

Surprises in sampling Gaussian mixtures using diffusion models

Gen Li^{*†}

Changxiao Cai^{*‡}

Yuting Wei[§]

March 23, 2025

Abstract

Diffusion models are distinguished by their exceptional generative performance, particularly in producing high-quality samples through iterative denoising. While current theory suggests that the number of denoising steps required for accurate sample generation should scale linearly with data dimension, this does not reflect the practical efficiency of widely used algorithms like Denoising Diffusion Probabilistic Models (DDPM). This paper investigates the effectiveness of diffusion models in learning Gaussian Mixture Models (GMMs). Our main result shows that, with perfect score estimates, DDPM requires at most $\tilde{O}(1/\varepsilon)$ iterations to achieve an ε -accurate distribution in total variation (TV) distance, independent of both the ambient dimension d and the number of components K , up to logarithmic factors. Furthermore, this result remains robust to score estimation errors. These findings highlight the remarkable effectiveness of diffusion models for GMMs, even in high-dimensional settings, shedding lights on their capabilities.

Contents

1	Introduction	2
1.1	Diffusion models and sampling efficiency	2
1.2	Learning GMMs using diffusion models	3
1.3	Other related works	3
1.4	Notation	4
2	Preliminaries for diffusion models	4
3	Main results	5
4	Analysis	7
4.1	Preliminaries	8
4.2	Step 1: Constructing auxiliary processes	9
4.3	Step 2: Bounding discretization error	10
4.4	Step 3: Relating to score estimation error	11
5	Discussion	12
A	Proof of technical lemmas	13
A.1	Proof of Lemma 2	13
A.2	Proof of Lemma 3	14
A.3	Proof of Lemma 4	17
A.4	Proof of Lemma 5	18

*The authors contributed equally. Corresponding author: Yuting Wei.

[†]Department of Statistics, The Chinese University of Hong Kong, genli@cuhk.edu.hk.

[‡]Department of Industrial and Operations Engineering, University of Michigan, cxcai@umich.edu.

[§]Department of Statistics and Data Science, University of Pennsylvania, ytwei@wharton.upenn.edu.

1 Introduction

Diffusion models have garnered significant attention for their remarkable generative capabilities, producing high-quality samples with enhanced stability (Diakonikolas et al., 2018; Dhariwal and Nichol, 2021; Song et al., 2020c; Ramesh et al., 2022). Compared to methods like Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), which generate samples in a single forward pass, diffusion models are designed to iteratively denoise samples over hundreds or thousands of steps. A prominent example is the widely used Denoising Diffusion Probabilistic Models (DDPM) sampler (Ho et al., 2020). The current theory suggests the number of denoising steps required for accurate sample generation should scale at least linearly with the data dimension (Chen et al., 2022; Benton et al., 2024) in order to learn the distribution accurately. While various acceleration schemes have been proposed in literature (see, e.g. Li and Cai (2024); Li et al. (2024a); Li and Jiao (2024); Wu et al. (2024b); Huang et al. (2024b,a); Taheri and Lederer (2025)), in practical applications such as high-resolution image synthesis, where the dimensionality of the data can be extremely large, DDPM often requires far fewer steps than predicted by theory while maintaining excellent sample quality.

This gap between theoretical complexity bounds and empirical performance has inspired a strand of recent research, investigating whether diffusion models have implicitly exploited structural properties of real-world data to circumvent worst-case complexity bounds. A growing line of works have shed light on this question by showing that diffusion models, in its original form, can automatically adapt to the intrinsic dimension of the target distribution without explicitly modeling its low-dimensional structure. Notably, prior work has examined cases where the data lies in low-dimensional linear spaces, low-dimensional manifolds, or distributions whose support have small covering number (Li and Yan, 2024a; Tang and Yang, 2024; Huang et al., 2024c; Potapchik et al., 2024; Liang et al., 2025). In this work, we take a different perspective, and explore this question by focusing on a fundamental and well-studied statistical model: Gaussian Mixture Models (GMMs). GMMs serve as a cornerstone of statistical modeling and have been widely used to approximate complex distributions. Formally, we consider the setting where the target distribution is a mixture of isotropic Gaussians:

$$X_0 \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \sigma^2 I_d), \quad (1)$$

where $\{\pi_k\}$ are mixture weights satisfying $\pi_k \in (0, 1)$ and $\sum_{k=1}^K \pi_k = 1$. The study of Gaussian Mixture Models (GMMs) dates back to Pearson (1894), and a vast body of literature has since explored various aspects of GMMs, including parameter estimation, distribution learning, information-theoretic limits, computational efficiency and etc. This paper studies the performance of diffusion models in their original form when they are used to learn a GMM. We refer readers to a more detailed exposition of related work in Section 1.3.

1.1 Diffusion models and sampling efficiency

In a nutshell, diffusion models consist of two processes: a forward process and a backward process. In the forward process, noise is gradually added to the data, transforming it into a noise-like distribution chosen *a priori* (e.g., a Gaussian distribution). Mathematically, given an initial sample $X_0 \in \mathbb{R}^d$ from the target distribution p_{data} , this transformation follows

$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} W_t, \quad t = 1, 2, \dots, T, \quad (2)$$

where $\alpha_t \in (0, 1), t \geq 1$ denotes the learning rates and $W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d), t \geq 1$ are i.i.d. standard Gaussian vectors in \mathbb{R}^d . In the backward process, starting from $Y_T \sim \mathcal{N}(0, I_d)$, diffusion models iteratively denoise Y_T to approximate p_{data} . Classical results from stochastic differential equations (SDE) theory (e.g. Anderson (1982); Haussmann and Pardoux (1986)) show that under mild conditions, recovering p_{data} is possible provided access to the (Stein) score function $s_t^*(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ for all $1 \leq t \leq T$, defined as

$$s_t^*(x) := \nabla \log p_{X_t}(x), \quad \forall x \in \mathbb{R}^d. \quad (3)$$

Given the complexity of developing a comprehensive end-to-end theory, a divide-and-conquer approach — pioneered by (Chen et al., 2022) — has become standard, separating the score learning phase (i.e., estimating

score functions reliably from training data) from the generative sampling phase (i.e., generating new data instances based on the estimated scores). The quality of the sampler in terms of its discrepancy to the target distribution depends on the errors from both phases. Over the past several years, the theoretical community has made significant progress in understanding both phases. Notably, for the sampling phase, convergence theory has been established for various samplers (Liu et al., 2022; Lee et al., 2023; Chen et al., 2023a; Li et al., 2023; Chen et al., 2023c; Tang and Zhao, 2024; Liang et al., 2024a; Huang et al., 2024a; Gao and Zhu, 2024), especially DDPM and Denoising Diffusion Implicit Models (DDIM) which are widely adopted in practice (Ho et al., 2020; Song et al., 2020a). For DDPM, Benton et al. (2024) establishes an iteration complexity of $\tilde{O}(d/\varepsilon^2)$ ¹ in Kullback Leibler (KL) divergence, and Li and Yan (2024b) shows a complexity of $\tilde{O}(d/\varepsilon)$ in total variation (TV) distance. When it comes to the DDIM sampler or the probability flow ODE, notably, an $\tilde{O}(d/\varepsilon)$ iteration complexity has been established in Li et al. (2024b).

1.2 Learning GMMs using diffusion models

In the context of GMMs, several recent works have contributed towards unraveling the capabilities of diffusion models. In particular, inspired by diffusion models, Shah et al. (2023) introduced an algorithm designed for GMMs that achieves polynomial time complexity in d , provided the component centers are well-separated. Liang et al. (2024b) established an iteration complexity of $\tilde{O}(d/\varepsilon^2)$ for obtaining an ε -accurate distribution measured in TV distance by analyzing the Lipschitz and second moments of GMMs. Additionally, Wu et al. (2024a); Chidambaram et al. (2024) investigated the role of guidance in diffusion models. Two exciting recent works (Chen et al., 2024; Gatmiry et al., 2024) proposed using piecewise polynomial regression to estimate the score functions, and they combined this with existing convergence result for DDPM to develop an end-to-end theory for DDPM. Notably, in these works, the number of diffusion steps scales also linearly with d . Further, Wang et al. (2024) explored diffusion models for mixtures of low-rank Gaussians. Despite these advancements, a fundamental question remains open: Can diffusion models achieve efficient sampling when the target distribution is a GMM?

A glimpse of our main contributions. This paper investigates learning a GMM without imposing the well-separated component assumption, a setting where parameter estimation is inherently challenging. Our main result provides a non-asymptotic characterization of DDPM’s iteration complexity for learning an ε -accurate distribution in TV distance. We prove that, given access to perfect score estimates, DDPM requires at most

$$\tilde{O}\left(\frac{1}{\varepsilon}\right),$$

number of iterations. Remarkably, this iteration complexity is independent of both the ambient dimension d and the number of components K , up to some logarithmic factors. Moreover, our result is robust to score estimation errors: the TV distance between the learned distribution and the target distribution scales proportionally to the score estimation error, modulo logarithmic factor. This leads to a surprising insight:

Even in ultra-high-dimensional settings, diffusion models remain highly effective in sampling from GMMs.

1.3 Other related works

Learning GMMs. GMMs are fundamental statistical models that bear a well-established body of research from both statistics and computer science communities. One major line of research focuses on parameter estimation with some separation conditions. Partial examples include Dasgupta (1999); Vempala and Wang (2004); Arora and Kannan (2005); Kalai et al. (2010); Hsu and Kakade (2013); Diakonikolas et al. (2018); Hopkins and Li (2018); Kothari et al. (2018); Liu and Li (2022).

Our work is more closely related to the density estimation perspective, where no separation conditions are imposed (e.g. Diakonikolas and Kane (2020); Moitra and Valiant (2010); Dwivedi et al. (2020); Bakshi et al. (2022); Ho and Nguyen (2016)). In this setting, parameter estimation is information-theoretically infeasible, yet accurate density estimation is still possible. The information theoretical limit of this problem

¹The definition for $O(\cdot)$ and $\tilde{O}(\cdot)$ notation can be found in Section 1.4.

is first characterized in [Ashtiani et al. \(2018\)](#) up to logarithmic factors, with a brute-forth algorithm that scales exponentially in both d and K . For one-dimensional Gaussian mixtures, [Chen \(1995\)](#); [Heinrich and Kahn \(2018\)](#); [Wu and Yang \(2020\)](#) obtained optimal estimation rates and practical algorithms, which were generalized to the high-dimensional case for mixtures of spherical Gaussians with a computationally efficient algorithm in [Doss et al. \(2023\)](#). Beyond finite mixtures, when mixing distribution is an arbitrary probability measure (e.g. [Genovese and Wasserman \(2000\)](#); [Ghosal and Van Der Vaart \(2001\)](#)), [Saha and Guntuboyina \(2020\)](#); [Polyanskiy and Wu \(2020\)](#); [Kim and Guntuboyina \(2022\)](#) established convergence rates and adaptivity, regarding the non-parametric maximum likelihood estimator, generalizing the one-dimension results in [Zhang \(2009\)](#).

Score estimation. As mentioned earlier, score estimation plays a crucial role in diffusion models. [Hyvärinen \(2005\)](#) introduced an integration-by-parts-based approach to simplify score estimation. More recently, [Song et al. \(2020b\)](#) proposed training neural networks to learn score functions by minimizing the score matching objective. The theoretical guarantees for score estimation using neural networks have been analyzed across various distributional settings, including sub-Gaussian distributions ([Cole and Lu, 2024](#)), graphical models ([Mei and Wu, 2023](#)), low-dimensional structured distributions ([Chen et al., 2023b](#); [Kwon et al., 2025](#); [De Bortoli, 2022](#)), and Besov function space ([Oko et al., 2023](#)). These guarantees are often achieved by designing neural architectures that well approximate the true score function. Other than neural networks, classical methods such as kernel-based approaches and empirical Bayes smoothing have also been studied for score estimation ([Cai and Li, 2025](#); [Wibisono et al., 2024](#); [Zhang et al., 2024](#); [Dou et al., 2024](#)). These methods have been shown to achieve minimax-optimal rates under some smoothness assumptions. Furthermore, [Feng et al. \(2024\)](#) demonstrated that statistical procedures based on score matching can achieve minimal asymptotic covariance for convex M-estimation.

1.4 Notation

For any a , the Dirac delta function $\delta_a(x)$ is defined as $\delta_a(x) = \infty$ if $x = a$ and $\delta_a(x) = 0$ otherwise. For positive integer $N > 0$, let $[N] := \{1, \dots, N\}$. In addition, given any matrix A , we use $\|A\|$, $\text{tr}(A)$, and $\det(A)$ to denote the spectral norm, trace, and determinant of the matrix, respectively. Next, we recall the definitions of the KL divergence and TV distance to measure the discrepancies between two distributions. Specifically, for random vectors X and Y with probability density functions p_X and p_Y , let

$$\begin{aligned} \text{KL}(X \parallel Y) &\equiv \text{KL}(p_X \parallel p_Y) = \int p_X(x) \log \left(\frac{p_X(x)}{p_Y(x)} \right) dx, \\ \text{TV}(X, Y) &\equiv \text{TV}(p_X, p_Y) = \frac{1}{2} \int |p_X(x) - p_Y(x)| dx. \end{aligned}$$

For any two functions $f(T)$, $g(T) > 0$, we write $f(T) \lesssim g(T)$ or $f(T) = O(g(T))$ to indicate $f(T) \leq Cg(T)$ for some absolute constant $C > 0$. We say $f(T) \asymp g(T)$ when $Cf(T) \leq g(T) \leq C'f(T)$ for some absolute constants $C' > C > 0$. The notation $\tilde{O}(\cdot)$ and $\tilde{\Omega}(\cdot)$ represent the respective bounds up to logarithmic factors. Finally, we write $f(T) = o(g(T))$ to denote that $\limsup_{T \rightarrow \infty} f(T)/g(T) = 0$.

