t}} X(\bar{\alpha}_t) + \left(\sqrt{1 - \bar{\alpha}_{t-1}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} \right) \epsilon_t(X_t) \\ + \left(\frac{\sqrt{\bar{\alpha}_{t-1}}}{\bar{\alpha}_t - \bar{\alpha}_{t+1}} \right) \left(\bar{\alpha}_t \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_{t-1}}} + \arcsin \sqrt{\bar{\alpha}_{t-1}} - \bar{\alpha}_t \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} - \arcsin \sqrt{\bar{\alpha}_t} \right) (\epsilon_{t+1}(X_{t+1}) - \epsilon_t(X_t)). \quad (29)$$

This is our new sampler for implementation.

4.2 Experimental results

We use pre-trained score functions from Huggingface (von Platen et al., 2022) for three datasets: CelebA-HQ, LSUN-Bedroom and LSUN-Churches. The same score functions are used in all the samplers. Note that we have not attempted to optimize the speed nor the performance using additional tricks, e.g., employing better score functions, but aim to corroborate our theoretical findings regarding the acceleration of the new samplers without training additional functions when the implementations are otherwise kept the same.

We first compare the vanilla DDIM-type sampler (cf. (28)) and the accelerated DDIM-type sampler (cf. (29)). To begin, Figure 1 illustrates the progress of the generated samples over different numbers of function evaluation (NFEs) (between 4 and 50) from the same random seed, using pre-trained scores from the LSUN-Churches dataset. Here, the NFE is the same as the number of diffusion steps since each step takes one score evaluation.



Figure 1: The progress of the generated samples over different numbers of NFEs (from 4 to 50), using pre-trained scores from the LSUN-Churches dataset. Top row: the vanilla DDIM-type sampler. Bottom row: the accelerated DDIM-type sampler (ours).

To further demonstrate the quality of the sampled images, Figure 2 provides examples of sampled images from the DDIM-type samplers, using pre-trained scores from CelebA-HQ, LSUN-Bedroom and LSUN-Churches datasets, respectively. It can be seen that the sampled images are crisper and less noisy from the accelerated DDIM-type sampler, compared with from the original one, indicating the effectiveness of our method.

5 Discussion

In this paper, we have developed novel strategies to achieve provable acceleration in score-based generative modeling. The proposed deterministic sampler achieves a convergence rate $1/T^2$ that substantially improves



Figure 2: Examples of sampled images from the DDIM-type samplers with 5 NFEs, using pre-trained scores from the LSUN-Churches, LSUN-Bedroom, and CelebA-HQ datasets. For each dataset, the top image is the original DDIM-type sampler, and the bottom image is the accelerated DDIM-type sampler (ours).

upon prior theory for the probability flow ODE approach, whereas the proposed stochastic sampler enjoys a converge rate $1/T$ that also significantly outperforms the convergence theory for the DDPM-type sampler. We have demonstrated the stability of these samplers, establishing non-asymptotic theoretical guarantees that hold in the presence of ℓ_2 -accurate score estimates. Our algorithm development for the deterministic case draws inspiration from higher-order ODE approximations in discrete time, which might shed light on understanding popular ODE-based samplers like the DPM-Solver. In comparison, the accelerated stochastic sampler is designed based on higher-order expansions of the conditional density.

Our findings further suggest multiple directions that are worthy of future exploration. For instance, our convergence theory remains sub-optimal in terms of the dependency on the problem dimension d , which calls for a more refined theory to sharpen dimension dependency. Additionally, given the conceptual similarity between our accelerated deterministic sampler and second-order ODE, it would be interesting to extend the algorithm and theory using ideas arising from third-order or even higher-order ODE. In particular, third-order ODE has been implemented in DPM-Solver-3, which is among the most effective DPM-Solvers in practice. Finally, it would be important to design higher-order solvers for SDE-based samplers, in order to unveil the degree of acceleration that can be achieved through high-order SDE.

Acknowledgements

Y. Wei is supported in part by the NSF grants DMS-2147546/2015447, CAREER award DMS-2143215, CCF-2106778, and the Google Research Scholar Award. The work of T. Efimov and Y. Chi is supported in part by the grants ONR N00014-19-1-2404, NSF DMS-2134080, ECCS-2126634 and FHWA 693JJ321C000013. 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.

References

- Anderson, B. D. (1982). Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326.
- Benton, J., De Bortoli, V., Doucet, A., and Deligiannidis, G. (2023a). Linear convergence bounds for diffusion models via stochastic localization. *arXiv preprint arXiv:2308.03686*.
- Benton, J., Deligiannidis, G., and Doucet, A. (2023b). Error bounds for flow matching methods. *arXiv preprint arXiv:2305.16860*.
- 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. (2023a). Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pages 4735–4763.
- Chen, H. and Ying, L. (2024). Convergence analysis of discrete diffusion model: Exact implementation through uniformization. *arXiv preprint arXiv:2402.08095*.
- Chen, S., Chewi, S., Lee, H., Li, Y., Lu, J., and Salim, A. (2023b). The probability flow ODE is provably fast. *Neural Information Processing Systems*.
- Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., and Zhang, A. R. (2022). 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. (2023c). Restoration-degradation beyond linear diffusions: A non-asymptotic analysis for DDIM-type samplers. In *International Conference on Machine Learning*, pages 4462–4484.
- 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.
- Gao, X., Nguyen, H. M., and Zhu, L. (2023). Wasserstein convergence guarantees for a general class of score-based generative models. *arXiv preprint arXiv:2311.11003*.
- Gao, X. and Zhu, L. (2024). Convergence analysis for general probability flow ODEs of diffusion models in Wasserstein distances. *arXiv preprint arXiv:2401.17958*.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2020). Generative adversarial networks. *Communications of the ACM*, 63(11):139–144.
- Guo, Z., Liu, J., Wang, Y., Chen, M., Wang, D., Xu, D., and Cheng, J. (2023). Diffusion models in bioinformatics and computational biology. *Nature Reviews Bioengineering*, pages 1–19.
- Hausmann, U. G. and Pardoux, E. (1986). Time reversal of diffusions. *The Annals of Probability*, pages 1188–1205.
- 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.
- Kazerouni, A., Aghdam, E. K., Heidari, M., Azad, R., Fayyaz, M., Hacihaliloglu, I., and Merhof, D. (2023). Diffusion models in medical imaging: A comprehensive survey. *Medical Image Analysis*, page 102846.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. *International Conference on Learning Representations*.
- 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., Huang, Z., and Wei, Y. (2024a). Towards a mathematical theory for consistency training in diffusion models. *arXiv preprint arXiv:2402.07802*.
- Li, G., Wei, Y., Chen, Y., and Chi, Y. (2023). Towards faster non-asymptotic convergence for diffusion-based generative models. *arXiv preprint arXiv:2306.09251*.
- Li, S., Chen, S., and Li, Q. (2024b). A good score does not lead to a good generative model. *arXiv preprint arXiv:2401.04856*.
- Liang, Y., Ju, P., Liang, Y., and Shroff, N. (2024). Non-asymptotic convergence of discrete-time diffusion models: New approach and improved rate. *arXiv preprint arXiv:2402.13901*.
- Liu, L., Ren, Y., Lin, Z., and Zhao, Z. (2022a). Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*.
- Liu, X., Wu, L., Ye, M., and Liu, Q. (2022b). Let us build bridges: Understanding and extending diffusion generative models. *arXiv preprint arXiv:2208.14699*.
- Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J. (2022a). DPM-Solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787.
- Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J. (2022b). DPM-Solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*.
- Luhman, E. and Luhman, T. (2021). Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*.
- Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S., Ho, J., and Salimans, T. (2023). On distillation of guided diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306.
- Pang, T., Xu, K., Li, C., Song, Y., Ermon, S., and Zhu, J. (2020). Efficient learning of generative models via finite-difference score matching. *Advances in Neural Information Processing Systems*, 33:19175–19188.
- Pidstrigach, J. (2022). Score-based generative models detect manifolds. *arXiv preprint arXiv:2206.01018*.
- Salimans, T. and Ho, J. (2021). Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*.
- 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. (2020). Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*.
- Song, Y., Dhariwal, P., Chen, M., and Sutskever, I. (2023). Consistency models. In *International Conference on Machine Learning*.
- Song, Y. and Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*.
- Tang, W. and Zhao, H. (2024a). Contractive diffusion probabilistic models. *arXiv preprint arXiv:2401.13115*.

- Tang, W. and Zhao, H. (2024b). Score-based diffusion models via stochastic differential equations—a technical tutorial. *arXiv preprint arXiv:2402.07487*.
- Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674.
- von Platen, P., Patil, S., Lozhkov, A., Cuenca, P., Lambert, N., Rasul, K., Davaadorj, M., and Wolf, T. (2022). Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>.
- Wang, Z., Zheng, H., He, P., Chen, W., and Zhou, M. (2022). Diffusion-gan: Training gans with diffusion. *arXiv preprint arXiv:2206.02262*.
- 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*.
- Wu, Y., Chen, M., Li, Z., Wang, M., and Wei, Y. (2024). Theoretical insights for diffusion guidance: A case study for gaussian mixture models. *arXiv preprint arXiv:arXiv:2403.01639*.
- Xiao, Z., Kreis, K., and Vahdat, A. (2021). Tackling the generative learning trilemma with denoising diffusion gans. *arXiv preprint arXiv:2112.07804*.
- Xue, S., Yi, M., Luo, W., Zhang, S., Sun, J., Li, Z., and Ma, Z.-M. (2023). SA-Solver: Stochastic Adams solver for fast sampling of diffusion models. *arXiv preprint arXiv:2309.05019*.
- Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Zhang, W., Cui, B., and Yang, M.-H. (2023). Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39.
- Zhang, Q. and Chen, Y. (2022). Fast sampling of diffusion models with exponential integrator. *arXiv preprint arXiv:2204.13902*.
- Zhang, Q., Tao, M., and Chen, Y. (2022). gddim: Generalized denoising diffusion implicit models. *arXiv preprint arXiv:2206.05564*.
- Zhao, W., Bai, L., Rao, Y., Zhou, J., and Lu, J. (2023). UniPC: A unified predictor-corrector framework for fast sampling of diffusion models. *arXiv preprint arXiv:2302.04867*.

A Preliminaries

Before delving into the proof, we make note of a couple of preliminary facts, primarily from [Li et al. \(2023\)](#).

A.1 Basic facts

Score functions. We first give some characterizations of the score function, which follow from [Li et al. \(2023\)](#), properties (38).

Lemma 1. *The true score function s_t^* is given by the conditional expectation below:*

$$\begin{aligned}
 s_t^*(x) &= \mathbb{E} \left[-\frac{1}{\sqrt{1-\bar{\alpha}_t}} W \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} W = x \right] = \frac{1}{1-\bar{\alpha}_t} \mathbb{E} [\sqrt{\bar{\alpha}_t} X_0 - x \mid \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} W = x] \\
 &= -\frac{1}{1-\bar{\alpha}_t} \underbrace{\int_{x_0} (x - \sqrt{\bar{\alpha}_t} x_0) p_{X_0|X_t}(x_0|x) dx_0}_{=: g_t(x)}.
 \end{aligned} \tag{30}$$

Also, the Jacobian matrix

$$J_t(x) := \frac{\partial g_t(x)}{\partial x} \tag{31}$$

associated with the function $g_t(\cdot)$ (defined in (30)) satisfies

$$J_t(x) = I_d + \frac{1}{1 - \bar{\alpha}_t} \left\{ \mathbb{E}[X_t - \sqrt{\bar{\alpha}_t} X_0 \mid X_t = x] \left(\mathbb{E}[X_t - \sqrt{\bar{\alpha}_t} X_0 \mid X_t = x] \right)^\top - \mathbb{E} \left[(X_t - \sqrt{\bar{\alpha}_t} X_0)(X_t - \sqrt{\bar{\alpha}_t} X_0)^\top \mid X_t = x \right] \right\}. \quad (32)$$

Learning rates. The learning rates $\{\alpha_t\}$ as specified in (13b) enjoy several properties that will be used multiple times in the analysis. We record several of these properties below, which have been proven in Li et al. (2023, properties (39)).

Lemma 2. *The learning rates specified in (13b) obey*

$$\alpha_t \geq 1 - \frac{c_1 \log T}{T} \geq \frac{1}{2}, \quad 1 \leq t \leq T \quad (33a)$$

$$\frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} \leq \frac{1 - \alpha_t}{2\alpha_t - \bar{\alpha}_t} \leq \frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}} \leq \frac{4c_1 \log T}{T}, \quad 2 \leq t \leq T \quad (33b)$$

$$1 \leq \frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}} \leq 1 + \frac{4c_1 \log T}{T}, \quad 2 \leq t \leq T \quad (33c)$$

$$\bar{\alpha}_T \leq \frac{1}{T^{c_2}}, \quad (33d)$$

provided that T is large enough. Here, c_1 is defined in (13b), and $c_2 \geq 1000$ is some large numerical constant. In addition, if $\frac{d(1-\alpha_t)}{\alpha_t - \bar{\alpha}_t} \lesssim 1$, then one has

$$\left(\frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} = 1 + \frac{d(1 - \alpha_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{d(d-2)(1 - \alpha_t)^2}{8(\alpha_t - \bar{\alpha}_t)^2} + O\left(d^3 \left(\frac{1 - \alpha_t}{\alpha_t - \bar{\alpha}_t} \right)^3\right). \quad (33e)$$

The forward process. Next, we gather several conditional tail bounds for the random vector X_0 of the forward process, which have been established in Li et al. (2023, Lemmas 1 and 2).

Lemma 3. *Suppose that there exists some numerical constant $c_R > 0$ obeying*

$$\mathbb{P}(\|X_0\|_2 \leq R) = 1 \quad \text{and} \quad R = T^{c_R}. \quad (34)$$

Consider any $y \in \mathbb{R}$, and let

$$\theta(y) := \max \left\{ \frac{-\log p_{X_t}(y)}{d \log T}, c_6 \right\} \quad (35)$$

for some large enough constant $c_6 \geq 2c_R + c_0$. Then for any quantity $c_5 \geq 2$, conditioned on $X_t = y$ one has

$$\|\sqrt{\bar{\alpha}_t} X_0 - y\|_2 \leq 5c_5 \sqrt{\theta(y) d(1 - \bar{\alpha}_t) \log T} \quad (36)$$

with probability at least $1 - \exp(-c_5^2 \theta(y) d \log T)$. In addition, it holds that

$$\mathbb{E} \left[\|\sqrt{\bar{\alpha}_t} X_0 - y\|_2 \mid X_t = y \right] \leq 12 \sqrt{\theta(y) d(1 - \bar{\alpha}_t) \log T}, \quad (37a)$$

$$\mathbb{E} \left[\|\sqrt{\bar{\alpha}_t} X_0 - y\|_2^2 \mid X_t = y \right] \leq 120 \theta(y) d(1 - \bar{\alpha}_t) \log T, \quad (37b)$$

$$\mathbb{E} \left[\|\sqrt{\bar{\alpha}_t} X_0 - y\|_2^3 \mid X_t = y \right] \leq 1040 (\theta(y) d(1 - \bar{\alpha}_t) \log T)^{3/2}, \quad (37c)$$

$$\mathbb{E} \left[\|\sqrt{\bar{\alpha}_t} X_0 - y\|_2^4 \mid X_t = y \right] \leq 10080 (\theta(y) d(1 - \bar{\alpha}_t) \log T)^2. \quad (37d)$$

Lemma 4. *For some $\theta > c_6$, assume that $\|y - x\|_2 \lesssim (1 - \alpha_{t+1}) \sqrt{\frac{\theta d \log T}{1 - \bar{\alpha}_t}}$ and $\log p_{X_t}(x) \geq -\theta d \log T$. Then we have*

$$p_{X_0 \mid X_{t+1}}(x_0 \mid \sqrt{\bar{\alpha}_{t+1}} y) = p_{X_0 \mid X_t}(x_0 \mid x) \left\{ 1 + \frac{(1 - \alpha_{t+1}) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\bar{\alpha}_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} \right\}$$

$$\int_{x_0} \left(\frac{(1 - \alpha_{t+1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} \right) p_{X_0 | X_t}(x_0 | x) dx_0 + O\left(\theta d \left(\frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_t}\right)^{3/2} \log T\right). \quad (38)$$

Proof. See Section A.2. \square

Proximity of p_{X_T} and q_{Y_T} . When the number of steps T is sufficiently large, the distribution of p_{X_T} and that of p_{Y_T} become exceedingly close, as asserted by the following lemma (see Li et al. (2023, Lemma 3)).

Lemma 5. For any large enough T , one has

$$\text{TV}(p_{X_T} \parallel p_{Y_T})^2 \leq \frac{1}{2} \text{KL}(p_{X_T} \parallel p_{Y_T}) \lesssim \frac{1}{T^{200}}. \quad (39)$$

Score estimation errors. Consider any vector $x \in \mathbb{R}^d$. For any $1 < t \leq T$, define

$$\varepsilon_{\text{score},t}(x) := \|s_t(x) - s_t^*(x)\|_2 \quad \text{and} \quad \varepsilon_{\text{Jacobi},t}(x) := \|J_{s_t}(x) - J_{s_t^*}(x)\|, \quad (40)$$

where we use J_{s_t} and $J_{s_t^*}$ to represent the Jacobian matrices of $s_t(\cdot)$ and $s_t^*(\cdot)$, respectively. We also have, under Assumptions 2 and 3, that

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} [\varepsilon_{\text{score},t}(X)] \leq \left(\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} [\varepsilon_{\text{score},t}(X)^2] \right)^{1/2} \leq \varepsilon_{\text{score}}, \quad (41a)$$

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim q_t} [\varepsilon_{\text{Jacobi},t}(X)] \leq \varepsilon_{\text{Jacobi}}. \quad (41b)$$

A.2 Proof of Lemma 4

To establish this lemma, we observe that

$$\begin{aligned} p_{X_0 | X_{t+1}}(x_0 | \sqrt{\alpha_{t+1}} y) &= \frac{p_{X_0}(x_0) p_{X_{t+1} | X_0}(y | x_0)}{\int_{x_0} p_{X_0}(x_0) p_{X_{t+1} | X_0}(y | x_0) dx_0} \\ &= \frac{p_{X_0}(x_0) \exp\left(-\frac{\|y - \sqrt{\alpha_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)}\right)}{\int_{x_0} p_{X_0}(x_0) \exp\left(-\frac{\|y - \sqrt{\alpha_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)}\right) dx_0} \\ &= \frac{p_{X_0}(x_0) \exp\left(-\frac{\|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)}\right) \left(1 - \frac{(1 - \alpha_{t+1}^{-1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} + O\left(\theta d \left(\frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_t}\right)^{3/2} \log T\right)\right)}{\int_{x_0} p_{X_0}(x_0) \exp\left(-\frac{\|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)}\right) \left(1 - \frac{(1 - \alpha_{t+1}^{-1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} + O\left(\theta d \left(\frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_t}\right)^{3/2} \log T\right)\right) dx_0} \\ &= p_{X_0 | X_t}(x_0 | x) \left\{ 1 + \frac{(1 - \alpha_{t+1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} \right. \\ &\quad \left. - \int_{x_0} \left(\frac{(1 - \alpha_{t+1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} - \frac{(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{1 - \bar{\alpha}_t} \right) p_{X_0 | X_t}(x_0 | x) dx_0 \right. \\ &\quad \left. + O\left(\theta d \left(\frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_t}\right)^{3/2} \log T\right) \right\}, \end{aligned}$$

which follows from the following property:

$$\begin{aligned} &\frac{\|y - \sqrt{\alpha_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)} - \frac{\|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)} \\ &= \frac{(1 - \alpha_{t+1}^{-1}) \|y - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)(\alpha_{t+1}^{-1} - \bar{\alpha}_t)} + \frac{\|y - \sqrt{\alpha_t} x_0\|_2^2 - \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)} \\ &= -\frac{(1 - \alpha_{t+1}) \|x - \sqrt{\alpha_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)^2} + \frac{2(x - \sqrt{\alpha_t} x_0)^\top (y - x)}{2(1 - \bar{\alpha}_t)} + O\left(\theta d \left(\frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_t}\right)^{3/2} \log T\right). \end{aligned}$$

B Analysis for the accelerated ODE sampler (proof of Theorem 1)

In this section, we present our non-asymptotic analysis for the accelerated ODE sampler. Considering the total variation distance is always upper bounded by 1, we can reasonably assume the following conditions throughout the proof, which are necessary for the claimed result eq. (16) to be non-trivial.

$$T \geq \sqrt{C_1} d^3 \log^3 T, \quad (42a)$$

$$\varepsilon_{\text{score}} \leq \frac{1}{C_1 \sqrt{d} \log^2 T}, \quad (42b)$$

$$\varepsilon_{\text{Jacobi}} \leq \frac{1}{C_1 d \log^2 T}. \quad (42c)$$

B.1 Main steps of the proof

Preparation. To begin with, let us introduce the following functions that help ease presentation:

$$\phi_t^*(x) := x - \frac{1 - \alpha_{t+1}}{2} s_t^*(x), \quad (43a)$$

$$\phi_t(x) := x - \frac{1 - \alpha_{t+1}}{2} s_t(x), \quad (43b)$$

$$\psi_t^*(x) := x + \frac{1 - \alpha_t}{2} s_t^*(x) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \left(s_t^*(x) - \sqrt{\alpha_{t+1}} s_{t+1}^*(\phi_t^*(x)) \right), \quad (43c)$$

$$\psi_t(x) := x + \frac{1 - \alpha_t}{2} s_t(x) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \left(s_t(x) - \sqrt{\alpha_{t+1}} s_{t+1}(\phi_t(x)) \right). \quad (43d)$$

Armed with these functions, one can equivalently rewrite our update rule (15) as follows:

$$Y_{t-1} = \Psi_t(Y_t, \Phi_t(Y_t)) = \frac{1}{\sqrt{\alpha_t}} \psi_t(Y_t). \quad (43e)$$

Additionally, for any point $y_T \in \mathbb{R}^d$, we introduce the corresponding sequence

$$y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \psi_t(y_t), \quad y_t^- = \sqrt{\alpha_{t+1}} \phi_t(y_t), \quad t = T, T-1, \dots \quad (44)$$

Furthermore, it is also useful to single out the following error-related quantities for any point $y_T \in \mathbb{R}^d$ and its associated sequence $\{y_t\}_{t=1}^{T-1}$ and $\{y_t^-\}_{t=2}^T$:

$$\xi_t(y_T) := \frac{\log T}{T} \left\{ d(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t+1}(y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t+1}(y_t^-)) \right\}; \quad (45a)$$

$$S_t(y_T) := \sum_{1 < k \leq t} \xi_k(y_k), \quad \text{for } t \geq 2, \quad \text{and} \quad S_1(y_T) = 0, \quad (45b)$$

where we recall the definitions of $\varepsilon_{\text{Jacobi},t}(\cdot)$ and $\varepsilon_{\text{score},t}$ in (40). To understand these quantities, note that if we start from a point y_T , then $\xi_t(y_T)$ reflects a certain weighted score estimation error in the t -th step, while $S_t(y_T)$ aggregates these weighted score estimation errors from the very beginning to the t -th iteration.

In addition, there are several objects that play crucial roles in the subsequent analysis, which we single out as follows; here and throughout, we suppress their dependency on x to streamline presentation.

$$A_t := \frac{1}{1 - \bar{\alpha}_t} \int p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0; \quad (46a)$$

$$B_t := \frac{1}{1 - \bar{\alpha}_t} \left\| \int p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right\|_2^2; \quad (46b)$$

$$C_t := \frac{1}{(1 - \bar{\alpha}_t)^2} \int p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^4 dx_0; \quad (46c)$$

$$D_t := \frac{1}{(1 - \bar{\alpha}_t)^2} \int p_{X_0 | X_t}(x_0 | x) \left(\langle g_t(x), x - \sqrt{\bar{\alpha}_t} x_0 \rangle \right)^2 dx_0; \quad (46d)$$

$$E_t := \frac{1}{(1 - \bar{\alpha}_t)^2} \int p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 \langle g_t(x), x - \sqrt{\bar{\alpha}_t} x_0 \rangle dx_0. \quad (46e)$$

With the above preparation in place, we can now readily proceed to our proof.

Step 1: bounding the density ratios. To begin with, we make the observation that

$$\begin{aligned} \frac{p_{Y_{t-1}}(y_{t-1})}{p_{X_{t-1}}(y_{t-1})} &= \frac{p_{\sqrt{\alpha_t} Y_{t-1}}(\sqrt{\alpha_t} y_{t-1})}{p_{\sqrt{\alpha_t} X_{t-1}}(\sqrt{\alpha_t} y_{t-1})} \\ &= \frac{p_{\psi_t(Y_t)}(\psi_t(y_t))}{p_{Y_t}(y_t)} \cdot \left(\frac{p_{\sqrt{\alpha_t} X_{t-1}}(\psi_t(y_t))}{p_{X_t}(y_t)} \right)^{-1} \cdot \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)}, \end{aligned} \quad (47)$$

which is an elementary identity that allows one to link the target density ratio $\frac{p_{Y_{t-1}}}{p_{X_{t-1}}}$ at the $(t-1)$ -th step with the density ratio $\frac{p_{Y_t}}{p_{X_t}}$ at the t -th step. This relation reveals the importance of bounding $\frac{p_{\psi_t(Y_t)}(\psi_t(y_t))}{p_{Y_t}(y_t)}$ and $\frac{p_{\sqrt{\alpha_t} X_{t-1}}(\psi_t(y_t))}{p_{X_t}(y_t)}$, towards which we resort to the following lemma.

Lemma 6. For any $x \in \mathbb{R}^d$, suppose that

$$-\frac{\log p_{X_t}(x)}{d \log T} \leq c_6 \quad (48)$$

for some large enough constant $c_6 \geq 2c_R + c_0$, and that $\frac{40c_1 \varepsilon_{\text{score},t}(x) \log^{\frac{3}{2}} T}{T} \leq \sqrt{d}$. Then one has

$$\frac{p_{X_{t+1}/\sqrt{\alpha_{t+1}}}(\phi_t(x))}{p_{X_t}(x)} \geq \exp \left(- \left(\varepsilon_{\text{score},t}(x) \sqrt{d \log T} + d \log T \right) \frac{c_7 \log T}{T} \right) \quad (49)$$

for some universal constant $c_7 > 0$. If, in addition, we have $C_{10} \frac{d \log^2 T + (\varepsilon_{\text{score},t}(x) + \varepsilon_{\text{score},t+1}(\Phi_t(x))) \sqrt{d \log^3 T}}{T} \leq 1$ for some large enough constant $C_{10} > 0$, then it holds that

$$\begin{aligned} \frac{p_{\sqrt{\alpha_t} X_{t-1}}(\psi_t(x))}{p_{X_t}(x)} &= 1 + \frac{(1 - \alpha_t)(d + B_t - A_t)}{2(1 - \bar{\alpha}_t)} \\ &\quad + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)^2} [d(d+2) + (4+2d)(B_t - A_t) - B_t^2 + C_t + 2D_t - 3E_t + A_t B_t] \\ &\quad + O \left(\frac{d^3 \log^6 T}{T^3} + \frac{\sqrt{d \log^3 T}}{T} \left(\varepsilon_{\text{score},t}(x) + \varepsilon_{\text{score},t+1}(\Phi_t(x)) \right) \right). \end{aligned} \quad (50a)$$

Moreover, for any random vector Y , one has

$$\begin{aligned} \frac{p_{\psi_t(Y)}(\psi_t(x))}{p_Y(x)} &= 1 + \frac{(1 - \alpha_t)(d + B_t - A_t)}{2(1 - \bar{\alpha}_t)} \\ &\quad + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)^2} [d(d+2) + (4+2d)(B_t - A_t) - B_t^2 + C_t + 2D_t - 3E_t + A_t B_t] \\ &\quad + O \left(\frac{d^6 \log^6 T}{T^3} + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t}(x) + \frac{d \log T}{T} \left(\varepsilon_{\text{Jacobi},t}(x) + \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x)) \right) \right), \end{aligned} \quad (50b)$$

provided that $C_{11} \frac{d^2 \log^2 T + \varepsilon_{\text{score},t}(x) \sqrt{d \log^3 T} + d(\varepsilon_{\text{Jacobi},t}(x) + \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x))) \log T}{T} \leq 1$ holds for some large enough constant $C_{11} > 0$.

Additionally, if $\frac{d^2 \log^2 T + \sqrt{d \log T} \varepsilon_{\text{score},t}(x) + d \varepsilon_{\text{Jacobi},t}(x) \log T}{T} \lesssim 1$, then we have

$$\frac{p_{\Phi_t(X_t)}(\Phi_t(x))}{p_{X_{t+1}}(\Phi_t(x))} = 1 + O \left(\frac{d^2 \log^4 T}{T^2} + \frac{d^6 \log^6 T}{T^3} + \frac{\sqrt{d \log T} \varepsilon_{\text{score},t}(x) + d \varepsilon_{\text{Jacobi},t}(x) \log T}{T} \right). \quad (51)$$

The proof of this lemma is postponed to Section B.2. Noteworthy, the main terms in (50a) and (50b) coincide, a crucial fact that allows one to focus on the lower-order term later on. Moreover, the relation (51) captures the effect of performing one iteration of the probability flow ODE sampler (as captured by the mapping $\Phi_t(\cdot)$).

Step 2: decomposing the TV distance of interest. We now move on to look at the TV distance of interest. Akin to Li et al. (2023), we first single out the following set:

$$\mathcal{E} := \left\{ y : q_1(y) > \max \{ p_1(y), \exp(-c_6 d \log T) \} \right\}, \quad (52)$$

where $c_6 > 0$ is some large enough numerical constant introduced in Lemma 6. The points in \mathcal{E} satisfy two properties: (i) $q_1(y) > p_1(y)$, and (ii) $q_1(y)$ is not too small, so that y falls within a more normal range (w.r.t. $p_{X_1}(\cdot)$).

Following similar calculations as in Li et al. (2023, Step 2 of Section 5.2) and invoking the properties of the forward process in Lemma 3, we can demonstrate that

$$\text{TV}(q_1, p_1) \leq \mathbb{E}_{Y_1 \sim p_1} \left[\left(\frac{q_1(Y_1)}{p_1(Y_1)} - 1 \right) \mathbb{1} \{ Y_1 \in \mathcal{E} \} \right] + \exp(-c_6 d \log T), \quad (53)$$

and hence it suffices to focus attention on what happens on the set \mathcal{E} . To proceed, for any point y_T , we define

$$\tau(y_T) := \max \left\{ 2 \leq t \leq T + 1 : S_{t-1}(y_T) \leq c_{14} \text{ and } -\log q_k(y_k) \leq c_\tau d \log T, \text{ for all } k < t \right\} \quad (54)$$

for some universal constant $c_\tau > 0$. We shall often abbreviate $\tau(y_T)$ as τ as long as it is clear from the context. Taking $\{y_t\}_{t=1}^{T-1}$ to be the associated sequence of our deterministic sampler initialized at y_T (cf. (44)), we can further define

$$\mathcal{I}_0 := \left\{ y_T : \tau(y_T) = T + 1 \right\}, \quad (55a)$$

$$\mathcal{I}_1 := \left\{ y_T : S_{\tau-1}(y_T) \leq c_{14} \text{ and } -\log q_\tau(y_\tau) > c_\tau d \log T \right\}, \quad (55b)$$

$$\mathcal{I}_2 := \left\{ y_T : c_{14} \leq S_\tau(y_T) \leq 2c_{14} \right\} \cap \mathcal{I}_1^c, \quad (55c)$$

$$\mathcal{I}_3 := \left\{ y_T : S_{\tau-1}(y_T) \leq c_{14}, \xi_\tau(y_T) \geq c_{14}, \frac{q_{\tau-1}(y_{\tau-1})}{p_{\tau-1}(y_{\tau-1})} \leq \frac{8q_\tau(y_\tau)}{p_\tau(y_\tau)} \right\} \cap \mathcal{I}_1^c, \quad (55d)$$

$$\mathcal{I}_4 := \left\{ y_T : S_{\tau-1}(y_T) \leq c_{14}, \xi_\tau(y_T) \geq c_{14}, \frac{q_{\tau-1}(y_{\tau-1})}{p_{\tau-1}(y_{\tau-1})} > \frac{8q_\tau(y_\tau)}{p_\tau(y_\tau)} \right\} \cap \mathcal{I}_1^c. \quad (55e)$$

As an immediate consequence of the above definitions, one has

$$\mathcal{I}_0 \cup \mathcal{I}_1 \cup \mathcal{I}_2 \cup \mathcal{I}_3 \cup \mathcal{I}_4 = \mathbb{R}^d.$$

In the following, we shall look at each of these sets separately, and combine the respective bounds to control the first term on the right-hand side of (53).

Step 3: coping with the set \mathcal{I}_0 . In order to obtain a useful bound when restricting attention to \mathcal{I}_0 (cf. (55a)), we resort to the following key lemma, whose proof is provided in Section B.3.