2 Preliminaries for diffusion models

Given training samples from a target distribution p_{data} on \mathbb{R}^d , diffusion models aim to generate new samples from p_{data} . Recall the forward process [\(2\)](#). If we define

$$\bar{\alpha}_t := \prod_{k=1}^t \alpha_k, \quad t = 1, 2, \dots, T, \tag{4}$$

the forward process can be expressed as a linear combination of the initial distribution and a Gaussian noise

$$X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \bar{W}_t, \quad t = 1, 2, \dots, T, \tag{5}$$

where $\bar{W}_t \sim \mathcal{N}(0, I_d)$ denotes a d -dimensional standard Gaussian random vector independent of X_0 . When $\bar{\alpha}_T$ is sufficiently small, X_T is well-approximated by a standard Gaussian distribution. Taking the continuum limit of (2), the process satisfies the stochastic differential equation (SDE):

$$dx_t = -\frac{1}{2}\beta_t X_t dt + \sqrt{\beta_t} dB_t \quad X_0 \sim p_{\text{data}}; \quad t \in [0, T] \quad (6)$$

for some function $\beta_t : [0, T] \rightarrow \mathbb{R}$, where $(B_t)_{t \in [0, T]}$ is a standard Brownian motion in \mathbb{R}^d .

Diffusion models seek to reverse the above process by iteratively denoising noisy samples generated from $\mathcal{N}(0, I_d)$, reconstructing data samples from p_{data} . From a continuous perspective, given a solution $(X_t)_{t \in [0, T]}$ to (6), classical SDE theory (Anderson, 1982; Haussmann and Pardoux, 1986) ensures that its time reversal $Y_t^{\text{SDE}} := X_{T-t}$ satisfies:

$$dY_t^{\text{SDE}} = \frac{1}{2}\beta_{T-t} \left(Y_t^{\text{SDE}} + 2\nabla \log p_{X_{T-t}}(Y_t^{\text{SDE}}) \right) dt + \sqrt{\beta_{T-t}} dB_t, \quad Y_0^{\text{SDE}} \sim p_{X_T}; \quad t \in [0, T]. \quad (7)$$

Here, p_{X_t} denotes the marginal distribution of X_t in the forward SDE (6).

Score learning/matching. It is clear from the continuous perspective, that the score function $s_t^*(x) := \nabla \log p_{X_t}(x)$ plays an important role in characterizing the reverse process. In fact, if $s_t^*(x)$ were known exactly, the reverse process would be uniquely identified. In practice, however, score functions must be learned from training samples. A natural approach is to estimate $s_t^*(x)$ within a pre-selected function class \mathcal{F} by minimizing the expected squared error:

$$\min_{s_t \in \mathcal{F}} \mathbb{E}_{X \sim p_{X_t}} \left[\|s_t(X) - \nabla \log p_{X_t}(X)\|^2 \right].$$

For Gaussian distributions, integration by parts allows reformulating this objective as (e.g., Hyvärinen (2005); Vincent (2011); Chen et al. (2022))

$$\min_{s_t: \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{W \sim \mathcal{N}(0, I_d), X_0 \sim p_{\text{data}}} \left[\left\| s_t(X_t) + \frac{1}{\sqrt{1 - \bar{\alpha}_t}} W \right\|_2^2 \right]. \quad (8)$$

Here, given the observed $X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} W$, one seeks to predict the independent noise W , a strategy known as score matching. This formulation is particularly useful for practical training since it does not require explicit knowledge of the score function $\nabla \log p_{X_t}$. Instead, it can be approximated using finite samples, making it more feasible for learning the score function from data.

The DDPM sampling procedure. To implement the sampling process, we must discretize the continuous dynamics and obtain score estimates at discrete time steps. Suppose that one obtains score estimates $\{s_t\}$ at $t = 1, \dots, T$. Equipped with these score estimates, the renowned DDPM algorithm Ho et al. (2020) is a stochastic sampler that recursively generates samples using the following update rule. Starting from $Y_T \sim \mathcal{N}(0, I_d)$, DDPM computes Y_{t-1} by

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + (1 - \alpha_t) s_t(Y_t) \right) + \sqrt{1 - \alpha_t} Z_t, \quad t = T, \dots, 1. \quad (9)$$

Here, $Z_1, \dots, Z_T \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ is a sequence of i.i.d. Gaussian random vectors in \mathbb{R}^d . In words, at each step, Y_{t-1} is a weighted sum of Y_t and its estimated score, plus an independent Gaussian noise.

3 Main results

In this section, we state our main results on the performances of DDPM when applied to GMMs (1) and discuss their consequences. Without loss of generality, we focus on the case where $\sigma = 1$ and therefore the covariance of each component is the identity matrix. Otherwise, our algorithm and analysis framework are readily extended to the general case by either rescaling the data accordingly, or adjusting the learning rates accordingly. We start by introducing some assumptions on the GMMs and the quality of our score estimates.

Assumption 1. We assume that each component of the GMM (1) satisfies

$$\|\mu_k\|_2 \leq T^{c_R}, \quad \forall k \in [K] \quad (10)$$

for some absolute constant $c_R > 0$.

This assumption requires that the mean of each component grows at most polynomially with the iteration number T . Expressing the boundedness condition in terms of T allows for cleaner and more concise convergence guarantees. Given that the constant c_R can be chosen arbitrarily large, this assumption allows each component to have exceedingly large mean value. Therefore, it holds true for most distributions that are encountered in practice.

Next, we assess the quality of score estimates by their averaged ℓ_2 accuracy. This form of estimation error matches naturally with training procedures such as the score matching mentioned above.

Assumption 2. We assume the score estimates $\{s_t\}_{t \in [T]}$ satisfy

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X_t \sim p_{X_t}} \left[\|s_t(X_t) - s_t^*(X_t)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2. \quad (11)$$

Notably, this assumption requires the mean squared estimation error averaged over time steps is bounded, rather than the error at any individual step. It is commonly assumed in the literature of diffusion models (e.g., Chen et al. (2022); Benton et al. (2024); Li and Yan (2024b)).

Convergence theory for DDPM. Before stating our main result, we introduce the learning rate schedule $\{\alpha_t\}_{t \in [T]}$. As adopted in previous works on diffusion models (e.g. Li and Cai (2024)), the learning rate sequence is defined iteratively using the cumulative products $\bar{\alpha}_t = \prod_{k=1}^t \alpha_k$ in Eq. (4). More specifically, define

$$\bar{\alpha}_T = \frac{1}{T^{c_0}}, \quad \text{and} \quad \bar{\alpha}_{t-1} = \bar{\alpha}_t + c_1 \frac{\log T}{T} \bar{\alpha}_t (1 - \bar{\alpha}_t), \quad t = T, \dots, 2, \quad (12)$$

where $c_0, c_1 > 0$ are absolute constants satisfying c_0, c_1 are sufficiently large and $c_1/c_0 > 4$. As shown in Lemma 1, this choice of the learning rates yields that property that

$$1 - \alpha_t \lesssim \frac{\log T}{T} \text{ for } t \geq 2, \quad \text{and} \quad 1 - \alpha_1 \leq T^{-c_1/4}.$$

With these assumptions and preparations, we are positioned to state our main result below. The proof of this result is provided in Section 4, with the proofs of auxiliary lemmas postponed to Section A.

Theorem 1. *Under Assumptions 1 and 2, the output Y_0 of the DDPM sampler (9) with the learning rate selected according to (12) satisfies*

$$\text{TV}(X_0, Y_0) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}. \quad (13)$$

In a nutshell, Theorem 1 guarantees that the sampling quality of DDPM, measured in TV distance, is governed by two components: the first accounts for the time discretization error arising from approximating the continuous SDE in Eq. (7) with a discrete procedure; the second component results from the score estimation error. As a result, given access to perfect score estimates, it only takes DDPM no larger than

$$\tilde{O}\left(\frac{1}{\varepsilon}\right),$$

number of iterations to yield a sampler that is ε -close to the target distribution in terms of TV distance. Notably, this iteration complexity is independent of both the ambient dimension d and the number of components K up to some logarithmic factors. In addition, our result is robust to score estimation error: the TV distance between our output distribution and the target distribution scales proportionally to the

$\varepsilon_{\text{score}}$, modulo logarithmic factor. Our result contrasts the common belief that diffusion models inherently require complexity scaling at least with d .

Theorem 1 is established with respect to the TV distance between X_0 and Y_0 , whereas, most theoretical results in diffusion models fail to directly handle the TV distance due to technical reasons. More specifically, most prior works consider the KL divergence which is a natural choice if Girsanov’s theorem is invoked to handle the discrepancy between the forward process and the process when imperfect score functions are concerned. Noteworthily, a recent line of literature (e.g. Li et al. (2023); Li and Yan (2024b); Li et al. (2024b)) enriches the toolbox of analyzing diffusion models by providing a framework of directly working with the TV distance.

To provide some intuition about why one should expect an iteration number independent of both d and K up to logarithmic factors, consider the Jacobian matrix $J_t(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ of the score function $s_t^*(x)$ with $J_t(x) = \frac{\partial s_t^*}{\partial x}$. By direct computation (as in Eq. (20)), it satisfies

$$J_t(x) = -I_d + \bar{\alpha}_t \left\{ \sum_{k=1}^K \pi_k^{(t)}(x) \mu_k \mu_k^\top - \left(\sum_{k=1}^K \pi_k^{(t)}(x) \mu_k \right) \left(\sum_{k=1}^K \pi_k^{(t)}(x) \mu_k \right)^\top \right\}, \quad \forall x \in \mathbb{R}^d. \quad (14)$$

Here, $\pi_k^{(t)}(x)$ denotes the probability of x lying in the cluster k at time t of the forward process

$$\pi_k^{(t)}(x) := \frac{\pi_k \exp\left(-\frac{1}{2}\|x - \sqrt{\bar{\alpha}_t} \mu_k\|_2^2\right)}{\sum_{i=1}^K \pi_i \exp\left(-\frac{1}{2}\|x - \sqrt{\bar{\alpha}_t} \mu_i\|_2^2\right)}, \quad \forall k \in [K], t \in [T]. \quad (15)$$

We prove in Lemma 5 that with high probability

$$\text{tr}(I_d + J_t(X_t)) \leq C_1 \log(KT), \quad (16)$$

for some absolute constant C_1 independent of the problem parameters. This relation, which does not generally hold, is the key to our dimension-free iteration complexity for GMMs. Since GMMs are widely used to approximate general distributions, if a given distribution can be well-approximated by a Gaussian mixture satisfying Eq. (16) also holds true, it is reasonable to expect a dimension-independent iteration complexity for that broader class of distributions as well.

Comparisons to prior literature. Alongside the seminar work (Chen et al., 2022) and the follow-up works (e.g. Lee et al. (2023); Chen et al. (2023a); Benton et al. (2024)), Theorem 1 investigates the algorithmic aspect of learning GMMs assuming efficient score estimation/matching. Among existing works, the most closely related to ours is Liang et al. (2024b), which established an iteration complexity of $\tilde{O}(d/\varepsilon^2)$ by analyzing the Lipschitz properties and second moments of GMMs. Compared to Li and Yan (2024b), which studied DDPM for general distributions under mild assumptions, and attained an $\tilde{O}(d/\varepsilon)$ iteration complexity, Liang et al. (2024b) did not demonstrate any adaptation of DDPM to GMMs. Our result, however, highlights the surprising adaptive property of diffusion models in this setting.

Beyond the algorithmic aspect, Chen et al. (2024); Gatmiry et al. (2024) developed an end-to-end theory by leveraging piecewise polynomial regression for score estimation and integrating it with existing convergence results on DDPM. The runtime and sample complexity of the resulting algorithms scale quasi-polynomially with K/ε or $\log(K/\varepsilon)$ depending on the covariance assumptions. Notably, the number of diffusion steps used in these two works still scales linearly with d . Our result serves as a complementary contribution to Chen et al. (2024); Gatmiry et al. (2024) by isolating the component of the iteration complexity that is independent of both d and K , up to a logarithmic factor.

4 Analysis

In this section, we describe our proof strategies for deriving Theorem 1. The proofs for auxiliary lemmas and facts are deferred to the appendix.

4.1 Preliminaries

Before proceeding with the main analysis, let us collect several key properties that shall be used in our later analysis.

To begin with, Lemma 1 below characterizes the behavior of the learning rates $(\alpha_t)_{t \in [T]}$ chosen in (12). The proof of this result can be found in Li and Cai (2024, Appendix B.1).

Lemma 1. *The learning rates $(\alpha_t)_{t \in [T]}$ specified in (12) satisfy that*

$$1 - \alpha_t \leq c_1 \frac{\log T}{T}, \quad t = 2, \dots, T, \quad (17a)$$

$$1 - \alpha_1 \leq \frac{1}{T^{c_1/4}}, \quad (17b)$$

where c_1 is defined in (12).

Next, in light of Assumption 1 on the GMM we considered, and the forward process (2), it is straightforward to verify each X_t is another mixture of Gaussian distribution with

$$X_t \sim \sum_{k=1}^K \pi_k \mathcal{N}(\sqrt{\alpha_t} \mu_k, I_d),$$

with the density function given by

$$p_{X_t}(x) = \sum_{k=1}^K \pi_k \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|x - \sqrt{\alpha_t} \mu_k\|_2^2\right). \quad (18)$$

Through direct computation, we can derive that its score function takes the following explicit form:

$$s_t^*(x) := \nabla \log p_{X_t}(x) = - \sum_{k=1}^K \pi_k^{(t)}(x) (x - \sqrt{\alpha_t} \mu_k) = -x + \sum_{k=1}^K \pi_k^{(t)}(x) \sqrt{\alpha_t} \mu_k, \quad (19)$$

where we recall in Eq. (15) that $\pi_k^{(t)}(x) : \mathbb{R}^d \rightarrow [0, 1]$ is defined as

$$\pi_k^{(t)}(x) := \frac{\pi_k \exp\left(-\frac{1}{2} \|x - \sqrt{\alpha_t} \mu_k\|_2^2\right)}{\sum_{i=1}^K \pi_i \exp\left(-\frac{1}{2} \|x - \sqrt{\alpha_t} \mu_i\|_2^2\right)}, \quad \forall k \in [K], t \in [T].$$

In addition, the Jacobian matrix $J_t(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ of $s_t^*(x)$ can be computed as

$$\begin{aligned} J_t(x) &:= \frac{\partial s_t^*(x)}{\partial x} = -I_d + \bar{\alpha}_t \sum_{k=1}^K \pi_k^{(t)} \left(\mu_k - \sum_{i=1}^K \pi_i^{(t)} \mu_i \right) \left(\mu_k - \sum_{i=1}^K \pi_i^{(t)} \mu_i \right)^\top \\ &= -I_d + \bar{\alpha}_t \left\{ \sum_{k=1}^K \pi_k^{(t)} \mu_k \mu_k^\top - \left(\sum_{k=1}^K \pi_k^{(t)} \mu_k \right) \left(\sum_{k=1}^K \pi_k^{(t)} \mu_k \right)^\top \right\}, \quad \forall x \in \mathbb{R}^d. \end{aligned} \quad (20)$$

As a remark, we note that $I_d + J_t(x) \succeq 0$ for any $t \in [T]$ and $x \in \mathbb{R}^d$.

Next, we introduce the event \mathcal{E}_t for each $t \in [T]$ as follows:

$$\begin{aligned} \mathcal{E}_t &:= \left\{ x \in \mathbb{R}^d : \text{tr}(I_d + J_t(x)) \leq C_1 \log(KT) \quad \text{and} \right. \\ &\quad \left. \sum_{k=1}^K \pi_k^{(t)} \exp\left(-\zeta_k^{(t)}(x)\right) \leq \exp\left(C_2(1 - \alpha_t)^2 \log^2(KT)\right) \right\}, \end{aligned} \quad (21)$$

for some absolute constants $C_1, C_2 > 0$, where we define $\zeta_k^{(t)}(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ for each $k \in [K]$:

$$\zeta_k^{(t)}(x) := \frac{1 - \alpha_t^2}{2\alpha_t^2} \left(\|x - \sqrt{\alpha_t} \mu_k\|_2^2 - \sum_{i=1}^K \pi_i^{(t)} \|x - \sqrt{\alpha_t} \mu_i\|_2^2 \right) + \frac{1 - \alpha_t}{\alpha_t^2} s_t^*(x)^\top \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} (\mu_i - \mu_k). \quad (22)$$

Finally, we extend the d -dimensional Euclidean space \mathbb{R}^d by adding a single point ∞ , to obtain $\mathbb{R}^d \cup \{\infty\}$. Intuitively, this set $\{\infty\}$ serves as a convenient way to capture all atypical points in the reverse process.

4.2 Step 1: Constructing auxiliary processes

To facilitate our main analysis, we introduce several auxiliary processes below. These processes are constructed only for analysis purpose and are not used in our sampling algorithm.

Sequence $(Y_t^*)_{t=0}^T$ using true scores. We begin by constructing an auxiliary reverse process $(Y_t^*)_{t=0}^T$ using the true score functions:

$$Y_T^* \sim \mathcal{N}(0, I_d), \quad Y_{t-1}^* := \frac{1}{\sqrt{\alpha_t}}(Y_t^* + (1 - \alpha_t)s_t^*(Y_t^*)) + \sqrt{1 - \alpha_t}Z_t; \quad \forall t = T, \dots, 1. \quad (23)$$

where $Z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ is a sequence of i.i.d. standard d -dimensional Gaussian random vectors independent of Y_t^* .

Sequence $(\bar{Y}_t^-, \bar{Y}_t)_{t=0}^T$. Next, we introduce two auxiliary sequences $(\bar{Y}_t^-)_{t=0}^T$ and $(\bar{Y}_t)_{t=0}^T$ to capture the discretization error modulo some low probability event. These sequences, together with Y_T implemented practically, form a Markov chain with the following transition structure:

$$Y_T \rightarrow \bar{Y}_T^- \rightarrow \bar{Y}_T \rightarrow \bar{Y}_{T-1}^- \rightarrow \bar{Y}_{T-1} \rightarrow \dots \rightarrow \bar{Y}_1^- \rightarrow \bar{Y}_1 \rightarrow \bar{Y}_0^- \rightarrow \bar{Y}_0. \quad (24)$$

- *Initialization.* For $t = T$, we define

$$\bar{Y}_T^- := \begin{cases} Y_T, & \text{if } Y_T \in \mathcal{E}_T, \\ \infty, & \text{otherwise.} \end{cases} \quad (25a)$$

The density of \bar{Y}_T^- satisfies

$$p_{\bar{Y}_T^-}(y) = p_{Y_T}(y)\mathbb{1}\{y \in \mathcal{E}_T\} + \mathbb{P}\{Y_T \notin \mathcal{E}_T\}\delta_\infty(y). \quad (25b)$$

- *Transition from \bar{Y}_t^- to \bar{Y}_t .* For $t = T, \dots, 0$, we define \bar{Y}_t as follows: conditional on $\bar{Y}_t^- = y$,

$$\bar{Y}_t := \begin{cases} y, & \text{with prob. } p_{X_t}(y)/p_{\bar{Y}_t^-}(y) \wedge 1, \\ \infty, & \text{with prob. } 1 - \{p_{X_t}(y)/p_{\bar{Y}_t^-}(y) \wedge 1\}. \end{cases} \quad (26a)$$

The conditional density of \bar{Y}_t given $\bar{Y}_t^- = y$ obeys

$$p_{\bar{Y}_t|\bar{Y}_t^-}(x|y) = \{p_{X_t}(y)/p_{\bar{Y}_t^-}(y) \wedge 1\}\delta_y(x) + (1 - \{p_{X_t}(y)/p_{\bar{Y}_t^-}(y) \wedge 1\})\delta_\infty(x). \quad (26b)$$

We make a critical implication of the above construction: for any $t \geq 0$, the density of \bar{Y}_t satisfies

$$p_{\bar{Y}_t}(y) = \{p_{X_t}(y)/p_{\bar{Y}_t^-}(y) \wedge 1\}p_{\bar{Y}_t^-}(y) = p_{X_t}(y) \wedge p_{\bar{Y}_t^-}(y), \quad \forall y \in \mathbb{R}^d. \quad (27)$$

- *Transition from \bar{Y}_t to \bar{Y}_{t-1}^- .* For each $t = T, \dots, 1$, we first draw a candidate sample

$$\tilde{Y}_{t-1} := \frac{1}{\sqrt{\alpha_t}}(\bar{Y}_t + (1 - \alpha_t)s_t^*(\bar{Y}_t)) + \sqrt{1 - \alpha_t}W_t, \quad (28a)$$

where $W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$, $t \geq 1$ is a sequence of i.i.d. standard Gaussian random vectors independent of $(Z_t)_{t=1}^T$, and then define

$$\bar{Y}_{t-1}^- := \begin{cases} \tilde{Y}_{t-1}, & \text{if } \bar{Y}_t \in \mathcal{E}_t \text{ and } \tilde{Y}_{t-1} \in \mathcal{E}_t, \\ \infty, & \text{otherwise.} \end{cases} \quad (28b)$$

The conditional density of \bar{Y}_{t-1}^- given $\bar{Y}_t = y$ satisfies: if $y \in \mathcal{E}_t$, then

$$p_{\bar{Y}_{t-1}^-|\bar{Y}_t}(x|y) = p_{Y_{t-1}^*|Y_t^*}(x|y)\mathbb{1}\{x \in \mathcal{E}_t\} + \mathbb{P}\{Y_{t-1}^* \notin \mathcal{E}_t | Y_t^* = y\}\delta_\infty(x); \quad (28c)$$

otherwise,

$$p_{\bar{Y}_{t-1}^-|\bar{Y}_t}(x|y) = \delta_\infty(x). \quad (28d)$$

Sequence $(\widehat{Y}_t^-, \widehat{Y}_t)_{t=0}^T$. Finally, we introduce two additional auxiliary sequences $(\widehat{Y}_t^-)_{t=0}^T$ and $(\widehat{Y}_t)_{t=0}^T$, which forms the following Markov chain together with Y_T :

$$Y_T \rightarrow \widehat{Y}_T^- \rightarrow \widehat{Y}_T \rightarrow \widehat{Y}_{T-1}^- \rightarrow \widehat{Y}_{T-1} \rightarrow \cdots \rightarrow \widehat{Y}_1^- \rightarrow \widehat{Y}_1 \rightarrow \widehat{Y}_0^- \rightarrow \widehat{Y}_0. \quad (29)$$

- *Initialization.* For $t = T$, we initialize $\widehat{Y}_T^- = \overline{Y}_T^-$.
- *Transition from \widehat{Y}_t^- to \widehat{Y}_t .* For $t = T, \dots, 0$, the conditional density of \widehat{Y}_t given $\widehat{Y}_t^- = y$ obeys

$$p_{\widehat{Y}_t | \widehat{Y}_t^-}(x | y) = p_{\overline{Y}_t | \overline{Y}_t^-}(x | y). \quad (30)$$

- *Transition from \widehat{Y}_t to \widehat{Y}_{t-1}^- .* For $t = T, \dots, 1$, the conditional density of \overline{Y}_{t-1}^- given $\overline{Y}_t = y$ satisfies: if $y \in \mathcal{E}_t$, then

$$p_{\widehat{Y}_{t-1}^- | \widehat{Y}_t}(x | y) = p_{Y_{t-1} | Y_t}(x | y) \mathbf{1}\{x \in \mathcal{E}_t\} + \mathbb{P}\{Y_{t-1} \notin \mathcal{E}_t | Y_t = y\} \delta_\infty(x); \quad (31a)$$

otherwise,

$$p_{\widehat{Y}_{t-1}^- | \widehat{Y}_t}(x | y) = \delta_\infty(x). \quad (31b)$$

The sequences $(\widehat{Y}_t^-)_{t=0}^T$ and $(\widehat{Y}_t)_{t=0}^T$ are constructed following the same principles as $(\overline{Y}_t^-)_{t=0}^T$ and $(\overline{Y}_t)_{t=0}^T$, with one key difference: the transition from \widehat{Y}_t to \widehat{Y}_{t-1}^- is computed using estimated score functions rather than the true score functions.

A crucial property. It is noteworthy that for any $t \geq 0$, the density of \widehat{Y}_t satisfies

$$p_{\widehat{Y}_t}(x) \leq p_{Y_t}(x), \quad \forall x \in \mathbb{R}^d, \quad (32)$$

and consequently, $p_{\widehat{Y}_t}(x) \geq p_{Y_t}(x)$ for $x = \infty$. To see this, we first note that the base case $t = T$ holds since $\widehat{Y}_T \stackrel{d}{=} Y_T$, which arises from $\widehat{Y}_T^- = \overline{Y}_T^-$ and $p_{\widehat{Y}_T | \widehat{Y}_T^-} = p_{\overline{Y}_T | \overline{Y}_T^-}$ by (30). Next, suppose that (32) holds for $t + 1$. Then for any $x \in \mathbb{R}^d$, one has

$$\begin{aligned} p_{\widehat{Y}_t}(x) &\stackrel{(i)}{=} \{p_{X_t}(x)/p_{\overline{Y}_t^-}(x) \wedge 1\} p_{\widehat{Y}_t^-}(x) \leq p_{\widehat{Y}_t^-}(x) = \int_{\mathbb{R}^d} p_{\widehat{Y}_t^- | \widehat{Y}_{t+1}}(x | y) p_{\widehat{Y}_{t+1}}(y) dy \\ &\stackrel{(ii)}{\leq} \int_{\mathbb{R}^d} p_{Y_t | Y_{t+1}}(x | y) p_{Y_{t+1}}(y) dy = p_{Y_t}(x), \end{aligned}$$

where (i) uses (30) and (26b); (ii) is true due to the induction hypothesis and (31a).