Lemma 7. *Consider any y_T , along with the deterministic sequences $\{y_{T-1}, \dots, y_1\}$ and $\{y_T^-, \dots, y_2^-\}$. Set $\tau = \tau(y_T)$ (cf. (54)). Then one has*

$$\frac{q_1(y_1)}{p_1(y_1)} = \left\{ 1 + O \left(\frac{d^6 \log^6 T}{T^2} + S_{\tau-1}(y_{\tau-1}) \right) \right\} \frac{q_{\tau-1}(y_{\tau-1})}{p_{\tau-1}(y_{\tau-1})}, \quad (56a)$$

$$\text{and} \quad \frac{q_k(y_k)}{2p_k(y_k)} \leq \frac{q_1(y_1)}{p_1(y_1)} \leq 2 \frac{q_k(y_k)}{p_k(y_k)}, \quad \forall k < \tau. \quad (56b)$$

With this lemma in mind, we are ready to cope with the set \mathcal{I}_0 . Taking $\tau(y_T) = T + 1$ in Lemma 7 yields

$$\begin{aligned}
& \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_1(Y_1)}{p_1(Y_1)} - 1 \right) \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_0\} \right] \\
& \stackrel{(i)}{=} \mathbb{E}_{Y_T \sim p_T} \left[\left(\left\{ 1 + O \left(\frac{d^6 \log^6 T}{T^2} + S_T(y_T) \right) \right\} \frac{q_T(Y_T)}{p_T(Y_T)} - 1 \right) \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_0\} \right] \\
& = \int \left\{ \left(1 + O \left(\frac{d^6 \log^6 T}{T^2} + S_T(y_T) \right) \right) q_T(y_T) - p_T(y_T) \right\} \mathbb{1} \{y_1 \in \mathcal{E}, y_T \in \mathcal{I}_0\} dy_T \\
& \stackrel{(ii)}{\leq} \int |q_T(y_T) - p_T(y_T)| dy_T + O \left(\frac{d^6 \log^6 T}{T^2} \right) \int q_T(y_T) dy_T + O \left(\sqrt{d \log^3 T} \varepsilon_{\text{score}} + (d \log T) \varepsilon_{\text{Jacobi}} \right) \\
& \stackrel{(iii)}{\lesssim} \frac{d^6 \log^6 T}{T^2} + \sqrt{d \log^3 T} \varepsilon_{\text{score}} + (d \log T) \varepsilon_{\text{Jacobi}}. \tag{57}
\end{aligned}$$

Here, (i) invokes Lemma 7, whereas (iii) holds since $\text{TV}(p_T, q_T) \lesssim T^{-100}$ (according to Lemma 6). To see why (ii) is valid, it suffices to make the following observation:

$$\begin{aligned}
& \int S_T(y_T) q_T(y_T) \mathbb{1} \{y_1 \in \mathcal{E}, y_T \in \mathcal{I}_0\} dy_T = \\
& = \frac{\log T}{T} \sum_{t=1}^T \int \left\{ d(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t+1}(y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t+1}(y_t^-)) \right\} \\
& \quad \cdot q_T(y_T) \mathbb{1} \{y_1 \in \mathcal{E}, y_T \in \mathcal{I}_0\} dy_T \\
& \stackrel{(iv)}{\leq} \frac{4 \log T}{T} \sum_{t=1}^T \int \left\{ d(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t+1}(y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t+1}(y_t^-)) \right\} \\
& \quad \cdot \frac{q_t(y_t)}{p_t(y_t)} p_T(y_T) \mathbb{1} \{y_1 \in \mathcal{E}, y_T \in \mathcal{I}_0\} dy_T \\
& \leq \frac{4 \log T}{T} \sum_{t=1}^T \mathbb{E}_{Y_T \sim p_T} \left[\left\{ d(\varepsilon_{\text{Jacobi},t}(Y_t) + \varepsilon_{\text{Jacobi},t+1}(Y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(Y_t) + \varepsilon_{\text{score},t+1}(Y_t^-)) \right\} \right. \\
& \quad \left. \cdot \frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right] \\
& = \frac{4 \log T}{T} \sum_{t=1}^T \mathbb{E}_{Y_t \sim p_t} \left[\left\{ d(\varepsilon_{\text{Jacobi},t}(Y_t) + \varepsilon_{\text{Jacobi},t+1}(Y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(Y_t) + \varepsilon_{\text{score},t+1}(Y_t^-)) \right\} \right. \\
& \quad \left. \cdot \frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right] \\
& = \frac{4 \log T}{T} \sum_{t=1}^T \mathbb{E}_{Y_t \sim q_t} \left[\left\{ d(\varepsilon_{\text{Jacobi},t}(Y_t) + \varepsilon_{\text{Jacobi},t+1}(Y_t^-)) + \sqrt{d \log T} (\varepsilon_{\text{score},t}(Y_t) + \varepsilon_{\text{score},t+1}(Y_t^-)) \right\} \right. \\
& \quad \left. \cdot \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right] \\
& \stackrel{(v)}{\lesssim} \frac{\log T}{T} \sum_{t=1}^T \mathbb{E}_{Y_t \sim q_t} \left[d \varepsilon_{\text{Jacobi},t}(Y_t) + \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) \right] \\
& \stackrel{(vi)}{\lesssim} (d \log T) \varepsilon_{\text{Jacobi}} + \sqrt{d \log^3 T} \varepsilon_{\text{score}}.
\end{aligned}$$

Here, (iv) follows from Lemma 7, while (vi) comes from (41). To understand why (v) is valid, let us denote the probability density of $\Phi(X_t)$ by q_t^- and, by referring to (51), we derive that

$$\mathbb{E}_{Y_t \sim q_t} \left[\varepsilon_{\text{score},t+1}(Y_t^-) \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right]$$

$$\begin{aligned}
&= \mathbb{E}_{Y_t^- \sim q_t^-} \left[\varepsilon_{\text{score},t+1}(Y_t^-) \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right] \\
&\lesssim \mathbb{E}_{Y_t^- \sim q_{t+1}} \left[\varepsilon_{\text{score},t+1}(Y_t^-) \right]. \tag{58a}
\end{aligned}$$

Similarly, we have

$$\mathbb{E}_{Y_t \sim q_t} \left[\varepsilon_{\text{Jacobi},t+1}(Y_t^-) \mathbb{1} \left\{ \sqrt{d \log T} \varepsilon_{\text{score},t}(Y_t) + d \varepsilon_{\text{Jacobi},t}(Y_t) \log T \lesssim T \right\} \right] \lesssim \mathbb{E}_{Y_t^- \sim q_{t+1}} \left[\varepsilon_{\text{Jacobi},t+1}(Y_t^-) \right]. \tag{58b}$$

Step 4: coping with the set \mathcal{I}_1 . In view of Lemma 7, the condition $S_{\tau-1}(y_{\tau-1}) \leq c_{14}$ implies that

$$\frac{q_1(y_1)}{p_1(y_1)} \leq \frac{2q_{\tau-1}(y_{\tau-1})}{p_{\tau-1}(y_{\tau-1})}.$$

This in turn allows one to obtain

$$\begin{aligned}
\mathbb{E}_{Y_T \sim p_T} \left[\frac{q_1(Y_1)}{p_1(Y_1)} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_1\} \right] &\leq 2 \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{\tau-1}(Y_1)}{p_{\tau-1}(Y_1)} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_1\} \right] \\
&= 2 \sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_1)}{p_{t-1}(Y_1)} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_1\} \mathbb{1} \{\tau = t\} \right] \\
&\leq 2 \sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_t \in \mathcal{J}_t\} \right], \tag{59}
\end{aligned}$$

where the last line comes from the definition of \mathcal{I}_1 (cf. (55b)), with \mathcal{J}_t defined as

$$\mathcal{J}_t := \left\{ y_t : -\log q_t(y_t) \geq c_\tau d \log T \right\}. \tag{60}$$

To bound the right-hand side of (59), we make note of the following identities:

$$1 = \mathbb{E}_{Y_t \sim p_t} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \right] = \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \right] = \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_t(Y_t)}{p_t(Y_t)} (\mathbb{1} \{Y_t \in \mathcal{J}_t\} + \mathbb{1} \{Y_t \in \mathcal{J}_t^c\}) \right], \tag{61a}$$

$$1 = \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right] = \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} (\mathbb{1} \{Y_t \in \mathcal{J}_t\} + \mathbb{1} \{Y_t \in \mathcal{J}_t^c\}) \right], \tag{61b}$$

which in turn imply that

$$\begin{aligned}
&\mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \mathbb{1} \{Y_t \in \mathcal{J}_t\} \right] \\
&= 1 - \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \mathbb{1} \{Y_t \in \mathcal{J}_t^c\} \right] \\
&= \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c\} \right] + \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \{Y_t \in \mathcal{J}_t\} \right] \\
&\leq \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c\} \right] + O\left(\exp(-c_6 d \log T)\right).
\end{aligned}$$

Here, the last line follows since

$$\begin{aligned}
\mathbb{E}_{Y_T \sim p_T} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \{Y_t \in \mathcal{J}_t\} \right] &= \mathbb{E}_{Y_t \sim p_t} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \{Y_t \in \mathcal{J}_t\} \right] = \mathbb{E}_{Y_t \sim q_t} \left[\mathbb{1} \{Y_t \in \mathcal{J}_t\} \right] \\
&\lesssim \exp(-c_6 d \log T),
\end{aligned}$$

provided that $c_{20} > 0$ is large enough. As a consequence, the right-hand side of (59) can be bounded by

$$\sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_t \in \mathcal{J}_t\} \right]$$

$$\leq \sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c\} \right] + O\left(T \exp(-c_6 d \log T)\right). \quad (62)$$

Moreover, the first term in the last line (62) can be decomposed as follows

$$\begin{aligned} & \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c\} \right] \\ &= \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c, \xi_t(Y_t) < c_{14}\} \right] \\ & \quad + \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c, \xi_t(Y_t) \geq c_{14}\} \right], \end{aligned} \quad (63)$$

leaving us with two terms to control.

- With regards to the first term on the right-hand side of (63), for Y_t satisfies $\log q_t(Y_t) \leq -c_\tau d \log T$ and $\xi_t(Y_t) < c_{14}$, we can directly apply Lemma 6 to obtain a one-step version of Lemma 7 to control the difference between the density ratio $\frac{q_t(Y_t)}{p_t(Y_t)}$ and $\frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})}$ as follows

$$\begin{aligned} \frac{p_{t-1}(y_{t-1})}{q_{t-1}(y_{t-1})} &= \frac{p_t(y_t)}{q_t(y_t)} \cdot \left\{ 1 + O\left(\frac{d^6 \log^6 T}{T^3}\right) + \right. \\ & \quad \left. O\left(\frac{(\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t}(\Phi_t(y_t)))\sqrt{d \log^3 T}}{T} + \frac{d \log T(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t}(\Phi_t(y_t)))}{T}\right) \right\} \end{aligned}$$

which is an intermediate step in the proof of Lemma 7 (referring (94) in Section B.3 for details). The above relation yields:

$$\begin{aligned} & \sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_t(Y_t)}{p_t(Y_t)} - \frac{q_{t-1}(Y_{t-1})}{p_{t-1}(Y_{t-1})} \right) \mathbb{1} \{Y_t \in \mathcal{J}_t^c, \xi_t(Y_t) < c_{14}\} \right] \\ & \lesssim \frac{d^6 \log^6 T}{T^2} + \sqrt{d \log^3 T} \varepsilon_{\text{score}} + d \log T \varepsilon_{\text{Jacobi}}. \end{aligned}$$

- When it comes to the second term on the right-hand side of (63), we can invoke similar arguments as in Li et al. (2023) (i.e., the arguments therein to bound \mathcal{I}_3), as well as the relation (58) for the score error of Y_t^- , to obtain

$$\sum_{t=2}^T \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_t(Y_t)}{p_t(Y_t)} \mathbb{1} \{Y_t \in \mathcal{J}_t^c, \xi_t(Y_t) \geq c_{14}\} \right] \lesssim \sqrt{d \log^3 T} \varepsilon_{\text{score}} + d \log T \varepsilon_{\text{Jacobi}}.$$

Putting all this together, we arrive at

$$\mathbb{E}_{Y_T \sim p_T} \left[\frac{q_1(Y_1)}{p_1(Y_1)} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_1\} \right] \lesssim \frac{d^6 \log^6 T}{T^2} + \sqrt{d \log^3 T} \varepsilon_{\text{score}} + d \log T \varepsilon_{\text{Jacobi}}. \quad (64)$$

Step 5: coping with the remaining sets. The analyses for $\mathcal{I}_2, \mathcal{I}_3$ and \mathcal{I}_4 are similar to Li et al. (2023). For the sake of brevity, we state the combined result in the lemma below and omit the proof.

Lemma 8. *It holds that*

$$\mathbb{E}_{Y_T \sim p_T} \left[\frac{q_1(Y_1)}{p_1(Y_1)} \mathbb{1} \{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_2 \cup \mathcal{I}_3 \cup \mathcal{I}_4\} \right] \lesssim \frac{d^6 \log^6 T}{T^2} + \sqrt{d \log^3 T} \varepsilon_{\text{score}} + (d \log T) \varepsilon_{\text{Jacobi}}. \quad (65)$$

Step 6: putting all pieces together. Given the preceding results on \mathcal{I}_0 to \mathcal{I}_4 , we can substitute the upper bounds derived in (57), (64) and (65) back into (53), which leads to the following conclusion:

$$\begin{aligned} \text{TV}(p_1, q_1) &\leq \mathbb{E}_{Y_T \sim p_T} \left[\left(\frac{q_1(Y_1)}{p_1(Y_1)} - 1 \right) \mathbb{1}\{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_0\} \right] + \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_1(Y_1)}{p_1(Y_1)} \mathbb{1}\{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_1\} \right] \\ &\quad + \mathbb{E}_{Y_T \sim p_T} \left[\frac{q_1(Y_1)}{p_1(Y_1)} \mathbb{1}\{Y_1 \in \mathcal{E}, Y_T \in \mathcal{I}_2 \cup \mathcal{I}_3 \cup \mathcal{I}_4\} \right] + O\left(\exp(-c_6 d \log T)\right) \\ &\lesssim \frac{d^6 \log^6 T}{T^2} + \sqrt{d \log^3 T \varepsilon_{\text{score}}} + (d \log T) \varepsilon_{\text{Jacobi}}. \end{aligned}$$

B.2 Proof of Lemma 6

B.2.1 Proof of property (49)

This property can be established in a similar way to Li et al. (2023, Lemma 4). Before proceeding, let us introduce the following vector:

$$\begin{aligned} v_t(x) &:= x - \phi_t(x) = x - \phi_t^*(x) + \phi_t^*(x) - \phi_t(x) \\ &= -\frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)} \int_{x_0} (x - \sqrt{\bar{\alpha}_t} x_0) p_{X_0|X_t}(x_0 | x) dx_0 + \frac{1 - \alpha_{t+1}}{2} (s_t(x) - s_t^*(x)), \end{aligned}$$

where we have invoked the definitions of ϕ_t^* and ϕ_t , as well as the property (30). For notational simplicity, we shall abbreviate $v = v_t(x)$ in the following analysis.

Recognizing that

$$X_t \stackrel{d}{=} \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} W \quad \text{with } W \sim \mathcal{N}(0, I_d)$$

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

$$p_{X_0|X_t}(x_0 | x) = \frac{p_{X_0}(x_0)}{p_{X_t}(x)} p_{X_t|X_0}(x | x_0) = \frac{p_{X_0}(x_0)}{p_{X_t}(x)} \cdot \frac{1}{(2\pi(1 - \bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)}\right).$$

Additionally, recalling that

$$\frac{1}{\sqrt{\alpha_{t+1}}} X_{t+1} \stackrel{d}{=} \frac{1}{\sqrt{\alpha_{t+1}}} \left(\sqrt{\bar{\alpha}_{t+1}} X_0 + \sqrt{1 - \bar{\alpha}_{t+1}} W \right) = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{\frac{1}{\alpha_{t+1}} - \bar{\alpha}_t} W,$$

one can demonstrate that

$$\begin{aligned} \frac{p_{X_{t+1}/\sqrt{\alpha_{t+1}}}(\phi_t(x))}{p_{X_t}(x)} &= \frac{1}{p_{X_t}(x)} \int_{x_0} p_{X_0}(x_0) \frac{1}{(2\pi(\alpha_{t+1}^{-1} - \bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|\phi_t(x) - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)}\right) dx_0 \\ &= \frac{1}{p_{X_t}(x)} \int_{x_0} p_{X_0}(x_0) \frac{1}{(2\pi(\alpha_{t+1}^{-1} - \bar{\alpha}_t))^{d/2}} \exp\left(-\frac{\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(1 - \bar{\alpha}_t)}\right) \\ &\quad \cdot \exp\left(-\frac{(1 - \alpha_{t+1}^{-1})\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|v\|_2^2 - 2v^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)}\right) dx_0 \\ &= \left(\frac{1 - \bar{\alpha}_t}{\alpha_{t+1}^{-1} - \bar{\alpha}_t}\right)^{d/2} \cdot \int_{x_0} p_{X_0|X_t}(x_0 | x) \cdot \\ &\quad \exp\left(\frac{(\alpha_{t+1}^{-1} - 1)\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|v\|_2^2 - 2v^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)}\right) dx_0. \end{aligned}$$

To lower bound the above expression, we can focus on controlling the second term within the exponential part of the integral over the set $\mathcal{G}_c^{\text{typical}}$ defined as follows, since the first term is always non-negative:

$$\mathcal{G}_c^{\text{typical}} := \left\{ x_0 : \|x - \sqrt{\bar{\alpha}_t} x_0\|_2 \leq 5c\sqrt{c_6 d(1 - \bar{\alpha}_t) \log T} \right\}. \quad (66)$$

Furthermore, we make the following observations:

- When (48) holds, Lemma 3 implies that

$$\mathbb{P} \left(\|\sqrt{\bar{\alpha}_t} X_0 - x\|_2 > 5c_5 \sqrt{c_6 d(1 - \bar{\alpha}_t) \log T} \mid X_t = x \right) \leq \exp(-c_5^2 c_6 d \log T) \quad (67a)$$

for any quantity $c_5 \geq 2$, provided that $c_6 \geq 2c_R + c_0$.