Error decomposition. In view of triangle's inequality, we can upper bound the TV distance between p_{X_0} and p_{Y_0} by two terms

$$\text{TV}(p_{X_0}, p_{Y_0}) \leq \text{TV}(p_{X_0}, p_{\overline{Y}_0}) + \text{TV}(p_{\overline{Y}_0}, p_{Y_0}), \quad (33)$$

where the first term acts in the role of discretization error modulo some low probability event, as \overline{Y}_t is defined using the true scores, whereas the second term captures the error caused by imperfect score estimation. In the sequel, we control each term separately.

4.3 Step 2: Bounding discretization error

In this section, we proceed to bound $\text{TV}(p_{X_0}, p_{\overline{Y}_0})$. Let us first define function $\Delta_t(x) : \mathbb{R}^d \rightarrow \mathbb{R}$, where for each $t = 0, \dots, T$:

$$\Delta_t(x) := p_{X_t}(x) - p_{\overline{Y}_t}(x), \quad \forall x \in \mathbb{R}^d. \quad (34)$$

In view of relation (27), one has $\Delta_t(x) \geq 0$ for all $t \geq 0$ and $x \in \mathbb{R}^d$. Applying the formula for the total variation $\text{TV}(p, q) = \int_{x: p(x) > q(x)} (p(x) - q(x)) dx$, we find that

$$\text{TV}(p_{X_0}, p_{\bar{Y}_0}) = \int_{\mathbb{R}^d \cup \{\infty\}} (p_{X_0}(x) - p_{\bar{Y}_0}(x)) \mathbb{1}\{p_{X_0}(x) > p_{\bar{Y}_0}(x)\} dx = \int_{\mathbb{R}^d} \Delta_0(x) dx. \quad (35)$$

Thus, it is sufficient to bound $\int \Delta_0(x) dx$, which shall be done using an inductive argument. We start with the base case, which is characterized by the Lemma 2 below.

Lemma 2. *It satisfies that*

$$\int_{\mathbb{R}^d} \Delta_T(x) dx \lesssim T^{-3}. \quad (36)$$

Proof. See Appendix A.1. □

In addition, Lemma 3 below establishes the inductive relationship between t and $t - 1$.

Lemma 3. *For all $t = T, \dots, 1$, one has*

$$\int_{\mathbb{R}^d} \Delta_{t-1}(x) dx - \int_{\mathbb{R}^d} \Delta_t(x) dx \lesssim (1 - \alpha_t)^2 \log^2(KT) + T^{-3}, \quad (37)$$

Proof. See Appendix A.2. □

Consequently, combining (36)–(37) with (35) leads to

$$\begin{aligned} \text{TV}(p_{X_0}, p_{\bar{Y}_0}) &= \int_{\mathbb{R}^d} \Delta_0(x) dx \leq \int_{\mathbb{R}^d} \Delta_T(x) dx + T \cdot O((1 - \alpha_t)^2 \log^2(KT)) + T \cdot O(T^{-3}) \\ &\lesssim \frac{1}{T^3} + \frac{\log^2(KT) \log^2 T}{T} + \frac{1}{T^2} \\ &\asymp \frac{\log^2(KT) \log^2 T}{T}, \end{aligned} \quad (38)$$

where the penultimate step uses $1 - \alpha_t \lesssim \log T/T$ by (17).

4.4 Step 3: Relating to score estimation error

Next, we control the term $\text{TV}(p_{\bar{Y}_0}, p_{Y_0})$. First, in view of basic calculations, we can write

$$\begin{aligned} \text{TV}(p_{\bar{Y}_0}, p_{Y_0}) &= \int_{\mathbb{R}^d} (p_{\bar{Y}_0}(x) - p_{Y_0}(x)) \mathbb{1}\{p_{\bar{Y}_0}(x) > p_{Y_0}(x)\} dx + \mathbb{P}\{\bar{Y}_0 = \infty\} \\ &\stackrel{(i)}{\leq} \int_{\mathbb{R}^d} (p_{\bar{Y}_0}(x) - p_{\hat{Y}_0}(x)) \mathbb{1}\{p_{\bar{Y}_0}(x) > p_{\hat{Y}_0}(x)\} dx + \mathbb{P}\{\bar{Y}_0 = \infty\} \\ &\stackrel{(ii)}{\leq} \text{TV}(p_{\bar{Y}_0}, p_{\hat{Y}_0}) + \text{TV}(p_{X_0}, p_{\bar{Y}_0}) \\ &\stackrel{(iii)}{\leq} \sqrt{\text{KL}(p_{\bar{Y}_0} \parallel p_{\hat{Y}_0})} + O\left(\frac{\log^2(KT) \log^2 T}{T}\right). \end{aligned} \quad (39)$$

where (i) arises from (32) that $p_{Y_0}(x) \geq p_{\hat{Y}_0}(x)$ for any $x \in \mathbb{R}^d$; (ii) uses $\mathbb{P}\{\bar{Y}_0 = \infty\} \leq \text{TV}(p_{X_0}, p_{\bar{Y}_0})$ since $X_0 \in \mathbb{R}^d$; (iii) applies Pinsker's inequality and (38).

To further control the right hand side of (39), it suffices to bound $\text{KL}(p_{\bar{Y}_0} \parallel p_{\hat{Y}_0})$ from above. Towards this, notice that

$$\begin{aligned} \text{KL}(p_{\bar{Y}_0} \parallel p_{\hat{Y}_0}) &\stackrel{(i)}{\leq} \text{KL}(p_{\bar{Y}_T^-, \bar{Y}_T, \dots, \bar{Y}_0^-, \bar{Y}_0} \parallel p_{\hat{Y}_T^-, \hat{Y}_T, \dots, \hat{Y}_0^-, \hat{Y}_0}) \\ &\stackrel{(ii)}{=} \text{KL}(p_{\bar{Y}_T^-} \parallel p_{\hat{Y}_T^-}) + \sum_{t=0}^{T-1} \mathbb{E}_{x_t \sim p_{\bar{Y}_t^-}} \left[\text{KL}(p_{\bar{Y}_t | \bar{Y}_t^- = x_t} \parallel p_{\hat{Y}_t | \hat{Y}_t^- = x_t}) \right] \end{aligned}$$

$$\begin{aligned}
& + \sum_{t=1}^T \mathbb{E}_{x_t \sim p_{\bar{Y}_t}} \left[\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) \right] \\
& \stackrel{\text{(iii)}}{=} \sum_{t=1}^T \mathbb{E}_{x_t \sim p_{\bar{Y}_t}} \left[\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) \right].
\end{aligned} \tag{40}$$

Here, (i) applies the data-processing inequality; (ii) uses the chain rule of KL divergence and the Markov property; (iii) is true since we initialize $\hat{Y}_T^- = \bar{Y}_T^-$ and the transition kernels from \hat{Y}_t^- to \hat{Y}_{t+1}^- are the same as those from \bar{Y}_t^- to \bar{Y}_{t+1}^- for all $t \geq 0$.

Note that $Y_{t-1}^* | Y_t^* \sim \mathcal{N}(\frac{1}{\sqrt{\alpha_t}}(Y_t^* + (1 - \alpha_t)s_t^*(Y_t^*)), (1 - \alpha_t)I_d)$ and $Y_{t-1} | Y_t \sim \mathcal{N}(\frac{1}{\sqrt{\alpha_t}}(Y_t + (1 - \alpha_t)s_t(Y_t)), (1 - \alpha_t)I_d)$. For any $x_t \in \mathcal{E}_t$, write $p(\cdot) = p_{Y_{t-1}^* | Y_t^*}(\cdot)$ and $q(\cdot) = p_{Y_{t-1} | Y_t}(\cdot)$. In view of the definitions (31a) and (26b), one has

$$\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) = \int_{\mathcal{E}_t} \log \frac{p(x)}{q(x)} p(x) dx + \log \frac{\int_{\mathcal{E}_t^c} p(x) dx}{\int_{\mathcal{E}_t^c} q(x) dx} \int_{\mathcal{E}_t^c} p(x) dx. \tag{41}$$

Now invoke Li and Yan (2024b, Lemma 6) to derive

$$\begin{aligned}
\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) & \leq \int_{\mathbb{R}^d} \log \frac{p(x)}{q(x)} p(x) dx \\
& = \text{KL}(p_{Y_{t-1}^* | Y_t^* = x_t} \parallel p_{Y_{t-1} | Y_t = x_t}) \\
& \stackrel{\text{(i)}}{=} \frac{1 - \alpha_t}{2\alpha_t} \|s_t(x_t) - s_t^*(x_t)\|_2^2 \\
& \stackrel{\text{(ii)}}{\lesssim} \frac{\log T}{T} \|s_t(x_t) - s_t^*(x_t)\|_2^2.
\end{aligned} \tag{42}$$

Here, (i) uses the formula of the KL divergence for two normal distributions; (ii) uses $1 - \alpha_t \lesssim \log T/T$ by (17). Meanwhile, for any $x_t \in \mathcal{E}_t^c$, we know from (31b) that

$$\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) = 0. \tag{43}$$

Combined with (40), the above bounds gives

$$\begin{aligned}
\text{KL}(p_{\bar{Y}_0} \parallel p_{\hat{Y}_0}) & \stackrel{\text{(i)}}{\leq} \sum_{t=1}^T \mathbb{E}_{x_t \sim p_{X_t}} \left[\text{KL}(p_{\bar{Y}_{t-1}^- | \bar{Y}_t = x_t} \parallel p_{\hat{Y}_{t-1}^- | \hat{Y}_t = x_t}) \right] \\
& \stackrel{\text{(ii)}}{\lesssim} \frac{\log T}{T} \sum_{t=1}^T \mathbb{E} \left[\|s_t(X_t) - s_t^*(X_t)\|_2^2 \right] \\
& \stackrel{\text{(iii)}}{\lesssim} \varepsilon_{\text{score}}^2 \log T,
\end{aligned} \tag{44}$$

where (i) arises from (43) and (27) that $p_{\bar{Y}_t}(x) \leq p_{X_t}(x)$ for all $x \in \mathbb{R}^d$; (ii) uses (42); (iii) follows from Assumption 2 on the score estimation. Substituting (44) into (39) leads to

$$\text{TV}(p_{Y_0}, p_{\bar{Y}_0}) \lesssim \frac{\log^2(KT) \log^2 T}{T} + \varepsilon_{\text{score}} \sqrt{\log T}. \tag{45}$$

In conclusion, putting relations (38) and (45) together with (33) completes the proof of Theorem 1.

5 Discussion

In summary, this paper explores the effectiveness of diffusion models in learning GMMs and presents new theoretical insights on how generative models implicitly exploit data structure to achieve efficient sampling.

While DDPM requires a number of iterations that scale linearly with data dimension in the worst case, our main result unveils a surprising efficiency of DDPM: it only requires an iteration complexity of $\tilde{O}(1/\varepsilon)$ to learn an ε -accurate distribution, independent of both the data dimension d and the number of mixture components K , up to logarithmic factors. This result suggests that diffusion models can efficiently learn structured distributions even in ultra-high-dimensional settings.

Before concluding, we highlight several promising directions for future investigation. First, since GMMs are widely used to approximate complex distributions, our findings suggest that if a distribution can be well-approximated by a GMM, it may also be efficiently learned using diffusion models. Formalizing this intuition and developing adaptive guarantees for learning general distributions would be a valuable extension. Additionally, while this paper focuses on mixtures of spherical Gaussians, an important next step is to analyze the iteration complexity of DDPM when applied to more general cases, such as mixtures with well-conditioned but arbitrary covariances, as considered in [Chen et al. \(2024\)](#). Finally, our analysis primarily addresses the sampling phase, leaving open the question of how score estimation efficiency is affected by the structure of GMMs. It remains a crucial direction to establish an end-to-end theory that integrates both score learning and sampling and fully unleashes the potential of diffusion models adapting to low-dimensional structure.

Acknowledgement

G. Li is supported in part by the Chinese University of Hong Kong Direct Grant for Research. Y. Wei is supported in part by the NSF grants CCF-2106778, CCF-2418156 and CAREER award DMS-2143215.

A Proof of technical lemmas

A.1 Proof of Lemma 2

Recalling that $\Delta_T(x) := p_{X_T}(x) - p_{\bar{Y}_T}(x) \geq 0$ for any $x \in \mathbb{R}^d$, we can derive

$$\begin{aligned} \int_{\mathbb{R}^d} \Delta_T(x) dx &= \int_{\mathbb{R}^d} (p_{X_T}(x) - p_{\bar{Y}_T}(x)) \mathbb{1}\{p_{X_T}(x) > p_{\bar{Y}_T}(x)\} dx \\ &\stackrel{(i)}{=} \int_{\mathbb{R}^d} (p_{X_T}(x) - p_{\bar{Y}_T}(x)) \mathbb{1}\{p_{X_T}(x) > p_{\bar{Y}_T}(x)\} dx \\ &\stackrel{(ii)}{=} \text{TV}(p_{X_T}, p_{\bar{Y}_T}) \\ &\leq \text{TV}(p_{X_T}, p_{Y_T}) + \text{TV}(p_{Y_T}, p_{\bar{Y}_T}), \end{aligned} \quad (46)$$

where (i) arises from (27) that $p_{\bar{Y}_T}(x) = p_{X_T}(x) \wedge p_{\bar{Y}_T}(x)$ for any $x \in \mathbb{R}^d$, (ii) uses the formula of the total variation $\text{TV}(p, q) = \int_{x: p(x) > q(x)} (p(x) - q(x)) dx$ and $X_T \in \mathbb{R}^d$.