- The we can makes use of the above relation to bound v as follows:

$$\begin{aligned} \|v\|_2 &\leq \frac{1 - \alpha_{t+1}}{2} \varepsilon_{\text{score},t}(x) + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)} \mathbb{E} [\|\sqrt{\bar{\alpha}_t} X_0 - x\|_2 \mid X_t = x] \\ &\leq \frac{1 - \alpha_{t+1}}{2} \varepsilon_{\text{score},t}(x) + \frac{6(1 - \alpha_{t+1})}{1 - \bar{\alpha}_t} \sqrt{c_6 d(1 - \bar{\alpha}_t) \log T}. \end{aligned} \quad (67b)$$

Therefore, we obtain that: for any $x_0 \in \mathcal{G}$,

$$\begin{aligned} \frac{\|v\|_2^2}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)} &\stackrel{(i)}{\leq} \frac{(1 - \alpha_{t+1})^2}{4(\alpha_{t+1}^{-1} - \bar{\alpha}_t)} \varepsilon_{\text{score},t}(x)^2 + \frac{36(1 - \alpha_{t+1})^2}{(1 - \bar{\alpha}_t)^2 (\alpha_{t+1}^{-1} - \bar{\alpha}_t)} c_6 d \log T \\ &\stackrel{(ii)}{\leq} \frac{2c_1^2 \log^2 T}{T^2} \varepsilon_{\text{score},t}(x)^2 + \frac{2304c_1^2}{T^2} c_6 d \log^3 T; \\ \left| \frac{v^\top (x - \sqrt{\bar{\alpha}_t} x_0)}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \right| &\stackrel{(iii)}{\leq} \frac{\|v\|_2 \|x - \sqrt{\bar{\alpha}_t} x_0\|_2}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \\ &\stackrel{(iv)}{\leq} \frac{20cc_1}{T} \varepsilon_{\text{score},t}(x) \sqrt{c_6 d(1 - \bar{\alpha}_t) \log^3 T} + \frac{240cc_1 c_6 d \log^2 T}{T} \end{aligned}$$

Here, (i) is due to (67); (ii) holds by the choice of learning rates in (33); (iii) follows from the Cauchy-Schwarz inequality; and (iv) comes from the definition of \mathcal{G} and (33). Moreover, (33) also guarantees that

$$\left(\frac{1 - \bar{\alpha}_t}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \right)^{d/2} = \left(1 - \frac{1 - \alpha_{t+1}}{1 - \bar{\alpha}_{t+1}} \right)^{\frac{d}{2}} \geq \exp \left(-\frac{4c_1 d \log T}{T} \right).$$

Combine the above relations to yield

$$\begin{aligned} &\frac{p_{X_{t+1}/\sqrt{\bar{\alpha}_{t+1}}}(\phi_t(x))}{p_{X_t}(x)} \\ &\geq \left(\frac{1 - \bar{\alpha}_t}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \right)^{d/2} \cdot \int_{x_0 \in \mathcal{G}} p_{X_0 | X_t}(x_0 | x) \cdot \exp \left(-\frac{\|v\|_2^2 - 2v^\top (x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_{t+1}^{-1} - \bar{\alpha}_t)} \right) dx_0 \\ &\geq \exp \left(-\left(20c_1 \varepsilon_{\text{score},t}(x) \sqrt{c_6 d \log T} + 240c_1 d \log T \right) \frac{c \log T}{T} \right) \\ &\quad \cdot \exp \left(-\frac{4c_1 d \log T}{T} - \frac{2304c_1^2}{T^2} c_6 d \log^3 T - \frac{2c_1 \varepsilon_{\text{score},t}(x)^2 \log^2 T}{T^2} \right) \\ &\geq \exp \left(-\left(20\varepsilon_{\text{score},t}(x) \sqrt{d \log T} + 300d \log T \right) \frac{cc_1 \log T}{T} \right) \end{aligned}$$

provided that $T \geq \frac{386}{5} c_1 c_6 \log T$ and $\frac{40c_1 \varepsilon_{\text{score},t}(x) \log^{\frac{3}{2}} T}{T} \leq \sqrt{d}$. Taking any fixed $c \geq 2$ we obtain the desired result.

B.2.2 Proof of property (50a)

Before embarking on the proof, we first single out some useful properties about ψ_t^* and s_t^* .

Lemma 9. *Under the same conditions as Lemma 6, it holds that*

$$\|\psi_t(x) - \psi_t^*(x)\|_2 = O\left(\frac{\log T}{T} \left\{ \varepsilon_{\text{score},t}(x) + \varepsilon_{\text{score},t+1}(\Phi_t(x)) \right\}\right), \quad (68a)$$

and

$$\begin{aligned} \frac{\partial \psi_t(x, \Phi_t(x))}{\partial x} &= \frac{\partial \psi_t^*(x, \Phi_t^*(x))}{\partial x} + \tilde{\zeta}_t \\ &= \left(I - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(x) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \frac{\partial (s_t^*(x) - \sqrt{\alpha_{t+1}} s_{t+1}^*(\Phi_t^*(x)))}{\partial x} \right) + \tilde{\zeta}_t \end{aligned} \quad (68b)$$

where the residual term $\tilde{\zeta}_t$ satisfies

$$\|\tilde{\zeta}_t\| = O\left(\frac{\log T}{T} \left\{ \varepsilon_{\text{score},t}(x)/\sqrt{d} + \varepsilon_{\text{Jacobi},t}(x) + \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x)) \right\}\right).$$

Additionally, we introduce the following notations for simplicity:

$$w = \Phi_t^*(x) = \sqrt{\alpha_{t+1}} \left(x - \frac{1 - \alpha_{t+1}}{2} s_t^*(x) \right), \quad (69a)$$

$$z = -(1 - \bar{\alpha}_t) s_t^*(x) = g_t(x). \quad (69b)$$

Then the characterization of s_t^* is summarized in the following lemma.

Lemma 10. *Under the same conditions as Lemma 6, equipped with the notation (69), we can write*

$$\begin{aligned} s_t^*(x) - \sqrt{\alpha_{t+1}} s_{t+1}^*(w) &= -\left(\frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^2} - \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} \|z\|_2^2 \right) z \\ &\quad - \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top z dx_0 \\ &\quad + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0 - z) dx_0 + \zeta_{s_t^*} \end{aligned} \quad (70)$$

where $\|\zeta_{s_t^*}\|_2 = O\left(\frac{(d(1 - \alpha_{t+1}) \log T)^{3/2}}{(1 - \bar{\alpha}_t)^2}\right)$, and

$$\frac{\partial (s_t^*(x) - \sqrt{\alpha_{t+1}} s_{t+1}^*(w))}{\partial x} = -\frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^2} J_t(x) + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} (H_1 + H_4 + H_2 - H_3) + \zeta_{J_t} \quad (71)$$

where $\|\zeta_{J_t}\| = O\left(d^2 \frac{(1 - \alpha_{t+1})^{3/2}}{(1 - \bar{\alpha}_{t+1})^{5/2}} \log^2 T\right)$. Here, we denote $J_t = \frac{\partial z}{\partial x}$, and

$$\begin{aligned} H_1 &:= \frac{\partial}{\partial x} \|z\|_2^2 z, \\ H_2 &:= \frac{\partial}{\partial x} \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0) dx_0, \\ H_3 &:= \frac{\partial}{\partial x} \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 z, \\ H_4 &:= \frac{\partial}{\partial x} \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top z dx_0. \end{aligned}$$

To streamline presentation, we leave the proofs of Lemma 9 and Lemma 10 to Section B.2.5.

Equipped with the relations in Lemma 10, we can derive that

$$\|s_t^*(x) - \sqrt{\alpha_{t+1}}s_{t+1}^*(w)\|_2 \lesssim (1 - \alpha_{t+1}) \left(\frac{d \log T}{1 - \bar{\alpha}_t} \right)^{3/2}, \quad (73a)$$

$$\left\| \frac{\partial(s_t^*(x) - \sqrt{\alpha_{t+1}}s_{t+1}^*(w))}{\partial x} \right\| \lesssim \frac{d^2(1 - \alpha_{t+1}) \log^2 T}{(1 - \bar{\alpha}_t)^2}. \quad (73b)$$

The proof of (73) is also deferred to Section B.2.5.

Next, let us introduce the following vectors:

$$\begin{aligned} u_t(x) &:= x - \psi_t(x), \\ u_t^*(x) &:= x - \psi_t^*(x). \end{aligned}$$

For notational simplicity, we shall abbreviate $u = u_t(x)$ and $u^* = u_t^*(x)$ in the following analysis. Akin to the calculations in Appendix B.2.1, we can obtain

$$\begin{aligned} p_{\sqrt{\alpha_t}X_{t-1}}(\psi_t(x)) &= p_{X_t}(x) \left(\frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \right)^{d/2} \\ &\cdot \int_{x_0} p_{X_0|X_t}(x_0|x) \exp \left(- \frac{(1 - \alpha_t) \|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) dx_0. \end{aligned} \quad (74)$$

Similarly, by focusing mainly on the following set given x :

$$\mathcal{G} := \{x_0 : \|x - \sqrt{\alpha_t}x_0\|_2 \lesssim \sqrt{d(1 - \bar{\alpha}_t) \log T}\}, \quad (75)$$

we can derive

$$\begin{aligned} &\int_{x_0} p_{X_0|X_t}(x_0|x) \exp \left(- \frac{(1 - \alpha_t) \|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) dx_0 = O(\exp(-c_8 d \log T)) \\ &+ \int_{x_0 \in \mathcal{E}} p_{X_0|X_t}(x_0|x) \exp \left(- \frac{(1 - \alpha_t) \|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right) dx_0 =: \text{RHS} \end{aligned} \quad (76)$$

for some numerical constant $c_8 > 0$. To further control the right-hand side above, recall that the learning rates are selected such that $\frac{1 - \alpha_t}{1 - \bar{\alpha}_{t-1}} \leq \frac{4c_1 \log T}{T}$ for $1 < t \leq T$ (see (33b)). In view of the Taylor expansion $e^{-x} = 1 - x + \frac{1}{2}x^2 + O(x^3)$ for $x \leq 1/2$, we can derive

$$\begin{aligned} \text{RHS} &= O(\exp(-c_8 d \log T)) + O \left(\frac{d^3 \log^6 T}{T^3} + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t}(x) + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x)) \right) \\ &+ \int_{x_0 \in \mathcal{E}} p_{X_0|X_t}(x_0|x) \left\{ 1 - \frac{(1 - \alpha_t) \|x - \sqrt{\alpha_t}x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} - \frac{\frac{(1 - \alpha_t)^2}{4(1 - \bar{\alpha}_t)^2} \|z\|_2^2 - 2u^{*\top}(x - \sqrt{\alpha_t}x_0)}{2(\alpha_t - \bar{\alpha}_t)} \right. \\ &\left. + \frac{(1 - \alpha_t)^2}{8(\alpha_t - \bar{\alpha}_t)^2(1 - \bar{\alpha}_t)^2} \left(\|x - \sqrt{\alpha_t}x_0\|_2^2 - z^\top(x - \sqrt{\alpha_t}x_0) \right)^2 \right\} dx_0. \end{aligned} \quad (77)$$

Here, we have made use of the following facts:

$$\begin{aligned} \left\| u - \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} z \right\|_2 &= \left\| x - \psi_t(x) + \frac{1 - \alpha_t}{2} s_t^*(x) \right\|_2 \\ &\stackrel{(i)}{\leq} \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \|s_t^*(x) - \sqrt{\alpha_{t+1}}s_{t+1}^*(\Phi_t^*(x))\|_2 + O \left(\frac{\log T}{T} \left\{ \varepsilon_{\text{score},t}(x) + \varepsilon_{\text{score},t+1}(\Phi_t(x)) \right\} \right) \\ &\stackrel{(ii)}{\lesssim} (1 - \alpha_t)^2 \left(\frac{d \log T}{1 - \bar{\alpha}_t} \right)^{3/2} + \frac{\log T}{T} \varepsilon_{\text{score},t}(x) + \frac{\log T}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x)), \end{aligned} \quad (78)$$

where (i) follows from (68a) in Lemma 9 and (ii) follows from (73).

Moreover, for any $x_0 \in \mathcal{E}$, using the definition of \mathcal{E} (cf. (75)) and combining it with the properties (33) of the learning rates, we reach

$$\frac{(1 - \alpha_t) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} = O\left(\frac{d \log^2 T}{T}\right).$$

As a result, we can derive

$$\begin{aligned} & \frac{\|u\|_2^2 - 2u^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} \\ & \stackrel{(i)}{=} \frac{\frac{(1-\alpha_t)^2}{4(1-\bar{\alpha}_t)^2} \|z\|_2^2 - 2u^{\star\top}(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)} + O\left(\frac{d^2 \log^5 T}{T^3} + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t}(x) + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x))\right) \\ & = \frac{z^\top(x - \sqrt{\bar{\alpha}_t} x_0)}{2(\alpha_t - \bar{\alpha}_t)(1 - \bar{\alpha}_t)} + O\left(\frac{d^2 \log^4 T}{T^2} + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t}(x) + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x))\right) \\ & = O\left(\frac{d \log^2 T}{T} + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t}(x) + \frac{\sqrt{d \log^3 T}}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x))\right), \end{aligned}$$

where (i) follows from (78) and Lemma 10. Taking the above results together and using the following basic properties regarding quantities A_t, \dots, E_t (defined in (46))

$$\begin{aligned} \int p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0 &= (1 - \bar{\alpha}_t) A_t, \\ \int p_{X_0 | X_t}(x_0 | x) \|z\|_2^2 dx_0 &= (1 - \bar{\alpha}_t) B_t, \\ \int p_{X_0 | X_t}(x_0 | x) u^{\star\top}(x - \sqrt{\bar{\alpha}_t} x_0) dx_0 &= \frac{1 - \alpha_t}{2} B_t + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)} [B_t - B_t^2 + D_t - E_t + A_t B_t], \\ \int p_{X_0 | X_t}(x_0 | x) \left(\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 - z^\top(x - \sqrt{\bar{\alpha}_t} x_0)\right)^2 dx_0 &= (1 - \bar{\alpha}_t)^2 [C_t + D_t - 2E_t], \end{aligned}$$

we arrive at

$$(77) = 1 - \frac{(1 - \alpha_t)(A_t - B_t)}{2(\alpha_t - \bar{\alpha}_t)} + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)^2} [-B_t^2 + C_t + 2D_t - 3E_t + A_t B_t] + O\left(\frac{d^3 \log^6 T}{T^3}\right).$$

Once again, we note that integrating over the set \mathcal{E} and over all possible x_0 only incurs a difference at most as large as $O(\exp(-c_8 d \log T))$.

Putting the preceding results together establishes the claimed property (50a).

B.2.3 Proof of property (50b)

Consider any random vector Y , and recall the basic transformation

$$p_{\psi_t(Y)}(\psi_t(x)) = \det\left(\frac{\partial \psi_t(x)}{\partial x}\right)^{-1} p_Y(x),$$

where $\frac{\partial \psi_t(x)}{\partial x}$ denotes the Jacobian matrix. It then comes down to controlling the quantity $\det\left(\frac{\partial \psi_t(x)}{\partial x}\right)^{-1}$.

Towards this end, note that the determinant of a matrix obeys

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

with the proviso that $d \|A\| \lesssim 1$. This relation taken together with $\frac{\partial \psi_t(x)}{\partial x} = I - \frac{\partial u^*}{\partial x} + \frac{\partial(u^* - u)}{\partial x}$ leads to

$$p_{\psi_t(Y)}(\psi_t(x)) = \det\left(\frac{\partial \psi_t(x)}{\partial x}\right)^{-1} p_Y(x)$$

$$= \left\{ 1 + \text{Tr}\left(\frac{\partial u^*}{\partial x}\right) + \frac{1}{2} \left[\text{Tr}\left(\frac{\partial u^*}{\partial x}\right)^2 + \left\| \frac{\partial u^*}{\partial x} \right\|_{\text{F}}^2 \right] + O\left(d^3 \left\| \frac{\partial u^*}{\partial x} \right\| + d \left\| \frac{\partial(u - u^*)}{\partial x} \right\| \right) \right\} p_Y(x). \quad (79)$$

To further control the right-hand side of the above display, let us first make note of several identities introduced in Lemma 10:

$$J_t = I + \frac{1}{1 - \bar{\alpha}_t} \left\{ \left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right) \left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right)^\top - \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top dx_0 \right\}; \quad (80a)$$

$$H_1 = \|z\|_2^2 J_t + 2zz^\top J_t; \quad (80b)$$

$$H_2 = \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0 I + 2 \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top dx_0 + \frac{1}{1 - \bar{\alpha}_t} \left(\left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0) \right) \left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right)^\top - \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top dx_0 \right); \quad (80c)$$

$$H_3 = \int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 dx_0 J_t + 2zz^\top + \frac{1}{1 - \bar{\alpha}_t} \left(\left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 \right) zz^\top - z \left(\int_{x_0} p_{X_0 | X_t}(x_0 | x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0) dx_0 \right)^\top \right); \quad (80d)$$

$$H_4 = \|z\|_2^2 I + zz^\top + \int_{x_0} p_{X_0 | X_t}(x_0 | x) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top J_t dx_0 + \frac{1}{1 - \bar{\alpha}_t} \int_{x_0} p_{X_0 | X_t}(x_0 | x) (z^\top (x - \sqrt{\bar{\alpha}_t} x_0)) (x - \sqrt{\bar{\alpha}_t} x_0) z^\top dx_0 - \frac{1}{1 - \bar{\alpha}_t} \int_{x_0} p_{X_0 | X_t}(x_0 | x) (z^\top (x - \sqrt{\bar{\alpha}_t} x_0)) (x - \sqrt{\bar{\alpha}_t} x_0) (x - \sqrt{\bar{\alpha}_t} x_0)^\top dx_0. \quad (80e)$$

The above identities can be directly verified through elementary calculation involving Gaussian integration and derivatives, which are omitted here for the sake of brevity.

Recall the definition of u^* that

$$\begin{aligned} \frac{\partial u^*}{\partial x} &= -\frac{1 - \alpha_t}{2} J_{s_t^*}(x) - \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \frac{\partial(s_t^*(x) - \sqrt{\alpha_{t+1}} s_{t+1}^*(w))}{\partial x} \\ &\stackrel{(i)}{=} \frac{1 - \alpha_t}{2(1 - \bar{\alpha}_t)} J_t(x) - \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} \\ &\quad \left(-\frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^2} J_t(x) + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} (H_1 + H_4 + H_2 - H_3) + \zeta_{J_t} \right), \end{aligned}$$

where $\|\zeta_{J_t}\| \lesssim d^2 \frac{(1 - \alpha_{t+1})^{3/2}}{(1 - \bar{\alpha}_{t+1})^{5/2}} \log^2 T$. Here, (i) follows from Lemma 10. Then, invoking (80) and the definitions of A_t to E_t gives

$$\left\| \frac{\partial u^*}{\partial x} \right\| \lesssim \frac{d(1 - \alpha_t) \log T}{1 - \bar{\alpha}_t}, \quad (81a)$$

$$\begin{aligned} \text{Tr}\left(\frac{\partial u^*}{\partial x}\right) &= \frac{(1 - \alpha_t)(d + B_t - A_t)}{2(1 - \bar{\alpha}_t)} \\ &\quad + \frac{(1 - \alpha_t)^2}{8(1 - \bar{\alpha}_t)^2} (d - 2A_t - A_t^2 + 3A_t B_t + 2B_t - 3B_t^2 + C_t + 4D_t - 3E_t - F_t), \end{aligned} \quad (81b)$$

$$\begin{aligned}\left\|\frac{\partial u^*}{\partial x}\right\|_{\mathbb{F}}^2 &= \frac{(1-\alpha_t)^2}{4(1-\bar{\alpha}_t)^2}\left\|\frac{\partial z}{\partial x}\right\|_{\mathbb{F}}^2 + O\left(d^5\left(\frac{1-\alpha_t}{\alpha_t-\bar{\alpha}_t}\right)^3\log^3 T\right) \\ &= \frac{(1-\alpha_t)^2}{4(1-\bar{\alpha}_t)^2}(d+2(B_t-A_t)+B_t^2+F_t-2D_t) + O\left(d^5\left(\frac{1-\alpha_t}{\alpha_t-\bar{\alpha}_t}\right)^3\log^3 T\right),\end{aligned}\quad (81c)$$

as long as $d^2\left(\frac{1-\alpha_t}{\alpha_t-\bar{\alpha}_t}\right)\log T \lesssim 1$. Here,

$$F_t(x) := \left\|\frac{1}{1-\bar{\alpha}_t}\int_{x_0} p_{X_0|X_t}(x_0|x)(x-\sqrt{\bar{\alpha}_t}x_0)(x-\sqrt{\bar{\alpha}_t}x_0)^\top dx_0\right\|_{\mathbb{F}}^2. \quad (81d)$$

Further note that $\frac{\partial(u^*-u)}{\partial x} = \tilde{\zeta}_t$, where $\tilde{\zeta}_t$ is the residual term defined in (68b) from Lemma 9, and satisfies

$$\|\tilde{\zeta}_t\| = O\left(\frac{\log T}{T}\left\{\varepsilon_{\text{score},t}(x)/\sqrt{d} + \varepsilon_{\text{Jacobi},t}(x) + \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x))\right\}\right). \quad (81e)$$

Substituting these results into inequality (79) leads to

$$\begin{aligned}p_{\psi_t(Y)}(\psi_t(x)) &= p_Y(x)\left\{1 + \frac{(1-\alpha_t)(d+B_t-A_t)}{2(1-\bar{\alpha}_t)} + O\left(d^6\left(\frac{1-\alpha_t}{\alpha_t-\bar{\alpha}_t}\right)^3\log^3 T\right)\right. \\ &\quad \left.+ O\left(\frac{\sqrt{d\log^3 T}}{T}\varepsilon_{\text{score},t}(x) + \frac{d\log T}{T}\left(\varepsilon_{\text{Jacobi},t}(x) + \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x))\right)\right)\right\}\end{aligned}\quad (82)$$

$$\left.+ \frac{(1-\alpha_t)^2}{8(1-\bar{\alpha}_t)^2}\left[d(d+2) + (4+2d)(B_t-A_t) - B_t^2 + C_t + 2D_t - 3E_t + A_tB_t\right]\right\}. \quad (83)$$

B.2.4 Proof of property (51)

Following similar arguments as in Li et al. (2023, relation (58a)), we can obtain

$$\begin{aligned}&\frac{p_{X_{t+1}/\sqrt{\alpha_{t+1}}}(\phi_t(x))}{p_{X_t}(x)} \\ &= \left(\frac{\alpha_{t+1}-\bar{\alpha}_{t+1}}{1-\bar{\alpha}_{t+1}}\right)^{d/2} \cdot \int_{x_0} p_{X_0|X_t}(x_0|x) \cdot \\ &\quad \exp\left(-\frac{(1-\alpha_{t+1}^{-1})\|x-\sqrt{\bar{\alpha}_t}x_0\|_2^2}{2(\alpha_{t+1}^{-1}-\bar{\alpha}_t)(1-\bar{\alpha}_t)} - \frac{\|v\|_2^2 - 2v^\top(x-\sqrt{\bar{\alpha}_t}x_0)}{2(\alpha_{t+1}^{-1}-\bar{\alpha}_t)}\right) dx_0 \\ &= \left(1 - \frac{d(1-\alpha_{t+1})}{2(1-\bar{\alpha}_{t+1})} + O\left(\frac{d^2(1-\alpha_{t+1})^2}{(1-\bar{\alpha}_{t+1})^2}\right)\right) \cdot \int_{x_0} p_{X_0|X_t}(x_0|x) \cdot \\ &\quad \exp\left(\frac{(1-\alpha_{t+1})\|x-\sqrt{\bar{\alpha}_t}x_0\|_2^2}{2(1-\bar{\alpha}_{t+1})(1-\bar{\alpha}_t)} - \frac{\|v\|_2^2 - 2v^\top(x-\sqrt{\bar{\alpha}_t}x_0)}{2(\alpha_{t+1}^{-1}-\bar{\alpha}_t)}\right) dx_0 \\ &= 1 - \frac{d(1-\alpha_{t+1})}{2(1-\bar{\alpha}_{t+1})} + O\left(c_6^2 d^2 \left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_{t+1}}\right)^2 \log^2 T + \varepsilon_{\text{score},t}(x)\sqrt{c_6 d \log T} \left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_{t+1}}\right)\right) + \\ &\quad \frac{(1-\alpha_{t+1})\left(\int_{x_0} p_{X_0|X_t}(x_0|x)\|x-\sqrt{\bar{\alpha}_t}x_0\|_2^2 dx_0 + \alpha_{t+1}\left\|\int_{x_0} p_{X_0|X_t}(x_0|x)(x-\sqrt{\bar{\alpha}_t}x_0) dx_0\right\|_2^2\right)}{2(1-\bar{\alpha}_{t+1})(1-\bar{\alpha}_t)}.\end{aligned}\quad (84)$$

Similarly, using the arguments as in Li et al. (2023, relation (58b)), we can deduce that

$$\begin{aligned}\frac{p_{\phi_t(X_t)}(\phi_t(x))}{p_{X_t}(x)} &= 1 - \frac{d(1-\alpha_{t+1})}{2(1-\bar{\alpha}_{t+1})} + \\ &\quad \frac{(1-\alpha_{t+1})\left(\int_{x_0} p_{X_0|X_t}(x_0|x)\|x-\sqrt{\bar{\alpha}_t}x_0\|_2^2 dx_0 + \alpha_{t+1}\left\|\int_{x_0} p_{X_0|X_t}(x_0|x)(x-\sqrt{\bar{\alpha}_t}x_0) dx_0\right\|_2^2\right)}{2(1-\bar{\alpha}_{t+1})(1-\bar{\alpha}_t)} + \\ &\quad O\left(c_6^2 d^2 \left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_{t+1}}\right)^2 \log^2 T + c_6^3 d^6 \log^3 T \left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_{t+1}}\right)^3 + (1-\alpha_{t+1})d\varepsilon_{\text{Jacobi},t}(x)\right).\end{aligned}\quad (85)$$

Consequently, it is readily seen that

$$\begin{aligned} \frac{p_{\Phi_t(X_t)}(\Phi_t(x))}{p_{X_{t+1}}(\Phi_t(x))} &= \frac{p_{\phi_t(X_t)}(\phi_t(x))}{p_{X_{t+1}/\sqrt{\alpha_{t+1}}}(\phi_t(x))} = \frac{p_{\phi_t(X_t)}(\phi_t(x))}{p_{X_t}(x)} \cdot \left(\frac{p_{X_{t+1}/\sqrt{\alpha_{t+1}}}(\phi_t(x))}{p_{X_t}(x)} \right)^{-1} \\ &= 1 + O\left(\frac{d^2 \log^4 T}{T^2} + \frac{d^6 \log^6 T}{T^3} + \frac{\sqrt{d \log T} \varepsilon_{\text{score},t}(x) + d \varepsilon_{\text{Jacobi},t}(x) \log T}{T} \right), \end{aligned}$$

thus completing the proof of Lemma 6.

B.2.5 Proof of additional lemmas

To establish Lemma 9 and Lemma 10, making use of Lemma 3, we first summarize the following norm properties of the score function s_t^* and the Jacobian matrix $J_{s_t^*}$ for x satisfying $\log p_{X_t}(x) \geq -c_6 d \log T$:

$$\|s_t^*(x_t)\|_2 \leq \frac{1}{(1-\bar{\alpha}_t)} \mathbb{E} [\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2 \mid X_t = x_t] \lesssim \sqrt{\frac{d \log T}{1-\bar{\alpha}_t}}, \quad (86a)$$

$$\|J_{s_t^*}(x)\| \lesssim \frac{1}{(1-\bar{\alpha}_t)^2} \mathbb{E} [\|x - \sqrt{\bar{\alpha}_t} X_0\|_2^2 \mid X_t = x] \asymp d \log T, \quad (86b)$$

$$\|\nabla_x u^\top J_{s_t^*}(x) u\|_2 \lesssim d^{3/2} \log^{3/2} T, \quad \text{for } u \in \mathbb{S}^{d-1}. \quad (86c)$$

The detailed calculation for the second property is presented in Li et al. (2023, Lemma 8). The third property follows a rationale akin to that for $J_{s_t^*} = \frac{\partial s_t^*(x)}{\partial x}$, and is therefore omitted here for the sake of brevity.

Proof of Lemma 9. To start with, it follows from the definitions of ψ_t and ψ_t^* (cf. (43)) that

$$\begin{aligned} \psi_t(x) - \psi_t^*(x) &= \left(\frac{1-\alpha_t}{2} + \frac{(1-\alpha_t)^2}{4(1-\alpha_{t+1})} \right) (s_t(x) - s_t^*(x)) \\ &\quad - \frac{(1-\alpha_t)^2 \sqrt{\alpha_{t+1}}}{4(1-\alpha_{t+1})} (s_{t+1}(\Phi_t(x)) - s_{t+1}^*(\Phi_t(x)) + s_{t+1}^*(\Phi_t(x)) - s_{t+1}^*(\Phi_t^*(x))). \end{aligned}$$

Armed with this relation, to derive (68a), we only need to control the following term:

$$\begin{aligned} (1-\alpha_{t+1}) \left\| s_{t+1}^*(\Phi_t(x)) - s_{t+1}^*(\Phi_t^*(x)) \right\|_2 &\stackrel{(i)}{\lesssim} \frac{\log T}{T} \left\| s_{t+1}^*(\Phi_t(x)) - s_{t+1}^*(\Phi_t^*(x)) \right\|_2 \\ &\stackrel{(ii)}{\lesssim} \frac{\log T}{T} d(\log T) \|\Phi_t(x) - \Phi_t^*(x)\|_2 \\ &\stackrel{(iii)}{\lesssim} \frac{d \log^3 T}{T^2} \varepsilon_{\text{score},t}(x). \end{aligned}$$

Here, (i) follows directly from the choice of learning rate in (33); (ii) holds by observing that both $\Phi_t(x)$ and $\Phi_t^*(x)$ remain within the typical set with $\log p_{X_{t+1}}(\Phi_t(x)), \log p_{X_{t+1}}(\Phi_t^*(x)) \geq -c_6 d \log T$ due to (49), and then invoking (86b); (iii) is due to the definition of $\varepsilon_{\text{score},t}(x)$ (cf. (40)). Combining the above bound with (40) and (33), we arrive at

$$\begin{aligned} \|\psi_t(x) - \psi_t^*(x)\| &\lesssim \frac{\log T}{T} \varepsilon_{\text{score},t}(x) + \frac{\log T}{T} \varepsilon_{\text{score},t+1}(\Phi(x)) + \frac{d \log^3 T}{T^2} \varepsilon_{\text{score},t}(x) \\ &\lesssim \frac{\log T}{T} (\varepsilon_{\text{score},t}(x) + \varepsilon_{\text{score},t+1}(\Phi(x))). \end{aligned}$$

For (68b), by direct calculation, we have

$$\frac{\partial \psi_t(x, \Phi_t(x))}{\partial x} - \frac{\partial \psi_t^*(x, \Phi_t^*(x))}{\partial x} = \left(\frac{1-\alpha_t}{2} + \frac{(1-\alpha_t)^2}{4(1-\alpha_{t+1})} \right) (J_{s_t}(x) - J_{s_t^*}(x))$$

$$+ \frac{\sqrt{\alpha_{t+1}}(1-\alpha_t)^2}{4(1-\alpha_{t+1})} \left(J_{s_{t+1}}(\Phi_t(x)) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t}(x) \right) - J_{s_{t+1}^*}(\Phi_t^*(x)) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t^*}(x) \right) \right).$$

The term in the first line can be directly bounded by the definition of $\varepsilon_{\text{Jacobi},t}(x)$ and (33) as follows

$$\left\| \left(\frac{1-\alpha_t}{2} + \frac{(1-\alpha_t)^2}{4(1-\alpha_{t+1})} \right) (J_{s_t}(x) - J_{s_t^*}(x)) \right\| \lesssim \frac{d \log T}{T} \varepsilon_{\text{Jacobi},t}(x). \quad (87)$$

Turning to the second line, we have

$$\begin{aligned} & J_{s_{t+1}}(\Phi_t(x)) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t}(x) \right) - J_{s_{t+1}^*}(\Phi_t^*(x)) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t^*}(x) \right) \\ &= \left(J_{s_{t+1}}(\Phi_t(x)) - J_{s_{t+1}^*}(\Phi_t^*(x)) \right) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t^*}(x) \right) + \frac{1-\alpha_{t+1}}{2} J_{s_{t+1}}(\Phi_t(x)) (J_{s_t^*}(x) - J_{s_t}(x)). \end{aligned} \quad (88)$$

To proceed, we further observe that

$$\|J_{s_{t+1}^*}(\Phi_t(x)) - J_{s_{t+1}^*}(\Phi_t^*(x))\| \lesssim (d \log T)^{\frac{3}{2}} \|\Phi_t(x) - \Phi_t^*(x)\|_2 \lesssim \frac{d^{\frac{3}{2}} \log^{\frac{5}{2}} T}{T} \varepsilon_{\text{score},t}(x),$$

which is obtained by invoking (86c), (33) and (40). This bound together with (86b) allows us to control the first term in (88) as follows

$$\begin{aligned} & \left\| \left(J_{s_{t+1}}(\Phi_t(x)) - J_{s_{t+1}^*}(\Phi_t^*(x)) \right) \left(I - \frac{1-\alpha_{t+1}}{2} J_{s_t^*}(x) \right) \right\| \\ & \lesssim \left\| \left(J_{s_{t+1}}(\Phi_t(x)) - J_{s_{t+1}^*}(\Phi_t(x)) \right) \right\| + \left\| \left(J_{s_{t+1}^*}(\Phi_t(x)) - J_{s_{t+1}^*}(\Phi_t^*(x)) \right) \right\| \\ & \lesssim \frac{d \log T}{T} \varepsilon_{\text{Jacobi},t+1}(\Phi_t(x)) + \frac{d^{\frac{3}{2}} \log^{\frac{5}{2}} T}{T} \varepsilon_{\text{score},t}(x). \end{aligned} \quad (89)$$

The second term in (88) can be controlled by (86b) and (40) as follows

$$\left\| \frac{1-\alpha_{t+1}}{2} J_{s_{t+1}}(\Phi_t(x)) (J_{s_t^*}(x) - J_{s_t}(x)) \right\| \lesssim \frac{\log T}{T} \left(\varepsilon_{\text{Jacobi},t+1}(\Phi_t(x)) + d \log T \right) \varepsilon_{\text{Jacobi},t}(x). \quad (90)$$