Consequently, it suffices to control the two quantities in (46) respectively.

- For the first term $\text{TV}(p_{X_T}, p_{Y_T})$ corresponding to the initialization error, we can derive

$$\begin{aligned} \text{KL}(p_{X_T} \parallel p_{Y_T}) &\stackrel{(i)}{\leq} \mathbb{E} \left[\text{KL}(p_{X_T}(\cdot \mid X_0) \parallel p_{Y_T}(\cdot)) \right] \\ &\stackrel{(ii)}{=} \frac{1}{2} \mathbb{E} \left[d(1 - \bar{\alpha}_T) - d + \|\sqrt{\bar{\alpha}_T} X_0\|_2^2 - d \log(1 - \bar{\alpha}_T) \right] \\ &\stackrel{(iii)}{\lesssim} T^{-c_0} \mathbb{E}[\|X_0\|_2^2] \stackrel{(iv)}{\lesssim} T^{-c_0} (T^{c_R} + d) \stackrel{(v)}{\lesssim} T^{-8}, \end{aligned} \quad (47)$$

where (i) arises from the convexity of the KL divergence; (ii) applies the KL divergence formula for two normal distributions; (iii) is due to the choice of the learning rate $\bar{\alpha}_T = T^{-c_0} = o(1)$ in (12) and $\log(1 - x) \geq -x$ for any $x \in [0, 1/2]$; (iv) holds due to Assumption 1 that

$$\mathbb{E}[\|X_0\|_2^2] = \sum_{k=1}^K \pi_k (\|\mu_k\|_2^2 + d) \leq T^{c_R} + d;$$

and (v) holds as long as T and c_0 are sufficiently large. It then follows from Pinsker's inequality that

$$\mathrm{TV}(p_{X_T}, p_{Y_T}) \leq \sqrt{\mathrm{KL}(p_{X_T} \parallel p_{Y_T})} \lesssim T^{-4}. \quad (48)$$

- We proceed to consider the second term $\mathrm{TV}(p_{Y_T}, p_{\bar{Y}_T^-})$. By the construction of \bar{Y}_T^- (see (25)), one can write

$$\begin{aligned} \mathrm{TV}(p_{Y_T}, p_{\bar{Y}_T^-}) &\stackrel{(i)}{=} \int_{\mathbb{R}^d} (p_{Y_T}(x) - p_{\bar{Y}_T^-}(x)) \mathbb{1}\{p_{Y_T}(x) > p_{\bar{Y}_T^-}(x)\} dx \\ &\stackrel{(ii)}{=} \int_{\mathcal{E}_T^c} p_{Y_T}(x) dx \\ &\stackrel{(iii)}{\leq} \int_{\mathcal{E}_T^c} p_{X_T}(x) dx + \mathrm{TV}(p_{X_T}, p_{Y_T}). \end{aligned}$$

Here, (i) uses the formula of the total variation $\mathrm{TV}(p, q) = \int_{x: p(x) > q(x)} (p(x) - q(x)) dx$; (ii) follows from $Y_T \in \mathbb{R}^d$ and (25b) that $p_{\bar{Y}_T^-}(x) = p_{Y_T}(x)$ if $x \in \mathcal{E}_T$ and $p_{\bar{Y}_T^-}(x) = 0$ if $x \in \mathbb{R}^d \setminus \mathcal{E}_T$; (iii) arises from the definition of total variation distance that $\mathrm{TV}(p, q) = \sup_B |p(B) - q(B)|$. As we shall see momentarily in Lemma 5 in Appendix A.4, one has

$$\int_{\mathcal{E}_T^c} p_{X_T}(x) \lesssim T^{-3}.$$

Combined with (48), this leads to

$$\mathrm{TV}(p_{Y_T}, p_{\bar{Y}_T^-}) \lesssim T^{-3} + T^{-4} \asymp T^{-3}. \quad (49)$$

In conclusion, substituting (48) and (49) into (46) completes the proof of Lemma 2.

A.2 Proof of Lemma 3

Fix an arbitrary $t \in [T]$. To analyze $\int_{\mathbb{R}^d} \Delta_t(x_t) dx_t$, let us first introduce a function $\Delta_{t \rightarrow t-1}(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ where

$$\Delta_{t \rightarrow t-1}(x) := \int_{x_t \in \mathcal{E}_t} p_{Y_{t-1}^* | Y_t^*}(x | x_t) \Delta_t(x_t) dx_t, \quad \forall x \in \mathbb{R}^d. \quad (50)$$

Note that in view of relation (27), $\Delta_t(x) \geq 0$ for all $x \in \mathbb{R}^d$ and therefore $\Delta_{t \rightarrow t-1}(x) \geq 0$ for all $x \in \mathbb{R}^d$. It is easily seen that

$$\int_{\mathbb{R}^d} \Delta_{t \rightarrow t-1}(x_{t-1}) dx_{t-1} = \int_{x_t \in \mathcal{E}_t} \int_{x_{t-1} \in \mathbb{R}^d} p_{Y_{t-1}^* | Y_t^*}(x_{t-1} | x_t) dx_{t-1} \Delta_t(x_t) dx_t \leq \int_{\mathbb{R}^d} \Delta_t(x_t) dx_t. \quad (51)$$

As a result, to upper bound $\int_{\mathbb{R}^d} \Delta_{t-1}(x) dx - \int_{\mathbb{R}^d} \Delta_t(x) dx$, it is sufficient to consider $\int_{\mathbb{R}^d} \Delta_{t-1}(x) dx - \int_{\mathbb{R}^d} \Delta_{t \rightarrow t-1}(x_{t-1}) dx_{t-1}$.

Towards this, we find it helpful to first make the following observation. For any $x_{t-1} \in \mathbb{R}^d$ such that $\Delta_{t-1}(x_{t-1}) > 0$, or equivalently, $p_{X_{t-1}}(x_{t-1}) > p_{\bar{Y}_{t-1}^-}(x_{t-1})$, we have

$$p_{X_{t-1}}(x_{t-1}) - \Delta_{t-1}(x_{t-1}) + \Delta_{t \rightarrow t-1}(x_{t-1}) = p_{\bar{Y}_{t-1}^-}(x_{t-1}) + \Delta_{t \rightarrow t-1}(x_{t-1}). \quad (52)$$

Here, we use the fact that $p_{\bar{Y}_{t-1}^-}(x_{t-1}) = p_{X_{t-1}}(x_{t-1}) \wedge p_{\bar{Y}_{t-1}^-}(x_{t-1}) = p_{\bar{Y}_{t-1}^-}(x_{t-1})$ since $p_{X_{t-1}}(x_{t-1}) > p_{\bar{Y}_{t-1}^-}(x_{t-1})$. To further control the right hand side, recall the definition of $\Delta_t(x)$ in (34) and the constructed transition kernel of $p_{\bar{Y}_{t-1}^- | \bar{Y}_t}$ in (28c). For any $x_{t-1} \in \mathbb{R}^d$, we arrive at

$$p_{\bar{Y}_{t-1}^-}(x_{t-1}) \geq \int_{x_t \in \mathcal{E}_t} p_{Y_{t-1}^* | Y_t^*}(x_{t-1} | x_t) p_{\bar{Y}_t}(x_t) dx_t$$

$$= \int_{x_t \in \mathcal{E}_t} p_{Y_{t-1}^* | Y_t^*}(x_{t-1} | x_t) p_{X_t}(x_t) dx_t - \Delta_{t \rightarrow t-1}(x_{t-1}). \quad (53)$$

As a result, we obtain

$$\begin{aligned} & p_{X_{t-1}}(x_{t-1}) - \Delta_{t-1}(x_{t-1}) + \Delta_{t \rightarrow t-1}(x_{t-1}) \\ & \geq \int_{x_t \in \mathcal{E}_t} p_{Y_{t-1}^* | Y_t^*}(x_{t-1} | x_t) p_{X_t}(x_t) dx_t \\ & \stackrel{(i)}{=} \int_{\mathcal{E}_t} \left(\frac{1}{2\pi(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t}x_{t-1} - x_t - (1-\alpha_t)s_t^*(x_t)\|_2^2}{2\alpha_t(1-\alpha_t)} \right) p_{X_t}(x_t) dx_t \\ & \stackrel{(ii)}{=} \int_{\mathcal{E}_t} \left(\frac{1}{2\pi(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t}x_{t-1} - u_t\|_2^2}{2\alpha_t(1-\alpha_t)} \right) \det(I_d + (1-\alpha_t)J_t(x_t))^{-1} p_{X_t}(x_t) du_t. \end{aligned} \quad (54)$$

where (i) uses (53); (ii) arises from (23) that $Y_{t-1}^* | Y_t^* \sim \mathcal{N}(\alpha_t^{-1/2}[Y_t^* + (1-\alpha_t)s_t^*(Y_t^*)], (1-\alpha_t)I_d)$; (ii) applies the change of variable

$$u_t := x_t + (1-\alpha_t)s_t^*(x_t).$$

To bound the integral in (54), we present Lemma 4 below.

Lemma 4. *For any $t \in [T]$, the following holds for any $x_t \in \mathcal{E}_t$:*

$$\begin{aligned} & \det(I_d + (1-\alpha_t)J_t(x_t))^{-1} p_{X_t}(x_t) \\ & = \left(\frac{1}{2\pi\alpha_t^2} \right)^{d/2} \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t}\mu_k\|_2^2}{2\alpha_t^2} \right). \end{aligned} \quad (55)$$

Proof. See Appendix A.3. □

Plugging (55) into (54) and invoking the inequality $\exp(x) \geq 1 + x$ leads to

$$\begin{aligned} & p_{X_{t-1}}(x_{t-1}) - \Delta_{t-1}(x_{t-1}) + \Delta_{t \rightarrow t-1}(x_{t-1}) \\ & \geq \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) \\ & \quad \cdot \int_{\mathcal{E}_t} \left(\frac{1}{4\pi^2\alpha_t^2(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t}x_{t-1} - u_t\|_2^2}{2\alpha_t(1-\alpha_t)} \right) \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t}\mu_k\|_2^2}{2\alpha_t^2} \right) du_t. \end{aligned} \quad (56)$$

To further control the right hand side, direct computations give that

$$\begin{aligned} & \int_{\mathbb{R}^d} \left(\frac{1}{4\pi^2\alpha_t^2(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t}x_{t-1} - u_t\|_2^2}{2\alpha_t(1-\alpha_t)} \right) \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t}\mu_k\|_2^2}{2\alpha_t^2} \right) du_t \\ & = \int_{\mathbb{R}^d} \left(\frac{1}{4\pi^2\alpha_t^2(1-\alpha_t)} \right)^{d/2} \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t}(1-\alpha_t)\mu_k - \alpha_t^{3/2}x_{t-1}\|_2^2}{2\alpha_t^2(1-\alpha_t)} - \frac{\|x_{t-1} - \sqrt{\alpha_t}/\alpha_t\mu_k\|_2^2}{2} \right) du_t \\ & \stackrel{(i)}{=} \int_{\mathbb{R}^d} \left(\frac{1}{2\pi} \right)^{d/2} \sum_{k=1}^K \pi_k \exp \left(- \frac{1}{2} \|x_{t-1} - \sqrt{\alpha_{t-1}}\mu_k\|_2^2 \right) du_t \\ & \stackrel{(ii)}{=} p_{X_{t-1}}(x_{t-1}). \end{aligned} \quad (57)$$

Here (i) is true as $u_t \mapsto (2\pi\alpha_t^2(1-\alpha_t))^{-d/2} \exp \left(- (2\alpha_t^2(1-\alpha_t))^{-1} \|u_t - \sqrt{\alpha_t}(1-\alpha_t)\mu_k - \alpha_t^{3/2}x_{t-1}\|_2^2 \right)$ is a density function and $\bar{\alpha}_t := \prod_{i=1}^t \alpha_i$, and (ii) arises from (18). Hence, if we define function $\delta_{t-1}(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ to capture the integral on set \mathcal{E}_t^c where

$$\delta_{t-1}(x) := \int_{\mathcal{E}_t^c} \left(\frac{1}{4\pi^2\alpha_t^2(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t}x - u_t\|_2^2}{2\alpha_t(1-\alpha_t)} \right) \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t}\mu_k\|_2^2}{2\alpha_t^2} \right) du_t,$$

then it obeys

$$\delta_{t-1}(x) = p_{X_{t-1}}(x) - \int_{\mathcal{E}_t} \left(\frac{1}{4\pi^2 \alpha_t^2 (1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t} x - u_t\|_2^2}{2\alpha_t (1-\alpha_t)} \right) \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t} \mu_k\|_2^2}{2\alpha_t^2} \right) du_t.$$

Combining this definition with relation (56), we obtain

$$\begin{aligned} & p_{X_{t-1}}(x_{t-1}) - \Delta_{t-1}(x_{t-1}) + \Delta_{t \rightarrow t-1}(x_{t-1}) \\ & \geq \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) (p_{X_{t-1}}(x_{t-1}) - \delta_{t-1}(x_{t-1})) \\ & \geq p_{X_{t-1}}(x_{t-1}) + O((1-\alpha_t)^2 \log^2(KT)) p_{X_{t-1}}(x_{t-1}) - O(1) \delta_{t-1}(x_{t-1}), \end{aligned}$$

or equivalently,

$$\Delta_{t-1}(x_{t-1}) \leq \Delta_{t \rightarrow t-1}(x_{t-1}) + O((1-\alpha_t)^2 \log^2(KT)) p_{X_{t-1}}(x_{t-1}) + O(1) \delta_{t-1}(x_{t-1}). \quad (58)$$

Here, we use (17) that $(1-\alpha_t)^2 \log^2(KT) \lesssim \log^2(KT) \log^2 T/T^2 = o(1)$ as long as T is large enough.

We claim that $\int_{\mathbb{R}^d} \delta_{t-1}(x) dx$ satisfies

$$\int_{\mathbb{R}^d} \delta_{t-1}(x) dx \lesssim T^{-3} + (1-\alpha_t)^2 \log^2(KT). \quad (59)$$

Therefore, substituting (59) and (56) into (58) and integrating over x_{t-1} yields

$$\begin{aligned} & \int_{\mathbb{R}^d} \Delta_{t-1}(x_{t-1}) dx_{t-1} \leq \int_{\mathbb{R}^d} \Delta_{t \rightarrow t-1}(x_{t-1}) dx_{t-1} + O((1-\alpha_t)^2 \log^2(KT)) \int_{\mathbb{R}^d} p_{X_{t-1}}(x_{t-1}) dx_{t-1} \\ & \quad + O(1) \int_{\mathbb{R}^d} \delta_{t-1}(x_{t-1}) dx_{t-1} \\ & \leq \int_{\mathbb{R}^d} \Delta_t(x_{t-1}) dx_{t-1} + O((1-\alpha_t)^2 \log^2(KT)) + O(T^{-3}), \end{aligned}$$

where the penultimate line uses (51)

This completes the proof of Lemma 3.

Proof of Claim (59). It remains to control $\int_{\mathbb{R}^d} \delta_{t-1}(x) dx$. To this end, the expression above allows us to derive

$$\begin{aligned} & \int_{\mathbb{R}^d} \delta_{t-1}(x_{t-1}) dx_{t-1} \\ & = 1 - \int_{\mathbb{R}^d} \int_{\mathcal{E}_t} \left(\frac{1}{4\pi^2 \alpha_t^2 (1-\alpha_t)} \right)^{d/2} \sum_{k=1}^K \pi_k \exp \left(- \frac{\|u_t - \sqrt{\alpha_t} \mu_k\|_2^2}{2\alpha_t^2} \right) \exp \left(- \frac{\|\sqrt{\alpha_t} x_{t-1} - u_t\|_2^2}{2\alpha_t (1-\alpha_t)} \right) du_t dx_{t-1} \\ & \stackrel{(i)}{=} 1 - \int_{\mathbb{R}^d} \int_{\mathcal{E}_t} \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) p_{X_t}(x_t) \left(\frac{1}{2\pi(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|\sqrt{\alpha_t} x_{t-1} - u_t\|_2^2}{2\alpha_t (1-\alpha_t)} \right) dx_t dx_{t-1} \\ & = 1 - \int_{\mathcal{E}_t} \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) p_{X_t}(x_t) \int_{\mathbb{R}^d} \left(\frac{1}{2\pi(1-\alpha_t)} \right)^{d/2} \exp \left(- \frac{\|x_{t-1} - u_t/\sqrt{\alpha_t}\|_2^2}{2(1-\alpha_t)} \right) dx_{t-1} dx_t \\ & \stackrel{(ii)}{=} 1 - \int_{\mathcal{E}_t} \exp \left(O((1-\alpha_t)^2 \log^2(KT)) \right) p_{X_t}(x_t) dx_t \\ & \stackrel{(iii)}{\leq} 1 - \int_{\mathcal{E}_t} p_{X_t}(x_t) dx_t + O((1-\alpha_t)^2 \log^2(KT)) \int_{\mathcal{E}_t} p_{X_t}(x_t) dx_t \\ & \leq \int_{\mathcal{E}_t^c} p_{X_t}(x_t) dx_t + O((1-\alpha_t)^2 \log^2(KT)). \quad (60) \end{aligned}$$

Here, (i) invokes Lemma 4, (ii) is true as $x_{t-1} \mapsto (2\pi(1-\alpha_t))^{-d/2} \exp(- (2(1-\alpha_t))^{-1} \|x_{t-1} - u_t/\sqrt{\alpha_t}\|_2^2)$ is a density function, and (iii) holds as $\exp(x) \geq 1+x$ for all $x \in \mathbb{R}^d$.

Finally, the right-hand-side of the above bound is controlled by Lemma 5 below.

Lemma 5. Recall the definition of \mathcal{E}_t in (21). For any $t \in [T]$, one has

$$\int_{\mathcal{E}_t^c} p_{X_t}(x_t) dx_t \lesssim T^{-3}. \quad (61)$$

Proof. See Appendix A.4. \square

Putting everything together completes the proof of Claim (59).

A.3 Proof of Lemma 4

Let us first derive two relations that are key for this proof. To start with, fix an arbitrary $x_t \in \mathcal{E}_t$. Recalling the definition that $u_t := x_t + (1 - \alpha_t)s_t^*(x_t)$, direct calculations yield

$$\begin{aligned} & \frac{1}{2\alpha_t^2} \|u_t - \sqrt{\alpha_t}\mu_k\|_2^2 \\ &= \frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{(1 - \alpha_t)}{\alpha_t^2} s_t^*(x_t)^\top (x_t - \sqrt{\alpha_t}\mu_k) + \frac{(1 - \alpha_t)^2}{2\alpha_t^2} \|s_t^*(x_t)\|_2^2 \\ &= \frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \sum_{i=1}^K \pi_i^{(t)} \|x_t - \sqrt{\alpha_t}\mu_i\|_2^2 + \left(\frac{(1 - \alpha_t)^2}{2\alpha_t^2} - \frac{1 - \alpha_t}{\alpha_t^2} \right) \|s_t^*(x_t)\|_2^2 \\ & \quad + \frac{1 - \alpha_t^2}{2\alpha_t^2} \left(\|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 - \sum_{i=1}^K \pi_i^{(t)} \|x_t - \sqrt{\alpha_t}\mu_i\|_2^2 \right) + \frac{1 - \alpha_t}{\alpha_t^2} s_t^*(x_t)^\top (s_t^*(x_t) + x_t - \sqrt{\alpha_t}\mu_k) \\ &\stackrel{(i)}{=} \frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \sum_{i=1}^K \pi_i^{(t)} \|x_t - \sqrt{\alpha_t}\mu_i\|_2^2 - \frac{1 - \alpha_t^2}{2\alpha_t^2} \|s_t^*(x_t)\|_2^2 + \zeta_k^{(t)}(x_t) \\ &\stackrel{(ii)}{=} \frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \left(\sum_{i=1}^K \pi_i^{(t)} \|x_t - \sqrt{\alpha_t}\mu_i\|_2^2 - \left\| \sum_{i=1}^K \pi_i^{(t)} (x_t - \sqrt{\alpha_t}\mu_i) \right\|_2^2 \right) + \zeta_k^{(t)}(x_t) \\ &\stackrel{(iii)}{=} \frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2 + \frac{(1 - \alpha_t^2)}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) + \zeta_k^{(t)}(x_t). \end{aligned} \quad (62)$$

Here, (i) uses the definition of $\zeta_k^{(t)}(x)$ in (22); (ii) arises from the expression of $s_t^*(x) = -\sum_{k=1}^K \pi_k^{(t)}(x - \sqrt{\alpha_t}\mu_k)$ in (19); (iii) uses the expression of $J_t(x)$ in (14) that

$$\begin{aligned} I_d + J_t(x) &= \sum_{k=1}^K \pi_k^{(t)} \left(\sqrt{\alpha_t}\mu_k - \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t}\mu_i \right) \left(\sqrt{\alpha_t}\mu_k - \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t}\mu_i \right)^\top \\ &= \sum_{k=1}^K \pi_k^{(t)} \left(x - \sqrt{\alpha_t}\mu_k - \sum_{i=1}^K \pi_i^{(t)} (x - \sqrt{\alpha_t}\mu_i) \right) \left(x - \sqrt{\alpha_t}\mu_k - \sum_{i=1}^K \pi_i^{(t)} (x - \sqrt{\alpha_t}\mu_i) \right)^\top \\ &= \sum_{k=1}^K \pi_k^{(t)} (x - \sqrt{\alpha_t}\mu_k) (x - \sqrt{\alpha_t}\mu_k)^\top - \left(\sum_{k=1}^K \pi_k^{(t)} (x - \sqrt{\alpha_t}\mu_k) \right) \left(\sum_{k=1}^K \pi_k^{(t)} (x - \sqrt{\alpha_t}\mu_k) \right)^\top. \end{aligned}$$

In addition, recall the definition of \mathcal{E}_t (cf. (21)). For any $x_t \in \mathcal{E}_t$, using (17) that $1 - \alpha_t \lesssim \log T/T$, we know that

$$\begin{aligned} \frac{1 - \alpha_t}{\alpha_t} \text{tr}(I_d + J_t(x_t)) &\lesssim (1 - \alpha_t) \text{tr}(I_d + J_t(x_t)) \lesssim \frac{\log(KT)}{T} = o(1), \\ \frac{1 - \alpha_t^2}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) &\lesssim (1 - \alpha_t) \text{tr}(I_d + J_t(x_t)) \lesssim \frac{\log(KT)}{T} = o(1), \end{aligned}$$

for large enough T . It follows that

$$\det \left(I_d + (\alpha_t^{-1} - 1)(I_d + J_t(x_t)) \right)$$

$$\begin{aligned}
&\stackrel{(i)}{=} 1 + \frac{1 - \alpha_t}{\alpha_t} \text{tr}(I_d + J_t(x_t)) + O\left(\frac{(1 - \alpha_t)^2}{\alpha_t^2} \text{tr}^2(I_d + J_t(x_t))\right) \\
&= 1 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) - \frac{(1 - \alpha_t)^2}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) + O\left(\frac{(1 - \alpha_t)^2}{\alpha_t^2} \text{tr}^2(I_d + J_t(x_t))\right) \\
&\stackrel{(ii)}{=} 1 + \frac{1 - \alpha_t^2}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) + O((1 - \alpha_t)^2 \log^2(KT)) \\
&= \exp\left(\frac{1 - \alpha_t^2}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) + O((1 - \alpha_t)^2 \log^2(KT))\right).
\end{aligned}$$

where (i) holds as $I_d + J_t(x_t) \succeq 0$ and $\det(I + \varepsilon A) = 1 + \text{tr}(A)\varepsilon + O(\varepsilon^2(\text{tr}^2(A) - \text{tr}(A^2)))$ for any matrix A and $\varepsilon > 0$; (ii) is true since $\alpha_t \gtrsim 1$ by (17) and $\text{tr}(I_d + J_t(x_t)) \lesssim \log(KT)$ by the choice of \mathcal{E}_t in (21). Consequently, we can derive

$$\begin{aligned}
\det(I_d + (1 - \alpha_t)J_t(x_t)) &= \det(\alpha_t I_d + (1 - \alpha_t)(I_d + J_t(x_t))) \\
&= \alpha_t^d \det(I_d + (\alpha_t^{-1} - 1)(I_d + J_t(x_t))) \\
&= \alpha_t^d \exp\left(\frac{(1 - \alpha_t^2)}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t)) + O((1 - \alpha_t)^2 \log^2(KT))\right), \tag{63}
\end{aligned}$$

As a consequence of the above two relations, we move on to prove Lemma 4. In view of relation (62), we arrive at

$$\begin{aligned}
&\left(\frac{1}{2\pi\alpha_t^2}\right)^{d/2} \sum_{k=1}^K \pi_k \exp\left(-\frac{1}{2\alpha_t^2} \|u_t - \sqrt{\alpha_t}\mu_k\|_2^2\right) \\
&= \left(\frac{1}{2\pi\alpha_t^2}\right)^{d/2} \exp\left(-\frac{(1 - \alpha_t^2)}{2\alpha_t^2} \text{tr}(I_d + J_t(x_t))\right) \sum_{k=1}^K \pi_k \exp\left(-\frac{1}{2} \|x_t - \sqrt{\alpha_t}\mu_k\|_2^2\right) \exp(-\zeta_k^{(t)}(x_t)) \\
&= \det(I_d + (1 - \alpha_t)J_t(x_t))^{-1} \exp\left(O((1 - \alpha_t)^2 \log^2(KT))\right) p_{X_t}(x_t) \sum_{k=1}^K \pi_k^t \exp(-\zeta_k^{(t)}(x_t))
\end{aligned}$$

where the last equality uses (63) and $\pi_k \exp(-\|x - \sqrt{\alpha_t}\mu_k\|_2^2/2) = \pi_k^{(t)}(2\pi)^{d/2} p_{X_t}(x)$ due to (18) and (15). To further control the right hand side, by the definition of \mathcal{E}_t in (21), it satisfies that

$$1 \leq \sum_{k=1}^K \pi_k^t \exp(-\zeta_k^{(t)}(x_t)) \leq \exp(C_2(1 - \alpha_t)^2 \log^2(KT)).$$

Therefore, we can conclude that

$$\begin{aligned}
&\left(\frac{1}{2\pi\alpha_t^2}\right)^{d/2} \sum_{k=1}^K \pi_k \exp\left(-\frac{1}{2\alpha_t^2} \|u_t - \sqrt{\alpha_t}\mu_k\|_2^2\right) \\
&= \det(I_d + (1 - \alpha_t)J_t(x_t))^{-1} \exp\left(O((1 - \alpha_t)^2 \log^2(KT))\right) p_{X_t}(x_t),
\end{aligned}$$

which completes the proof of Lemma 4.