Substituting (89) and (90) into (88), together with (87), we obtain

$$\begin{aligned} \left\| \frac{\partial \psi_t(x, \Phi_t(x))}{\partial x} - \frac{\partial \psi_t^*(x, \Phi_t^*(x))}{\partial x} \right\| & \lesssim \frac{d \log T}{T} \varepsilon_{\text{Jacobi},t}(x) + \frac{d \log T}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x)) + \frac{d^{\frac{3}{2}} \log^{\frac{7}{2}} T}{T^2} \varepsilon_{\text{score},t}(x) \\ & \quad + \frac{\log^2 T}{T^2} \left(\varepsilon_{\text{Jacobi},t+1}(\Phi_t(x)) + d \log T \right) \varepsilon_{\text{Jacobi},t}(x) \\ & \lesssim \frac{d \log T}{T} \varepsilon_{\text{Jacobi},t}(x) + \frac{d \log T}{T} \varepsilon_{\text{score},t+1}(\Phi_t(x)) + \frac{\log T}{d^{\frac{1}{2}} T} \varepsilon_{\text{score},t}(x) \end{aligned}$$

where the last inequality invokes conditions for d and T in (42). \square

Proof of Lemma 10. To begin with, applying (86a) to w leads to

$$\left\| \frac{1}{\sqrt{\alpha_{t+1}}} w - x \right\|_2 = (1-\alpha_{t+1}) \|s_t^*(x)\|_2 \lesssim (1-\alpha_{t+1}) \sqrt{\frac{d \log T}{1-\alpha_t}}. \quad (91)$$

Then denoting $\hat{w} := w/\sqrt{\alpha_{t+1}}$, we can apply Lemma 4 to (\hat{w}, x) to obtain

$$K_1 := \int_{x_0} p_{X_0 | X_{t+1}}(x_0 | w) (w/\sqrt{\alpha_{t+1}} - \sqrt{\alpha_t} x_0) dx_0$$

$$\begin{aligned}
&= \widehat{w} - x + \int_{x_0} p_{X_0|X_t}(x_0|x) \left\{ 1 + \frac{(1-\alpha_{t+1})\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2\alpha_{t+1}(1-\bar{\alpha}_t)^2} - \frac{(x - \sqrt{\bar{\alpha}_t}x_0)^\top(\widehat{w} - x)}{1-\bar{\alpha}_t} \right. \\
&\quad - \int_{x_0} \left(\frac{(1-\alpha_{t+1})\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2\alpha_{t+1}(1-\bar{\alpha}_t)^2} - \frac{(x - \sqrt{\bar{\alpha}_t}x_0)^\top(\widehat{w} - x)}{1-\bar{\alpha}_t} \right) p_{X_0|X_t}(x_0|x) dx_0 \\
&\quad \left. + O\left(d\left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_t}\right)^{3/2} \log T\right) \right\} (x - \sqrt{\bar{\alpha}_t}x_0) dx_0.
\end{aligned}$$

Plugging $\widehat{w} - x = -\frac{1-\alpha_{t+1}}{2}s_t^*(x) = \frac{1-\alpha_{t+1}}{2(1-\bar{\alpha}_t)}z$ into the above equation and combining with the expression of $s_t^*(x)$ in (30), we can obtain

$$\begin{aligned}
K_1 &= \left(1 + \frac{1-\alpha_{t+1}}{2(1-\bar{\alpha}_t)} + \frac{1-\alpha_{t+1}}{2(1-\bar{\alpha}_t)^2}\|z\|_2^2\right)z + \frac{1-\alpha_{t+1}}{2(1-\bar{\alpha}_t)^2} \int_{x_0} p_{X_0|X_t}(x_0|x) \left\{ \|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2 \right. \\
&\quad \left. - (x - \sqrt{\bar{\alpha}_t}x_0)^\top z - \int_{x_0} \|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2 p_{X_0|X_t}(x_0|x) dx_0 \right\} (x - \sqrt{\bar{\alpha}_t}x_0) dx_0 + \zeta_{K_1},
\end{aligned}$$

where the residual term obeys

$$\|\zeta_{K_1}\|_2 \lesssim \frac{(d(1-\alpha_{t+1})\log T)^{3/2}}{1-\bar{\alpha}_t}.$$

Then we can immediately establish the first claim (70) by recognizing that

$$s_t^*(x) - \sqrt{\alpha_{t+1}}s_{t+1}^*(w) = \frac{1}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} K_1 - \frac{1}{1-\bar{\alpha}_t} z.$$

Similarly, one sees that

$$\begin{aligned}
&K_2 \\
&:= \int_{x_0} p_{X_0|X_{t+1}}(x_0|w) (w/\sqrt{\alpha_{t+1}} - \sqrt{\bar{\alpha}_t}x_0) (w/\sqrt{\alpha_{t+1}} - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0 \\
&= \int_{x_0} p_{X_0|X_{t+1}/\sqrt{\alpha_{t+1}}}(x_0|\widehat{w}) \left[(\widehat{w} - x) (\widehat{w} - \sqrt{\bar{\alpha}_t}x_0)^\top + (\widehat{w} - \sqrt{\bar{\alpha}_t}x_0) (\widehat{w} - x)^\top - (\widehat{w} - x) (\widehat{w} - x)^\top \right] dx_0 \\
&\quad + \int_{x_0} p_{X_0|X_t}(x_0|x) \left\{ 1 + \frac{(1-\alpha_{t+1})\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2\alpha_{t+1}(1-\bar{\alpha}_t)^2} - \frac{(x - \sqrt{\bar{\alpha}_t}x_0)^\top(\widehat{w} - x)}{1-\bar{\alpha}_t} + O\left(d\left(\frac{1-\alpha_{t+1}}{1-\bar{\alpha}_t}\right)^{3/2} \log T\right) \right. \\
&\quad \left. - \int_{x_0} \left(\frac{(1-\alpha_{t+1})\|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2}{2\alpha_{t+1}(1-\bar{\alpha}_t)^2} - \frac{(x - \sqrt{\bar{\alpha}_t}x_0)^\top(\widehat{w} - x)}{1-\bar{\alpha}_t} \right) p_{X_0|X_t}(x_0|x) dx_0 \right\} \\
&\quad \cdot (x - \sqrt{\bar{\alpha}_t}x_0) (x - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0 \\
&= \int_{x_0} p_{X_0|X_t}(x_0|x) (x - \sqrt{\bar{\alpha}_t}x_0) (x - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0 \\
&\quad + \frac{1-\alpha_{t+1}}{1-\bar{\alpha}_t} z z^\top + \frac{1-\alpha_{t+1}}{2(1-\bar{\alpha}_t)^2} \int_{x_0} p_{X_0|X_t}(x_0|x) \left\{ \|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2 - (x - \sqrt{\bar{\alpha}_t}x_0)^\top z \right. \\
&\quad \left. - \int_{x_0} \|x - \sqrt{\bar{\alpha}_t}x_0\|_2^2 p_{X_0|X_t}(x_0|x) dx_0 - \|z\|_2^2 \right\} (x - \sqrt{\bar{\alpha}_t}x_0) (x - \sqrt{\bar{\alpha}_t}x_0)^\top dx_0 + \zeta_{K_2},
\end{aligned}$$

where the residual term ζ_{K_2} satisfies

$$\|\zeta_{K_2}\| \lesssim d^2 \frac{(1-\alpha_{t+1})^{3/2}}{(1-\bar{\alpha}_t)^{1/2}} \log^2 T.$$

Then the second claim (71) immediately follows by recognizing

$$\frac{\partial(s_t^*(x) - \sqrt{\alpha_{t+1}}s_{t+1}^*(w))}{\partial x}$$

$$\begin{aligned}
&= \frac{\sqrt{\alpha_{t+1}}}{1 - \bar{\alpha}_{t+1}} J_{t+1}(w) \frac{\partial w}{\partial x} - \frac{1}{1 - \bar{\alpha}_t} J_t(x) \\
&= \frac{1}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \left(I + \frac{1}{\alpha_{t+1}^{-1} - \bar{\alpha}_t} \left(K_1 K_1^\top - K_2 \right) \right) \left(I + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)} J_t(x) \right) - \frac{1}{1 - \bar{\alpha}_t} J_t(x) \\
&= -\frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^2} J_t(x) + \frac{1 - \alpha_{t+1}}{2(1 - \bar{\alpha}_t)^3} (H_1 + H_4 + H_2 - H_3) + \zeta_{J_t},
\end{aligned}$$

where the residual term ζ_{J_t} satisfies

$$\|\zeta_{J_t}\| \lesssim d^2 \frac{(1 - \alpha_{t+1})^{3/2}}{(1 - \bar{\alpha}_{t+1})^{5/2}} \log^2 T.$$

□

Proof of properties (73). To prove these properties, we first note that Lemma 3 implies that

$$\|z\|_2 \leq \mathbb{E} \left[\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2 \mid X_t = x \right] \lesssim \sqrt{d(1 - \bar{\alpha}_t) \log T}, \quad (92a)$$

$$\begin{aligned}
&\left\| \int_{x_0} p_{X_0|X_t}(x_0 \mid x) (x - \sqrt{\bar{\alpha}_t} x_0) (X_t - \sqrt{\bar{\alpha}_t} x_0)^\top z \, dx_0 \right\|_2 \\
&\lesssim \|z\|_2 \mathbb{E} \left[\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2^2 \mid X_t = x \right] \lesssim \left(\frac{d \log T}{1 - \bar{\alpha}_t} \right)^{3/2}, \quad (92b)
\end{aligned}$$

$$\begin{aligned}
&\left\| \int_{x_0} p_{X_0|X_t}(x_0 \mid x) \|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2^2 (x - \sqrt{\bar{\alpha}_t} x_0 - z) \, dx_0 \right\|_2 \\
&\leq \left\| \mathbb{E} \left[\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2^3 \mid X_t = x \right] \right\|_2 + \|z\|_2 \left\| \mathbb{E} \left[\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2^2 \mid X_t = x \right] \right\|_2 \\
&\lesssim \left(\frac{d \log T}{1 - \bar{\alpha}_t} \right)^{3/2}. \quad (92c)
\end{aligned}$$

Substituting the above bounds into (70) yields the first claim of (73). Similarly, the second claim follows by applying Lemma 3 and utilizing (80) in the context of (71). □

B.3 Proof of Lemma 7

To begin with, it follows from the definition (54) of $\tau(y_T)$ that

$$-\log q_t(y_t) \leq c_\tau d \log T, \quad \forall t < \tau(y_T).$$

Our proof is mainly built upon Lemma 6. Specifically, combining Lemma 3, (33) and the definition (54) of $\tau(y_T)$ gives

$$|B_t| \leq |A_t| \lesssim \frac{1}{1 - \bar{\alpha}_t} \cdot d(1 - \bar{\alpha}_t) \log T \asymp d \log T \quad (93a)$$

$$|C_t| \lesssim \frac{1}{(1 - \bar{\alpha}_t)^2} d^2 (1 - \bar{\alpha}_t)^2 \log^2 T \asymp d^2 \log^2 T \quad (93b)$$

$$|D_t| \leq \frac{\|g_t(x)\|_2^2}{(1 - \bar{\alpha}_t)^2} \int p_{X_0|X_t}(x_0 \mid x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2 \, dx_0 \lesssim d^2 \log^2 T \quad (93c)$$

$$|E_t| \leq \frac{\|g_t(x)\|_2^2}{(1 - \bar{\alpha}_t)^2} \int p_{X_0|X_t}(x_0 \mid x) \|x - \sqrt{\bar{\alpha}_t} x_0\|_2^3 \, dx_0 \lesssim d^2 \log^2 T \quad (93d)$$

for all $t < \tau(y_T)$. As a consequence, the properties (50a) and (50b) in Lemma 6 tell us that

$$\frac{p_{\sqrt{\bar{\alpha}_t} Y_{t-1}}(\psi_t(y_t))}{p_{Y_t}(y_t)} \left(\frac{p_{\sqrt{\bar{\alpha}_t} X_{t-1}}(\psi_t(y_t))}{p_{X_t}(y_t)} \right)^{-1} = \frac{p_{\psi_t(Y_t)}(\psi_t(y_t))}{p_{Y_t}(y_t)} \left(\frac{p_{\sqrt{\bar{\alpha}_t} X_{t-1}}(\psi_t(y_t))}{p_{X_t}(y_t)} \right)^{-1}$$

$$= 1 + O\left(\frac{d^6 \log^6 T}{T^3} + \frac{(\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t}(\Phi_t(y_t)))\sqrt{d \log^3 T}}{T} + \frac{d \log T(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t}(\Phi_t(y_t)))}{T}\right)$$

for all $t < \tau(y_T)$. Give that $y_{t-1} = \frac{1}{\sqrt{\alpha_t}}\psi_t(y_t)$, one can make use of the relation (47) and derive

$$\frac{p_{t-1}(y_{t-1})}{q_{t-1}(y_{t-1})} = \frac{p_t(y_t)}{q_t(y_t)}. \quad (94)$$

$$\left\{1 + O\left(\frac{d^6 \log^6 T}{T^3} + \frac{(\varepsilon_{\text{score},t}(y_t) + \varepsilon_{\text{score},t}(\Phi_t(y_t)))\sqrt{d \log^3 T}}{T} + \frac{d \log T(\varepsilon_{\text{Jacobi},t}(y_t) + \varepsilon_{\text{Jacobi},t}(\Phi_t(y_t)))}{T}\right)\right\}$$

for any $t < \tau(y_T)$. If we employ the shorthand notation $\tau = \tau(y_T)$, then it can be seen that

$$\frac{q_1(y_1)}{p_1(y_1)} = \left\{1 + O\left(\frac{d^6 \log^6 T}{T^2} + S_{\tau-1}(y_{\tau-1})\right)\right\} \frac{q_{\tau-1}(y_{\tau-1})}{p_{\tau-1}(y_{\tau-1})}$$

$$\in \left[\frac{p_{\tau-1}(y_{\tau-1})}{2q_{\tau-1}(y_{\tau-1})}, \frac{2p_{\tau-1}(y_{\tau-1})}{q_{\tau-1}(y_{\tau-1})}\right]. \quad (95a)$$

Repeating this argument also yields

$$\frac{q_t(y_t)}{2p_t(y_t)} \leq \frac{q_1(y_1)}{p_1(y_1)} \leq \frac{2q_t(y_t)}{p_t(y_t)}, \quad \forall t < \tau. \quad (95b)$$

C Analysis for the accelerated DDPM sampler (proof of Theorem 2)

In this section, we turn to the accelerated stochastic sampler and present the proof of Theorem 2.

C.1 Main steps of the proof

Preparation. First, we find it convenient to introduce the following mapping

$$\widehat{\mu}_t^*(x_t) = \frac{1}{\sqrt{\alpha_t}}(x_t + (1 - \alpha_t)s_t^*(x_t)). \quad (96)$$

For any t , introduce the following auxiliary sequences: $H_T \sim \mathcal{N}(0, I_d)$, and

$$H_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left\{H_t + \sqrt{\frac{1 - \alpha_t}{2}}Z_t + (1 - \alpha_t)s_t^*(H_t) - \frac{(1 - \alpha_t)^{3/2}}{\sqrt{2(1 - \bar{\alpha}_t)}}J_t(H_t)Z_t + \sqrt{\frac{1 - \alpha_t}{2}}Z_t^+\right\} \quad (97)$$

$$= \widehat{\mu}_t^*(H_t) + \sqrt{\frac{1 - \alpha_t}{2\alpha_t}}\left(Z_t - \frac{1 - \alpha_t}{1 - \bar{\alpha}_t}J_t(H_t)Z_t + Z_t^+\right) \quad (98)$$

for $t = T, \dots, 1$. We shall also adopt the following notation throughout for notational convenience:

$$\widehat{x}_t := \frac{1}{\sqrt{\alpha_t}}x_t. \quad (99)$$

Step 1: decomposing the KL divergence of interest. Applying Pinsker's inequality and repeating the arguments as in Li et al. (2023, Section 5.3) lead to the following elementary decompositions:

$$\text{TV}(p_{X_1}, p_{Y_1}) \leq \sqrt{\frac{1}{2}\text{KL}(p_{X_1} \parallel p_{Y_1})}, \quad (100)$$

$$\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{t=2}^T \mathbb{E}_{x \sim q_t} \left[\text{KL}\left(p_{X_{t-1} | X_t}(\cdot | x) \parallel p_{Y_{t-1} | Y_t}(\cdot | x)\right) \right]. \quad (101)$$

In particular, the term $\text{KL}(p_{X_T} \parallel p_{Y_T})$ can be readily bounded by Lemma 5 as follows:

$$\text{KL}(p_{X_T} \parallel p_{Y_T}) \lesssim \frac{1}{T^{200}}.$$

As a result, it suffices to bound $\text{KL}(p_{X_{t-1}|X_t}(\cdot|x) \parallel p_{Y_{t-1}|Y_t}(\cdot|x))$ for each $1 < t \leq T$ separately, which we shall accomplish next.

Step 2: bounding the conditional distributions $p_{X_{t-1}|X_t}$ and $p_{H_{t-1}|H_t}$. We now compare the two conditional distributions $p_{X_{t-1}|X_t}$ and $p_{H_{t-1}|H_t}$.