A.4 Proof of Lemma 5

Recalling the definition of \mathcal{E}_t in expression (21), we have

$$\begin{aligned}
\mathbb{P}(X_t \in \mathcal{E}_t^c) &\leq \mathbb{P}\left(\text{tr}(I_d + J_t(X_t)) \geq C_1 \log(KT)\right) \\
&+ \mathbb{P}\left(\sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(X_t)) < 1 \text{ or } \sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(X_t)) > \exp(C_2(1 - \alpha_t)^2 \log^2(KT))\right). \tag{64}
\end{aligned}$$

In the following, we bound the two terms on the right respectively.

Before proceeding, we make the following observation. Fix an arbitrary $t \geq 1$. For each $k \in [K]$, we define the event

$$\mathcal{T}_k := \left\{ x \in \mathbb{R}^d : |(x - \sqrt{\bar{\alpha}_t} \mu_k)^\top \sqrt{\bar{\alpha}_t} (\mu_i - \mu_k)| \leq C_5 \sqrt{\bar{\alpha}_t \log(KT)} \|\mu_i - \mu_k\|_2 \text{ for all } i \in [K] \right\} \quad (65)$$

for some absolute constant $C_5 > 0$. Note that if we let $Z_k \sim \mathcal{N}(\sqrt{\bar{\alpha}_t} \mu_k, I_d)$ be a Gaussian random vector in \mathbb{R}^d , which implies that $(Z_k - \sqrt{\bar{\alpha}_t} \mu_k)^\top \sqrt{\bar{\alpha}_t} (\mu_i - \mu_k) \sim \mathcal{N}(0, \bar{\alpha}_t \|\mu_i - \mu_k\|_2^2)$, the standard Gaussian concentration inequality guarantees that

$$\mathbb{P}\{Z_k \notin \mathcal{T}_k\} \lesssim T^{-3}, \quad (66)$$

provided C_5 is large enough.

Bounding the first term in Eq. (64)

Let us begin with the first event $\{\text{tr}(I_d + J_t(X_t)) \leq C_1 \log(KT)\}$. As $X_t \sim \sum_{k=1}^K \pi_k \mathcal{N}(\sqrt{\bar{\alpha}_t} \mu_k, I_d)$, it is easily seen that

$$\begin{aligned} & \mathbb{P}\left\{\text{tr}(I_d + J_t(X_t)) > C_1 \log(KT)\right\} \\ &= \sum_{k=1}^K \pi_k \mathbb{P}\left\{\text{tr}(I_d + J_t(Z_k)) > C_1 \log(KT)\right\} \\ &\leq \sum_{k=1}^k \pi_k \mathbb{P}\left\{\text{tr}(I_d + J_t(Z_k)) > C_1 \log(KT)\right\} \mathbb{1}\{\pi_k \geq 1/(KT^3)\} + \sum_{k=1}^K \pi_k \mathbb{1}\{\pi_k < 1/(KT^3)\} \\ &\leq \sum_{k=1}^k \pi_k \mathbb{P}\left\{\text{tr}(I_d + J_t(Z_k)) > C_1 \log(KT)\right\} \mathbb{1}\{\pi_k \geq 1/(KT^3)\} + T^{-3}. \end{aligned}$$

We claim that for any $k \in [K]$ such that $\pi_k \geq 1/(KT^3)$, one has

$$\mathcal{T}_k \subset \left\{ x \in \mathbb{R}^d : \text{tr}(I_d + J_t(x)) \leq C_1 \log(KT) \right\}. \quad (67)$$

It then immediately follows from (66) that

$$\begin{aligned} \mathbb{P}\left\{\text{tr}(I_d + J_t(X_t)) > C_1 \log(KT)\right\} &\leq \sum_{k=1}^K \pi_k \mathbb{P}\{Z_k \notin \mathcal{T}_k\} \mathbb{1}\{\pi_k \geq 1/(KT^3)\} + T^{-3} \\ &\lesssim T^{-3} \sum_{k=1}^K \pi_k + T^{-3} \asymp T^{-3}. \end{aligned} \quad (68)$$

Proof of claim (67). It is sufficient to establish the claim (67). Towards this, fix an arbitrary $k \in [K]$ such that $\pi_k \geq 1/(KT^3)$. For any $x \in \mathcal{T}_k$, we know that for all $i \in [K]$,

$$\begin{aligned} \pi_i^{(t)} &\leq \frac{\pi_i}{\pi_k} \exp\left(-\frac{1}{2}\|x - \sqrt{\bar{\alpha}_t} \mu_i\|_2^2 + \frac{1}{2}\|x - \sqrt{\bar{\alpha}_t} \mu_k\|_2^2\right) \wedge 1 \\ &= \frac{\pi_i}{\pi_k} \exp\left(-\frac{1}{2}\bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 + (x - \sqrt{\bar{\alpha}_t} \mu_k)^\top \sqrt{\bar{\alpha}_t} (\mu_i - \mu_k)\right) \wedge 1 \\ &\leq \exp\left(-\frac{1}{2}\bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 + C_5 \sqrt{\bar{\alpha}_t \log(KT)} \|\mu_i - \mu_k\|_2 + 3 \log(KT)\right) \wedge 1, \end{aligned}$$

where the last line holds due to the definition of \mathcal{T}_k . As a result, for any $i \in [K]$ satisfying $\sqrt{\bar{\alpha}_t} \|\mu_i - \mu_k\|_2 > 6C_5 \sqrt{\log(KT)}$, one has

$$\pi_i^{(t)} \leq \exp\left(-\frac{1}{6}\bar{\alpha}_t \|\mu_i - \mu_k\|_2^2\right)$$

as long as $C_5 \geq \sqrt{2}/2$. This further implies that

$$\begin{aligned} \pi_i^{(t)} \bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 &\leq \bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 \exp\left(-\frac{1}{6}\bar{\alpha}_t \|\mu_i - \mu_k\|_2^2\right) \\ &\leq \exp\left(-\frac{1}{12}\bar{\alpha}_t \|\mu_i - \mu_k\|_2^2\right) \leq \exp(-3C_5^2 \log(KT)), \end{aligned} \quad (69)$$

provided T is large enough. Meanwhile, for any $i \in [K]$ obeying $\sqrt{\bar{\alpha}_t} \|\mu_i - \mu_k\|_2 \leq 6C_5 \sqrt{\log(KT)}$, we can simply upper bound

$$\pi_i^{(t)} \bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 \leq \pi_i^{(t)} \cdot 36C_5^2 \log(KT). \quad (70)$$

Denote the set $\mathcal{F}_k := \{i \in [K] : \sqrt{\bar{\alpha}_t} \|\mu_i - \mu_k\|_2 \leq 6C_5 \sqrt{\log(KT)}\}$. Using the expression of J_t (cf. (14)), we conclude that

$$\begin{aligned} \text{tr}(I_d + J_t(x)) &= \sum_{i=1}^K \pi_i^{(t)} \bar{\alpha}_t \left\| \mu_i - \sum_{k=1}^K \pi_k^{(t)} \mu_k \right\|_2^2 \\ &\stackrel{(i)}{\leq} \sum_{i=1}^K \pi_i^{(t)} \bar{\alpha}_t \|\mu_i - \mu_k\|_2^2 \\ &\stackrel{(ii)}{\leq} 36C_5^2 \log(KT) \sum_{i \in \mathcal{F}_k} \pi_i^{(t)} + \sum_{i \in \mathcal{F}_k^c} \exp(-3C_5^2 \log(KT)) \\ &\leq 36C_5^2 \log(KT) \log(KT) + K \exp(-3C_5^2 \log(KT)) \\ &\leq C_1 \log(KT) \end{aligned} \quad (71)$$

provided C_5 and C_1/C_5^2 are large enough. Here, (i) is true since $\sum_{k=1}^K \pi_k^{(t)} \mu_k$ is the minimizer of the function $x \mapsto \sum_{i=1}^K \pi_i^{(t)} \|\mu_i - x\|_2^2$ and (ii) uses (69)–(70). This establishes the claim (67).

Bounding the second term in Eq. (64)

Next, let analyze the second event $\{1 \leq \sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(X_t)) \leq \exp(C_2(1-\alpha_t)^2 \log^2(KT))\}$. We first establish the lower bound of 1. For any $x \in \mathbb{R}^d$, given $\sum_k \pi_k^{(t)} = 1$, direct calculation shows that

$$\begin{aligned} &\sum_{k=1}^K \pi_k^{(t)} \left(\|x - \sqrt{\bar{\alpha}_t} \mu_k\|_2^2 - \sum_{i=1}^K \pi_i^{(t)} \|x - \sqrt{\bar{\alpha}_t} \mu_i\|_2^2 \right) + \sum_{k=1}^K \pi_k^{(t)} \sum_{i=1}^K \pi_i^{(t)} (\mu_i - \mu_k) \\ &= \sum_{k=1}^K \pi_k^{(t)} \|x - \sqrt{\bar{\alpha}_t} \mu_k\|_2^2 - \sum_{i=1}^K \pi_i^{(t)} \|x - \sqrt{\bar{\alpha}_t} \mu_i\|_2^2 + \sum_{i=1}^K \pi_i^{(t)} \mu_i - \sum_{k=1}^K \pi_k^{(t)} \mu_k = 0. \end{aligned}$$

Combined with the definition of $\zeta_k^{(t)}(x)$ in (22), this yields

$$\sum_{k=1}^K \pi_k^{(t)} \zeta_k^{(t)}(x) = 0.$$

We can then apply Jensen's inequality to obtain that for any $x \in \mathbb{R}^d$,

$$\sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(x)) \geq \exp\left(-\sum_{k=1}^K \pi_k^{(t)} \zeta_k^{(t)}(x)\right) = 1. \quad (72)$$

Recall the definition of \mathcal{T}_k in expression (65). To bound the second term in Eq. (64), it suffices to prove that for any $k \in [K]$ such that $\pi_k \geq 1/(KT^3)$,

$$\mathcal{T}_k \subset \left\{ x \in \mathbb{R}^d : \sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(x)) \leq \exp(C_2(1-\alpha_t)^2 \log^2(KT)) \right\}. \quad (73)$$

Indeed, assuming (73) holds, one can apply the same reasoning as that for (68) to obtain

$$\begin{aligned} & \mathbb{P} \left\{ \sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(X_t)) > \exp(C_2(1-\alpha_t)^2 \log^2(KT)) \right\} \\ & \leq \sum_{k=1}^K \pi_k \mathbb{P}\{Z_k \notin \mathcal{T}_k\} \mathbb{1}\{\pi_k \geq 1/(KT^3)\} + T^{-3} \lesssim T^{-3}, \end{aligned} \quad (74)$$

Taking this collectively with relations (68), and (64) completes the proof of Lemma 5. Now it is only left for us to prove inequality (73).