Towards this end, let us first introduce the set below:

$$\mathcal{E} := \left\{ (x_t, x_{t-1}) \mid -\log p_{X_t}(x_t) \leq \frac{1}{2}c_6 d \log T, \|x_{t-1} - \hat{x}_t\|_2 \leq c_3 \sqrt{d(1-\alpha_t) \log T} \right\} \quad (102)$$

with \hat{x}_t defined in (99), and we would like to evaluate both $p_{H_{t-1}|H_t}$ and $p_{X_{t-1}|X_t}$ over the set \mathcal{E} . Regarding $p_{H_{t-1}|H_t}$, we have the following lemma.

Lemma 11. *For every $(x_t, x_{t-1}) \in \mathcal{E}$ as defined in (102), we have*

$$p_{H_{t-1}|H_t}(x_{t-1}|x_t) \propto \exp \left\{ -\frac{\alpha_t}{2(1-\alpha_t)} \left\| \left(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t) \right)^{-1} (x_{t-1} - \hat{\mu}_t^*(x_t)) \right\|_2^2 + O\left(\frac{d^3 \log^5 T}{T^2}\right) \right\}. \quad (103)$$

Turning to $p_{X_{t-1}|X_t}$ over the set \mathcal{E} , we can invoke Li et al. (2023, Lemma 12) to derive the following result.

Lemma 12. *There exists some large enough numerical constant $c_\zeta > 0$ such that: for every $(x_t, x_{t-1}) \in \mathcal{E}$,*

$$p_{X_{t-1}|X_t}(x_{t-1}|x_t) = \frac{1}{(2\pi^{\frac{1-\alpha_t}{\alpha_t}})^{d/2} |\det(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t))|} \cdot \exp \left(-\frac{\alpha_t}{2(1-\alpha_t)} \left\| \left(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t) \right)^{-1} (x_{t-1} - \hat{\mu}_t^*(x_t)) \right\|_2^2 + \zeta_t(x_{t-1}, x_t) \right) \quad (104)$$

holds for some residual term $\zeta_t(x_{t-1}, x_t)$ obeying

$$|\zeta_t(x_{t-1}, x_t)| \leq c_\zeta \frac{d^3 \log^{4.5} T}{T^{3/2}}. \quad (105)$$

Moving beyond the set \mathcal{E} , it suffices to bound the log density ratio $\log \frac{p_{X_{t-1}|X_t}}{p_{H_{t-1}|H_t}}$ for all pairs (x_t, x_{t-1}) , which can be accomplished in a way similar to Li et al. (2023, Lemma 13).

Lemma 13. *For all $(x_t, x_{t-1}) \in \mathbb{R}^d \times \mathbb{R}^d$, we have*

$$\log \frac{p_{X_{t-1}|X_t}(x_{t-1}|x_t)}{p_{H_{t-1}|H_t}(x_{t-1}|x_t)} \leq T^{c_0+2c_R+2} \left\{ \|x_{t-1} - \hat{x}_t\|_2^2 + \|x_t\|_2^2 + 1 \right\}, \quad (106)$$

where c_0 is defined in (13b).

Equipped with Lemmas 11 to 13, one can readily repeat similar arguments as in Li et al. (2023, Step 3, Theorem 3) to derive the following result:

Lemma 14. *For any $1 < t \leq T$, one has*

$$\mathbb{E}_{x_t \sim q_t} \left[\text{KL}(p_{X_{t-1}|X_t}(\cdot|x_t) \parallel p_{H_{t-1}|H_t}(\cdot|x_t)) \right] \lesssim \left(\frac{d^3 \log^{4.5} T}{T^{3/2}} \right)^2. \quad (107)$$

Step 3: quantifying the KL divergence between $p_{H_{t-1}|H_t}$ and $p_{Y_{t-1}|Y_t}$. In the previous step, we have quantified the KL divergence between $p_{X_{t-1}|X_t}$ and $p_{H_{t-1}|H_t}$. Recognizing that H_{t-1} is a first-order approximation of Y_{t-1} using the true score function, we still need to look at the influence of the score estimation errors, for which we resort to the lemma below.

Lemma 15. *For any $1 < t \leq T$, one has*

$$\begin{aligned} & \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1}|X_t}(\cdot | x_t) \parallel p_{Y_{t-1}|Y_t}(\cdot | x_t) \right) \right] - \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1}|X_t}(\cdot | x_t) \parallel p_{H_{t-1}|H_t}(\cdot | x_t) \right) \right] \\ & \lesssim \exp(-c_{20} d \log T) + \frac{d \log^3 T}{T} \mathbb{E}_{X_t \sim q_t} [\varepsilon_{\text{score},t}(X_t)^2] + \frac{d^5 \log^7 T}{T^3}. \end{aligned} \quad (108)$$

Step 4: putting all this together. We are now ready to complete the proof. Substituting (107) and (108) into the decomposition (101) yields

$$\begin{aligned} \text{KL}(p_{X_1} \parallel p_{Y_1}) & \leq \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{t=1}^{T-1} \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1}|X_t}(\cdot | x_t) \parallel p_{H_{t-1}|H_t}(\cdot | x_t) \right) \right] \\ & \quad + \sum_{t=1}^{T-1} \left\{ \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1}|X_t}(\cdot | x_t) \parallel p_{Y_{t-1}|Y_t}(\cdot | x_t) \right) \right] - \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1}|X_t}(\cdot | x_t) \parallel p_{H_{t-1}|H_t}(\cdot | x_t) \right) \right] \right\} \\ & \lesssim \text{KL}(p_{X_T} \parallel p_{Y_T}) + \sum_{2 \leq t \leq T} \frac{d^6 \log^9 T}{T^3} + \frac{d \log^3 T}{T} \sum_{t=2}^T \mathbb{E}_{X_t \sim q_t} [\varepsilon_{\text{score},t}(X_t)^2] \\ & \asymp \frac{d^6 \log^9 T}{T^2} + d \varepsilon_{\text{score}}^2 \log^3 T, \end{aligned}$$

thereby concluding the proof of Theorem 2.

C.2 Proof of Lemma 11

To begin with, we observe that

$$p_{H_{t-1}|H_t}(x_{t-1} | x_t) \propto \exp \left(-\frac{\alpha_t}{1-\alpha_t} (x_{t-1} - \hat{\mu}_t^*(x_t))^\top \text{Var} \left(Z_t - \frac{1-\alpha_t}{1-\bar{\alpha}_t} J_t(H_t) Z_t + Z_t^+ \mid H_t = x_t \right)^{-1} (x_{t-1} - \hat{\mu}_t^*(x_t)) \right).$$

It is easy to verify that

$$\begin{aligned} & \text{Var} \left(Z_t - \frac{1-\alpha_t}{1-\bar{\alpha}_t} J_t(H_t) Z_t + Z_t^+ \mid H_t = x_t \right) \\ & = 2 \left(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t) \right) \left(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t) \right)^\top + \frac{(1-\alpha_t)^2}{2(1-\bar{\alpha}_t)^2} J_t(x_t) J_t(x_t)^\top. \end{aligned}$$

For any $(x_{t-1}, x_t) \in \mathcal{E}$, we can deduce that

$$\|J_t(x_t)\| \lesssim d \log T, \quad (109a)$$

$$\begin{aligned} \|x_{t-1} - \hat{\mu}_t^*(x_t)\|_2 & \leq \|x_{t-1} - \hat{x}_t\|_2 + \frac{1-\alpha_t}{\sqrt{\alpha_t}(1-\bar{\alpha}_t)} \mathbb{E} [\|x_t - \sqrt{\alpha_t} X_0\|_2 \mid X_t = x_t] \\ & \stackrel{(i)}{\lesssim} \sqrt{d(1-\alpha_t) \log T} + \sqrt{\frac{d \log T}{1-\bar{\alpha}_t}} (1-\alpha_t) \asymp \sqrt{d(1-\alpha_t) \log T}, \end{aligned} \quad (109b)$$

where (109a) follows (86b) and (i) both arises from Lemma 3. Taking the above relations together and using the relation (13b), we arrive at

$$\frac{1-\alpha_t}{(1-\bar{\alpha}_t)^2} \|J_t(x_t)\|^2 \|x_{t-1} - \hat{\mu}_t^*(x_t)\|_2^2 \lesssim \frac{d^3 \log^5 T}{T^2},$$

which completes the proof.

C.3 Proof of Lemma 13

According to the expression (103), one has

$$H_{t-1} | H_t = x_t \sim \mathcal{N}\left(\widehat{\mu}_t^*(x_t), \underbrace{\frac{1-\alpha_t}{\alpha_t} \left(I - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)} J_t(x_t)\right)^2 + \frac{(1-\alpha_t)^3}{4\alpha_t(1-\bar{\alpha}_t)^2} J_t(x_t)^2}_{=:\Sigma(\widehat{x}_t)}\right).$$

In order to quantify the above density of interest, we first bound the Jacobian matrix $J_t(x)$ defined in (31). On the one hand, the expression (32) tells us that $J_t(x) \preceq I_d$ for any x , given that the term within the curly bracket in (32) is a negative covariance matrix. On the other hand, $J_t(x)$ can be lower bounded by

$$\begin{aligned} J_t(x) &\succeq -\frac{1}{1-\bar{\alpha}_t} \mathbb{E}\left[(X_t - \sqrt{\bar{\alpha}_t} X_0)(X_t - \sqrt{\bar{\alpha}_t} X_0)^\top \mid X_t = x\right] \\ &\succeq -\frac{\mathbb{E}\left[\|X_t - \sqrt{\bar{\alpha}_t} X_0\|_2^2 \mid X_t = x\right]}{1-\bar{\alpha}_t} I_d \succeq -\frac{2\|x\|_2^2 + 2T^{2c_R}}{1-\bar{\alpha}_t} I_d \\ &\succeq -T^{c_0+1}(\|x\|_2^2 + T^{2c_R}) I_d, \end{aligned}$$

where the second line applies the assumption that $\|X_0\|_2 \leq T^{c_R}$, and the last line invokes the choice (13b). As a consequence, we obtain

$$\Sigma(\widehat{x}_t) \succeq \frac{1-\alpha_t}{\alpha_t} \left(1 - \frac{1-\alpha_t}{2(1-\bar{\alpha}_t)}\right)^2 I_d = \frac{1-\alpha_t}{4\alpha_t} \left(\frac{1-\bar{\alpha}_t + \alpha_t - \bar{\alpha}_t}{1-\bar{\alpha}_t}\right)^2 I_d \succeq \frac{1-\alpha_t}{4\alpha_t} I_d \succeq \frac{1-\alpha_t}{4} I_d; \quad (110a)$$

$$\Sigma(\widehat{x}_t) \preceq \frac{1-\alpha_t}{\alpha_t} T^{2c_0+2} (2\|\widehat{x}_t\|_2^4 + 2T^{4c_R}) I_d + \frac{(1-\alpha_t)^3}{4\alpha_t(1-\bar{\alpha}_t)^2} I_d \preceq 4T^{2c_0+2} (\|\widehat{x}_t\|_2^4 + T^{4c_R}) I_d. \quad (110b)$$

With the above relations in mind, we are ready to bound the density function $p_{H_{t-1} | H_t}(x_{t-1} | x_t)$ for any $x_t, x_{t-1} \in \mathbb{R}^d$. It is seen from (103) that

$$\begin{aligned} \log \frac{1}{p_{H_{t-1} | H_t}(x_{t-1} | x_t)} &= \frac{(x_{t-1} - \widehat{\mu}_t^*(x_t))^\top (\Sigma(\widehat{x}_t))^{-1} (x_{t-1} - \widehat{\mu}_t^*(x_t))}{2} + \frac{1}{2} \log \det(\Sigma(\widehat{x}_t)) + \frac{d}{2} \log(2\pi) \\ &\leq \frac{2\|x_{t-1} - \widehat{\mu}_t^*(x_t)\|_2^2}{1-\alpha_t} + \frac{d}{2} \log\left(8\pi T^{2c_0+2} (\|\widehat{x}_t\|_2^4 + T^{4c_R})\right) \\ &\leq 2T^{c_0+1} \left\{2\|x_{t-1} - \widehat{x}_t\|_2^2 + \|x_t\|_2^2 + T^{2c_R}\right\} + \frac{d}{2} \log\left(8\pi T^{2c_0+2} (\|\widehat{x}_t\|_2^4 + T^{4c_R})\right) \\ &\leq T^{c_0+2c_R+2} \left\{\|x_{t-1} - \widehat{x}_t\|_2^2 + \|x_t\|_2^2 + 1\right\}, \end{aligned}$$

where the second inequality results from (110), and the third inequality makes use of (13b) and the fact that

$$\begin{aligned} \|x_{t-1} - \widehat{\mu}_t^*(x_t)\|_2^2 &\leq 2\|x_{t-1} - \widehat{x}_t\|_2^2 + 2\|\widehat{x}_t - \widehat{\mu}_t^*(x_t)\|_2^2 \\ &= 2\|x_{t-1} - \widehat{x}_t\|_2^2 + 2\left(\frac{1-\alpha_t}{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_t)}\right)^2 \left\|\int_{x_0} p_{X_0 | X_t}(x_0 | x_t) (x_t - \sqrt{\bar{\alpha}_t} x_0) dx_0\right\|_2^2 \\ &\leq 2\|x_{t-1} - \widehat{x}_t\|_2^2 + \frac{2(1-\alpha_t)^2}{\alpha_t(1-\bar{\alpha}_{t-1})^2} \sup_{x_0: \|x_0\|_2 \leq T^{c_R}} \|x_t - \sqrt{\bar{\alpha}_t} x_0\|_2^2 \\ &\leq 2\|x_{t-1} - \widehat{x}_t\|_2^2 + \frac{64c_1^2 \log^2 T}{T^2} \left(2\|x_t\|_2^2 + 2\bar{\alpha}_t T^{2c_R}\right) \\ &\leq 2\|x_{t-1} - \widehat{x}_t\|_2^2 + \|x_t\|_2^2 + T^{2c_R}. \end{aligned} \quad (111)$$

Given that $\log \frac{p_{X_{t-1} | X_t}(x_{t-1} | x_t)}{p_{H_{t-1} | H_t}(x_{t-1} | x_t)} \leq \log \frac{1}{p_{H_{t-1} | H_t}(x_{t-1} | x_t)}$, we have concluded the proof.

C.4 Proof of Lemma 14

Firstly, it follows from Lemma 11 and Lemma 12 that: for any $(x_t, x_{t-1}) \in \mathcal{E}$,

$$\frac{p_{X_{t-1}|X_t}(x_{t-1}|x_t)}{p_{H_{t-1}|H_t}(x_{t-1}|x_t)} = \exp\left(O\left(\frac{d^3 \log^{4.5} T}{T^{3/2}}\right)\right) \quad (112)$$

$$= 1 + O\left(\frac{d^3 \log^{4.5} T}{T^{3/2}}\right) \in \left[\frac{1}{2}, 2\right], \quad (113)$$

which further allows one to derive

$$\begin{aligned} & \mathbb{E}_{x_t \sim q_t} \left[\text{KL}\left(p_{X_{t-1}|X_t}(\cdot|x_t) \parallel p_{H_{t-1}|H_t}(\cdot|x_t)\right) \right] \\ &= \left(\int_{\mathcal{E}} + \int_{\mathcal{E}^c} \right) p_{X_t}(x_t) p_{X_{t-1}|X_t}(x_{t-1}|x_t) \log \frac{p_{X_{t-1}|X_t}(x_{t-1}|x_t)}{p_{H_{t-1}|H_t}(x_{t-1}|x_t)} dx_{t-1} dx_t, \\ &\stackrel{(i)}{=} \int_{\mathcal{E}} p_{X_t}(x_t) \left\{ p_{X_{t-1}|X_t}(x_{t-1}|x_t) - p_{H_{t-1}|H_t}(x_{t-1}|x_t) \right. \\ &\quad \left. + p_{X_{t-1}|X_t}(x_{t-1}|x_t) \cdot O\left(\left(\frac{p_{H_{t-1}|H_t}(x_{t-1}|x_t)}{p_{X_{t-1}|X_t}(x_{t-1}|x_t)} - 1\right)^2\right) \right\} dx_{t-1} dx_t \\ &\quad + \int_{\mathcal{E}^c} p_{X_t}(x_t) p_{X_{t-1}|X_t}(x_{t-1}|x_t) \log \frac{p_{X_{t-1}|X_t}(x_{t-1}|x_t)}{p_{H_{t-1}|H_t}(x_{t-1}|x_t)} dx_{t-1} dx_t \\ &\stackrel{(ii)}{=} \int_{\mathcal{E}} p_{X_t}(x_t) \left\{ p_{X_{t-1}|X_t}(x_{t-1}|x_t) - p_{H_{t-1}|H_t}(x_{t-1}|x_t) + p_{X_{t-1}|X_t}(x_{t-1}|x_t) O\left(\frac{d^6 \log^9 T}{T^3}\right) \right\} dx_{t-1} dx_t \\ &\quad + \int_{\mathcal{E}^c} p_{X_t}(x_t) p_{X_{t-1}|X_t}(x_{t-1}|x_t) \left\{ 2T(\|x_t\|_2^2 + \|x_{t-1} - \hat{x}_t\|_2^2 + T^{2c_R}) \right\} dx_{t-1} dx_t. \end{aligned} \quad (114)$$

Here, (i) invokes the basic fact that: if $\left|\frac{p_Y(x)}{p_X(x)} - 1\right| < \frac{1}{2}$, then the Taylor expansion gives

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

and in (ii) we apply (113) and Lemma 13.

Next, we would like to bound each term on the right-hand side of (114) separately. In view of the definition of the set \mathcal{E} (cf. (102)), one has

$$\begin{aligned} \mathbb{P}((X_t, X_{t-1}) \notin \mathcal{E}) &= \int_{(x_t, x_{t-1}) \notin \mathcal{E}} p_{X_{t-1}}(x_{t-1}) p_{X_t|X_{t-1}}(x_t|x_{t-1}) dx_{t-1} dx_t \\ &= \int_{(x_t, x_{t-1}) \notin \mathcal{E}} p_{X_{t-1}}(x_{t-1}) \frac{1}{(2\pi(1-\alpha_t))^{d/2}} \exp\left(-\frac{\|x_t - \sqrt{\alpha_t} x_{t-1}\|_2^2}{2(1-\alpha_t)}\right) dx_{t-1} dx_t \\ &\leq \exp(-c_3 d \log T), \end{aligned} \quad (115)$$

and similarly,

$$\int_{(x_{t-1}, x_t) \notin \mathcal{E}} p_{X_t}(x_t) p_{X_{t-1}|X_t}(x_{t-1}|x_t) \left(2T(\|x_t\|_2^2 + \|x_{t-1} - \hat{x}_t\|_2^2) + T^{2c_R}\right) dx_{t-1} dx_t \leq \exp(-c_3 d \log T). \quad (116)$$

In addition, for every (x_t, x_{t-1}) obeying $\|x_{t-1} - x_t/\sqrt{\alpha_t}\|_2 > c_3 \sqrt{d(1-\alpha_t) \log T}$ and $-\log p_{X_t}(x_t) \leq \frac{1}{2} c_6 d \log T$, it follows from the definition (96) of $\hat{\mu}_t^*(\cdot)$ that

$$\|x_{t-1} - \hat{\mu}_t^*(x_t)\|_2 = \left\| x_{t-1} - \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1-\alpha_t}{\sqrt{\alpha_t}(1-\bar{\alpha}_t)} \mathbb{E}\left[x_t - \sqrt{\alpha_t} X_0 \mid X_t = x_t\right] \right\|_2 \quad (117)$$

$$\begin{aligned}
&\geq \left\| x_{t-1} - \frac{1}{\sqrt{\alpha_t}} x_t \right\|_2 - \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \mathbb{E} \left[\|x_t - \sqrt{\bar{\alpha}_t} X_0\|_2 \mid X_t = x_t \right] \\
&\geq c_3 \sqrt{d(1 - \alpha_t) \log T} - 6\bar{c}_5 \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \sqrt{d \log T} \\
&= \left(c_3 - 6\bar{c}_5 \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \right) \sqrt{d(1 - \alpha_t) \log T} \geq \frac{c_3}{2} \sqrt{d(1 - \alpha_t) \log T}, \tag{118}
\end{aligned}$$

where the third line results from (37a) in Lemma 3, and the last line applies (33) and holds true as long as c_3 is large enough. Taking this result together with Lemma 11 reveals that: for any x_t obeying $-\log p_{X_t}(x_t) \leq \frac{1}{2} c_6 d \log T$, one has

$$\int_{x_{t-1}: \|x_{t-1} - x_t / \sqrt{\alpha_t}\|_2 > c_3 \sqrt{d(1 - \alpha_t) \log T}} p_{H_{t-1} | H_t}(x_{t-1} | x_t) dx_{t-1} \leq \exp \left(- \frac{c_3}{2} d \log T \right). \tag{119}$$

Combine (115) and (119) to arrive at

$$\begin{aligned}
&\left| \int_{\mathcal{E}} p_{X_t}(x_t) \left\{ p_{X_{t-1} | X_t}(x_{t-1} | x_t) - p_{H_{t-1} | H_t}(x_{t-1} | x_t) \right\} dx_{t-1} dx_t \right| \leq \mathbb{P}((X_t, X_{t-1}) \notin \mathcal{E}) \\
&\quad + \int_{\log p_{X_t}(x_t) \leq \frac{1}{2} c_6 d \log T, \|x_{t-1} - x_t / \sqrt{\alpha_t}\|_2 > c_3 \sqrt{d(1 - \alpha_t) \log T}} p_{X_t}(x_t) p_{H_{t-1} | H_t}(x_{t-1} | x_t) dx_{t-1} dx_t \\
&\leq 2 \exp \left(- \frac{c_3}{2} d \log T \right). \tag{120}
\end{aligned}$$

To finish up, plugging (116) and (120) into (114) yields: for each $t \geq 2$,

$$\mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1} | X_t}(\cdot | x_t) \parallel p_{Y_{t-1}^* | Y_t}(\cdot | x_t) \right) \right] \lesssim \frac{d^6 \log^9 T}{T^3} + 3 \exp \left(- \frac{c_3}{2} d \log T \right) \lesssim \frac{d^6 \log^9 T}{T^3}. \tag{121}$$

C.5 Proof of Lemma 15

We first introduce the following notation:

$$\begin{aligned}
\mu_t(x_t, z_t) &:= \frac{1}{\sqrt{\alpha_t}} \left(x_t + \sqrt{\frac{1 - \alpha_t}{2}} z_t + (1 - \alpha_t) s_t \left(x_t + \sqrt{\frac{1 - \alpha_t}{2}} z_t \right) \right); \\
\mu_t^*(x_t, z_t) &:= \frac{1}{\sqrt{\alpha_t}} \left(x_t + \sqrt{\frac{1 - \alpha_t}{2}} z_t + (1 - \alpha_t) s_t^*(x_t) - \frac{(1 - \alpha_t)^{3/2}}{\sqrt{2}(1 - \bar{\alpha}_t)} J_t(x_t) z_t \right).
\end{aligned}$$

In the sequel, we shall use μ_t and μ_t^* to denote $\mu_t(x_t, z_t)$ and $\mu_t^*(x_t, z_t)$, respectively, for simplicity, as long as it is clear from the context. It is observed that

$$\begin{aligned}
&\mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1} | X_t}(\cdot | x_t) \parallel p_{Y_{t-1} | Y_t}(\cdot | x_t) \right) \right] - \mathbb{E}_{x_t \sim q_t} \left[\text{KL} \left(p_{X_{t-1} | X_t}(\cdot | x_t) \parallel p_{H_{t-1} | H_t}(\cdot | x_t) \right) \right] \\
&= \int p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) \log \frac{p_{H_{t-1} | H_t}(x_{t-1} | x_t)}{p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)} dx_{t-1} dx_t \\
&= \int p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) p_{Z_t | H_{t-1}, H_t}(z_t | x_{t-1}, x_t) \log \frac{p_{H_{t-1} | H_t}(x_{t-1} | x_t)}{p_{Y_{t-1} | Y_t}(x_{t-1} | x_t)} dz_t dx_{t-1} dx_t \\
&\stackrel{(i)}{\leq} \int p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) p_{Z_t | H_{t-1}, H_t}(z_t | x_{t-1}, x_t) \log \frac{p_{H_{t-1} | H_t, Z_t}(x_{t-1} | x_t, z_t)}{p_{Y_{t-1} | Y_t, Z_t}(x_{t-1} | x_t, z_t)} dz_t dx_{t-1} dx_t \\
&\stackrel{(ii)}{=} \int p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) p_{Z_t | H_{t-1}, H_t}(z_t | x_{t-1}, x_t) \frac{\alpha_t}{(1 - \alpha_t)} \left(\|x_{t-1} - \mu_t\|_2^2 - \|x_{t-1} - \mu_t^*\|_2^2 \right) dz_t dx_{t-1} dx_t \\
&= \underbrace{\int p_{X_t}(x_t) p_{X_{t-1} | X_t}(x_{t-1} | x_t) p_{Z_t | H_{t-1}, H_t}(z_t | x_{t-1}, x_t) \frac{\alpha_t}{(1 - \alpha_t)} \|\mu_t^* - \mu_t\|_2^2 dz_t dx_{t-1} dx_t}_{\mathcal{H}_1}
\end{aligned}$$

$$+ \underbrace{\int p_{X_t}(x_t) p_{X_{t-1}|X_t}(x_{t-1}|x_t) p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t) \frac{2\alpha_t}{(1-\alpha_t)} (\mu_t^* - \mu_t)^\top (x_{t-1} - \mu_t^*) dz_t dx_{t-1} dx_t}_{\mathcal{H}_2}.$$

Here, (i) follows the property of KL divergence that

$$\int p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t) \log \frac{p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t)}{p_{Z_t|Y_{t-1},Y_t}(z_t|x_{t-1},x_t)} dz_t \geq 0,$$

whereas (ii) results from the following expressions:

$$\begin{aligned} p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) &\propto \exp\left(-\frac{\alpha_t}{(1-\alpha_t)} \left\| (x_{t-1} - \mu_t^*(x_t)) \right\|_2^2\right) \\ p_{Y_{t-1}|Y_t,Z_t}(x_{t-1}|x_t,z_t) &\propto \exp\left(-\frac{\alpha_t}{(1-\alpha_t)} \left\| (x_{t-1} - \mu_t(x_t)) \right\|_2^2\right). \end{aligned}$$

To bound \mathcal{H}_1 , we first note that

$$\begin{aligned} \frac{1}{1-\alpha_t} \left\| \mu_t - \mu_t^* \right\|_2^2 &= (1-\alpha_t) \cdot \\ &\left\| s_t\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) - s_t^*\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) + s_t^*\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) - s_t^*(x_t) + \frac{(1-\alpha_t)^{1/2}}{\sqrt{2}(1-\bar{\alpha}_t)} J_t(x_t) z_t \right\|_2^2 \\ &\leq (1-\alpha_t) \left\| s_t\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) - s_t^*\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) \right\|_2^2 + \\ &\quad (1-\alpha_t) \left\| s_t^*\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) - s_t^*(x_t) + \frac{(1-\alpha_t)^{1/2}}{\sqrt{2}(1-\bar{\alpha}_t)} J_t(x_t) z_t \right\|_2^2 \\ &\lesssim \frac{\log T}{T} \varepsilon_{\text{score},t}\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right)^2 + \frac{d^5 \log^7 T}{T^3} \end{aligned}$$

where the last inequality follows from the definition (40), the relation (13b), and the fact that

$$\begin{aligned} (1-\alpha_t) &\left\| s_t^*\left(x_t + \sqrt{\frac{1-\alpha_t}{2}} z_t\right) - s_t^*(x_t) + \frac{(1-\alpha_t)^{1/2}}{\sqrt{2}(1-\bar{\alpha}_t)} J_t(x_t) z_t \right\|_2^2 \\ &= \frac{(1-\alpha_t)^2}{2(1-\bar{\alpha}_t)^2} \left\| \int_0^1 \left(J_t(x_t) - J_t\left(x_t + \gamma \sqrt{\frac{1-\alpha_t}{2}} z_t\right) \right) z_t d\gamma \right\|_2^2 \lesssim \frac{d^5 \log^7 T}{T^3}. \end{aligned}$$

Here, the last inequality holds by invoking the property (86c) that for $(x, x_{t-1}) \in \mathcal{E}$,

$$\|J_t(x) - J_t(x_t)\| \leq \sup_{u \in \mathbb{S}^{d-1}} |u^\top (J_t(x) - J_t(x_t)) u| \lesssim d^{3/2} \|x - x_t\|_2 \log^{3/2} T. \quad (122)$$

For the case with $(x, x_{t-1}) \notin \mathcal{E}$, this term will decay exponentially fast and can be bounded analogously. Furthermore, we observe that

$$p_{\Phi_t(X_t, Z_t)}(x) = (\pi(2(1-\bar{\alpha}_t) + 1 - \alpha_t))^{-d/2} \int p_{X_0}(x_0) \exp\left(-\frac{\|x - \sqrt{\bar{\alpha}_t} x_0\|_2^2}{2(1-\bar{\alpha}_t) + 1 - \alpha_t}\right) dx_0 \asymp p_{X_t}(x),$$

which in turn implies that

$$\begin{aligned} \mathcal{H}_1 &\leq \left(1 + O\left(\frac{d^3 \log^{4.5} T}{T^{3/2}}\right)\right) \int p_{X_t}(x_t) p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) p_{Z_t}(z_t) \frac{\alpha_t}{1-\alpha_t} \left\| \mu_t - \mu_t^* \right\|_2^2 dx_t dx_{t-1} dz_t \\ &\lesssim \mathbb{E}_{x^+ \sim \Phi_t(X_t, Z_t)} \left[\frac{\log T}{T} \varepsilon_{\text{score},t}(x^+)^2 \right] + \frac{d^5 \log^7 T}{T^3} \end{aligned}$$

$$\asymp \frac{d \log^3 T}{T} \mathbb{E}_{X_t \sim q_t} [\varepsilon_{\text{score},t}(X_t)^2] + \frac{d^5 \log^7 T}{T^3}$$

We then decompose \mathcal{H}_2 as follows

$$\begin{aligned} \mathcal{H}_2 &= \int p_{X_t}(x_t) (p_{X_{t-1}|X_t}(x_{t-1}|x_t) - p_{H_{t-1}|H_t}(x_{t-1}|x_t)) p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t) \\ &\quad \cdot \frac{2\alpha_t}{(1-\alpha_t)} (\mu_t^* - \mu_t)^\top (x_{t-1} - \mu_t^*) dz_t dx_{t-1} dx_t \\ &+ \int p_{X_t}(x_t) p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) p_{Z_t}(z_t) \cdot \frac{2\alpha_t}{(1-\alpha_t)} (\mu_t^* - \mu_t)^\top (x_{t-1} - \mu_t^*) dz_t dx_{t-1} dx_t \\ &\stackrel{(i)}{=} \int p_{X_t}(x_t) (p_{X_{t-1}|X_t}(x_{t-1}|x_t) - p_{H_{t-1}|H_t}(x_{t-1}|x_t)) p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t) \\ &\quad \cdot \frac{2\alpha_t}{(1-\alpha_t)} (\mu_t^* - \mu_t)^\top (x_{t-1} - \mu_t^*) dz_t dx_{t-1} dx_t \\ &= \left(\int_{\mathcal{E}} + \int_{\mathcal{E}^c} \right) p_{X_t}(x_t) (p_{X_{t-1}|X_t}(x_{t-1}|x_t) - p_{H_{t-1}|H_t}(x_{t-1}|x_t)) p_{Z_t|H_{t-1},H_t}(z_t|x_{t-1},x_t) \\ &\quad \cdot \frac{2\alpha_t}{(1-\alpha_t)} (\mu_t^* - \mu_t)^\top (x_{t-1} - \mu_t^*) dz_t dx_{t-1} dx_t, \end{aligned}$$

where (i) follows the fact that $\mathbb{E}[H_{t-1} - \mu_t^* | H_t, Z_t] = 0$. In the following, we mainly focus on the term $\int_{\mathcal{E}}$ denoted as \mathcal{K}_1 , since the other term can be bounded similarly as (Li et al., 2023, Lemma 10) and is exponentially small.

$$\begin{aligned} \mathcal{K}_1 &\stackrel{(i)}{\lesssim} \frac{d^3 \log^{4.5} T}{T^{3/2}} \int_{\mathcal{E}} p_{X_t}(x_t) p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) P_{Z_t}(z_t) \|x_{t-1} - \mu_t^*\|_2 \frac{1}{1-\alpha_t} \|\mu_t - \mu_t^*\|_2 dx_{t-1} dx_t \\ &\stackrel{(ii)}{\lesssim} \frac{d^3 \log^{4.5} T}{T^{3/2}} \sqrt{\mathcal{K}_2 \mathcal{K}_3}. \end{aligned} \tag{123}$$

Here, we have

$$\begin{aligned} \mathcal{K}_2 &= \int_{\mathcal{E}} p_{X_t}(x_t) p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) P_{Z_t}(z_t) \|x_{t-1} - \mu_t^*\|_2^2 dx_{t-1} dx_t dz_t \\ &\leq \frac{d(1-\alpha_t)}{\alpha_t} \lesssim \frac{d \log T}{T}; \\ \mathcal{K}_3 &= \int_{\mathcal{E}} p_{X_t}(x_t) p_{H_{t-1}|H_t,Z_t}(x_{t-1}|x_t,z_t) P_{Z_t}(z_t) \frac{1}{(1-\alpha_t)^2} \|\mu_t - \mu_t^*\|_2^2 dx_{t-1} dx_t dz_t \\ &\lesssim \mathbb{E}_{X_t \sim q_t} [\varepsilon_{\text{score},t}(X_t)^2] + \frac{d^5 \log^6 T}{T^2}. \end{aligned}$$

Therefore, we arrive at

$$\mathcal{K}_1 \lesssim \frac{d^{3.5} \log^5 T}{T^2} \mathbb{E}_{X_t \sim q_t} [\varepsilon_{\text{score},t}(X_t)^2] + \frac{d^6 \log^8 T}{T^3}.$$

Taking the above bounds on \mathcal{H}_1 and \mathcal{K}_1 together completes the proof.