Proof of inequality (73). To this end, recall the definitions of $\zeta_k^{(t)}(x)$ and $s_t^*(x)$ in (22) and (19), respectively. By some basic algebra, $\zeta_k^{(t)}(x)$ can be written as

$$\begin{aligned} \zeta_k^{(t)}(x) &= \frac{1-\alpha_t^2}{2\alpha_t^2} \sum_{i=1}^K \pi_i^{(t)} \left(\|x - \sqrt{\alpha_t} \mu_k\|_2^2 - \|x - \sqrt{\alpha_t} \mu_i\|_2^2 \right) \\ & \quad + \frac{1-\alpha_t}{\alpha_t^2} \left(-x + \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} \mu_i \right)^\top \left(\sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} (\mu_i - \mu_k) \right) \\ &= \frac{1-\alpha_t^2}{2\alpha_t^2} \sum_{i=1}^K \pi_i^{(t)} \left(-\frac{1}{2} \alpha_t \|\mu_i - \mu_k\|_2^2 + (x - \sqrt{\alpha_t} \mu_k)^\top \sqrt{\alpha_t} (\mu_i - \mu_k) \right) \\ & \quad - \frac{1-\alpha_t}{\alpha_t^2} \sum_{i=1}^K \pi_i^{(t)} (x - \sqrt{\alpha_t} \mu_k)^\top \sqrt{\alpha_t} (\mu_i - \mu_k) + \frac{1-\alpha_t}{\alpha_t^2} \left\| \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} (\mu_i - \mu_k) \right\|_2^2. \end{aligned}$$

For any $x \in \mathcal{T}_k$, one can obtain

$$\begin{aligned} |\zeta_k^{(t)}(x)| &\lesssim (1-\alpha_t) \sum_{i=1}^K \pi_i^{(t)} \left(-\frac{1}{2} \alpha_t \|\mu_i - \mu_k\|_2^2 + C_5 \sqrt{\alpha_t \log(KT)} \|\mu_i - \mu_k\|_2 \right) \\ & \quad + (1-\alpha_t) \sqrt{\log(KT)} \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} \|\mu_i - \mu_k\|_2 + (1-\alpha_t) \sum_{i=1}^K \pi_i^{(t)} \alpha_t \|\mu_i - \mu_k\|_2^2 \\ &\asymp (1-\alpha_t) \sum_{i=1}^K \pi_i^{(t)} \alpha_t \|\mu_i - \mu_k\|_2^2 + (1-\alpha_t) \sqrt{\log(KT)} \sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} \|\mu_i - \mu_k\|_2. \end{aligned} \quad (75)$$

where the first inequality holds due to $1-\alpha_t \lesssim \log T/T$ in (17), the definition of \mathcal{T}_k in (65), and Jensen's inequality. Using the same argument as that for (71), it can be easily seen that

$$\sum_{i=1}^K \pi_i^{(t)} \sqrt{\alpha_t} \|\mu_i - \mu_k\|_2 \lesssim \sqrt{\log(KT)}. \quad (76)$$

Plugging (76) and (71) into (75) demonstrates that

$$|\zeta_k^{(t)}(x)| \lesssim (1-\alpha_t) \log(KT) = o(1), \quad (77)$$

since $1-\alpha_t \lesssim \log T/T$ as in (17).

As a consequence, for any $x \in \mathcal{T}_k$, we find that

$$\begin{aligned} \sum_{k=1}^K \pi_k^{(t)} \exp(-\zeta_k^{(t)}(x)) &= \sum_{k=1}^K \pi_k^{(t)} \left(1 - \zeta_k^{(t)}(x) + \frac{1}{2} (\zeta_k^{(t)}(x))^2 + o\left((\zeta_k^{(t)}(x))^2\right) \right) \\ &= 1 + \frac{1}{2} \sum_{k=1}^K \pi_k^{(t)} (\zeta_k^{(t)}(x))^2 + \sum_{k=1}^K \pi_k^{(t)} o\left((\zeta_k^{(t)}(x))^2\right) \\ &= 1 + O\left((1 - \alpha_t)^2 \log^2(KT)\right) \\ &\leq \exp\left(C_2(1 - \alpha_t)^2 \log^2(KT)\right) \end{aligned}$$

as long as C_2 is sufficiently large. This establishes the claim (73), thereby leads to (74).

References

- Anderson, B. D. (1982). Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326. [1.1](#), [2](#)
- Arora, S. and Kannan, R. (2005). Learning mixtures of separated nonspherical Gaussians. *The Annals of Applied Probability*, 15(1A):69 – 92. [1.3](#)
- Ashtiani, H., Ben-David, S., Harvey, N., Liaw, C., Mehrabian, A., and Plan, Y. (2018). Nearly tight sample complexity bounds for learning mixtures of gaussians via sample compression schemes. *Advances in Neural Information Processing Systems*, 31. [1.3](#)
- Bakshi, A., Diakonikolas, I., Jia, H., Kane, D. M., Kothari, P. K., and Vempala, S. S. (2022). Robustly learning mixtures of k arbitrary gaussians. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1234–1247. [1.3](#)
- Benton, J., Bortoli, V., Doucet, A., and Deligiannidis, G. (2024). Nearly d-linear convergence bounds for diffusion models via stochastic localization. [1](#), [1.1](#), [3](#), [3](#)
- Cai, C. and Li, G. (2025). Minimax optimality of the probability flow ode for diffusion models. *arXiv preprint arXiv:2503.09583*. [1.3](#)
- 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. PMLR. [1.1](#), [3](#)
- Chen, J. (1995). Optimal rate of convergence for finite mixture models. *The Annals of Statistics*, pages 221–233. [1.3](#)
- Chen, M., Huang, K., Zhao, T., and Wang, M. (2023b). Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In *International Conference on Machine Learning*, pages 4672–4712. PMLR. [1.3](#)
- Chen, S., Chewi, S., Lee, H., Li, Y., Lu, J., and Salim, A. (2023c). The probability flow ode is provably fast. *Advances in Neural Information Processing Systems*, 36:68552–68575. [1.1](#)
- 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*. [1](#), [1.1](#), [2](#), [3](#), [3](#)
- Chen, S., Kontonis, V., and Shah, K. (2024). Learning general gaussian mixtures with efficient score matching. *arXiv preprint arXiv:2404.18893*. [1.2](#), [3](#), [5](#)
- Chidambaram, M., Gattmiry, K., Chen, S., Lee, H., and Lu, J. (2024). What does guidance do? a fine-grained analysis in a simple setting. *arXiv preprint arXiv:2409.13074*. [1.2](#)

- Cole, F. and Lu, Y. (2024). Score-based generative models break the curse of dimensionality in learning a family of sub-gaussian distributions. In *The Twelfth International Conference on Learning Representations*. [1.3](#)
- Dasgupta, S. (1999). Learning mixtures of gaussians. In *40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039)*, pages 634–644. IEEE. [1.3](#)
- De Bortoli, V. (2022). Convergence of denoising diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2208.05314*. [1.3](#)
- Dhariwal, P. and Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794. [1](#)
- Diakonikolas, I. and Kane, D. M. (2020). Small covers for near-zero sets of polynomials and learning latent variable models. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 184–195. IEEE. [1.3](#)
- Diakonikolas, I., Kane, D. M., and Stewart, A. (2018). List-decodable robust mean estimation and learning mixtures of spherical gaussians. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1047–1060. [1](#), [1.3](#)
- Doss, N., Wu, Y., Yang, P., and Zhou, H. H. (2023). Optimal estimation of high-dimensional gaussian location mixtures. *The Annals of Statistics*, 51(1):62–95. [1.3](#)
- Dou, Z., Kotekal, S., Xu, Z., and Zhou, H. H. (2024). From optimal score matching to optimal sampling. *arXiv preprint arXiv:2409.07032*. [1.3](#)
- Dwivedi, R., Ho, N., Khamaru, K., Wainwright, M. J., Jordan, M. I., and Yu, B. (2020). Singularity, misspecification and the convergence rate of em. *The Annals of Statistics*, 48(6):3161–3182. [1.3](#)
- Feng, O. Y., Kao, Y.-C., Xu, M., and Samworth, R. J. (2024). Optimal convex m -estimation via score matching. *arXiv preprint arXiv:2403.16688*. [1.3](#)
- Gao, X. and Zhu, L. (2024). Convergence analysis for general probability flow odes of diffusion models in wasserstein distances. *arXiv preprint arXiv:2401.17958*. [1.1](#)
- Gatmiry, K., Kelner, J., and Lee, H. (2024). Learning mixtures of gaussians using diffusion models. *arXiv preprint arXiv:2404.18869*. [1.2](#), [3](#)
- Genovese, C. R. and Wasserman, L. (2000). Rates of convergence for the gaussian mixture sieve. *The Annals of Statistics*, 28(4):1105–1127. [1.3](#)
- Ghosal, S. and Van Der Vaart, A. W. (2001). Entropies and rates of convergence for maximum likelihood and bayes estimation for mixtures of normal densities. *The Annals of Statistics*, 29(5):1233–1263. [1.3](#)
- Hausmann, U. G. and Pardoux, E. (1986). Time reversal of diffusions. *The Annals of Probability*, pages 1188–1205. [1.1](#), [2](#)
- Heinrich, P. and Kahn, J. (2018). Strong identifiability and optimal minimax rates for finite mixture estimation. [1.3](#)
- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851. [1](#), [1.1](#), [2](#)
- Ho, N. and Nguyen, X. (2016). Convergence rates of parameter estimation for some weakly identifiable finite mixtures. [1.3](#)
- Hopkins, S. B. and Li, J. (2018). Mixture models, robustness, and sum of squares proofs. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1021–1034. [1.3](#)

- Hsu, D. and Kakade, S. M. (2013). Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 11–20. [1.3](#)
- Huang, D. Z., Huang, J., and Lin, Z. (2024a). Convergence analysis of probability flow ODE for score-based generative models. *arXiv preprint arXiv:2404.09730*. [1](#), [1.1](#)
- Huang, X., Zou, D., Dong, H., Zhang, Y., Ma, Y.-A., and Zhang, T. (2024b). Reverse transition kernel: A flexible framework to accelerate diffusion inference. *arXiv preprint arXiv:2405.16387*. [1](#)
- Huang, Z., Wei, Y., and Chen, Y. (2024c). Denoising diffusion probabilistic models are optimally adaptive to unknown low dimensionality. *arXiv preprint arXiv:2410.18784*. [1](#)
- Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4). [1.3](#), [2](#)
- Kalai, A. T., Moitra, A., and Valiant, G. (2010). Efficiently learning mixtures of two gaussians. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 553–562. [1.3](#)
- Kim, A. K. and Guntuboyina, A. (2022). Minimax bounds for estimating multivariate gaussian location mixtures. *Electronic Journal of Statistics*, 16(1):1461–1484. [1.3](#)
- Kothari, P. K., Steinhardt, J., and Steurer, D. (2018). Robust moment estimation and improved clustering via sum of squares. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1035–1046. [1.3](#)
- Kwon, H. K., Kim, D., Ohn, I., and Chae, M. (2025). Nonparametric estimation of a factorizable density using diffusion models. *arXiv preprint arXiv:2501.01783*. [1.3](#)
- 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. PMLR. [1.1](#), [3](#)
- Li, G. and Cai, C. (2024). Provable acceleration for diffusion models under minimal assumptions. *arXiv preprint arXiv:2410.23285*. [1](#), [3](#), [4.1](#)
- Li, G., Huang, Y., Efimov, T., Wei, Y., Chi, Y., and Chen, Y. (2024a). Accelerating convergence of score-based diffusion models, provably. *arXiv preprint arXiv:2403.03852*. [1](#)
- Li, G. and Jiao, Y. (2024). Improved convergence rate for diffusion probabilistic models. *arXiv preprint arXiv:2410.13738*. [1](#)
- 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*. [1.1](#), [3](#)
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2024b). A sharp convergence theory for the probability flow odes of diffusion models. *arXiv preprint arXiv:2408.02320*. [1.1](#), [3](#)
- Li, G. and Yan, Y. (2024a). Adapting to unknown low-dimensional structures in score-based diffusion models. *arXiv preprint arXiv:2405.14861*. [1](#)
- Li, G. and Yan, Y. (2024b). $O(d/T)$ convergence theory for diffusion probabilistic models under minimal assumptions. *arXiv preprint arXiv:2409.18959*. [1.1](#), [3](#), [3](#), [3](#), [4.4](#)
- Liang, J., Huang, Z., and Chen, Y. (2025). Low-dimensional adaptation of diffusion models: Convergence in total variation. *arXiv preprint arXiv:2501.12982*. [1](#)
- Liang, Y., Ju, P., Liang, Y., and Shroff, N. (2024a). Non-asymptotic convergence of discrete-time diffusion models: New approach and improved rate. *arXiv preprint arXiv:2402.13901*. [1.1](#)
- Liang, Y., Shi, Z., Song, Z., and Zhou, Y. (2024b). Unraveling the smoothness properties of diffusion models: A gaussian mixture perspective. *arXiv preprint arXiv:2405.16418*. [1.2](#), [3](#)

- Liu, A. and Li, J. (2022). Clustering mixtures with almost optimal separation in polynomial time. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1248–1261. [1.3](#)
- Liu, X., Wu, L., Ye, M., and Liu, Q. (2022). Let us build bridges: Understanding and extending diffusion generative models. *arXiv preprint arXiv:2208.14699*. [1.1](#)
- Mei, S. and Wu, Y. (2023). Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. *arXiv preprint arXiv:2309.11420*. [1.3](#)
- Moitra, A. and Valiant, G. (2010). Settling the polynomial learnability of mixtures of gaussians. In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 93–102. IEEE. [1.3](#)
- Oko, K., Akiyama, S., and Suzuki, T. (2023). Diffusion models are minimax optimal distribution estimators. In *International Conference on Machine Learning*, pages 26517–26582. PMLR. [1.3](#)
- Pearson, K. (1894). Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London. A*, 185:71–110. [1](#)
- Polyanskiy, Y. and Wu, Y. (2020). Self-regularizing property of nonparametric maximum likelihood estimator in mixture models. *arXiv preprint arXiv:2008.08244*. [1.3](#)
- Potapchik, P., Azangulov, I., and Deligiannidis, G. (2024). Linear convergence of diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2410.09046*. [1](#)
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3. [1](#)
- Saha, S. and Guntuboyina, A. (2020). On the nonparametric maximum likelihood estimator for gaussian location mixture densities with application to gaussian denoising. *The Annals of Statistics*, 48(2):738–762. [1.3](#)
- Shah, K., Chen, S., and Klivans, A. (2023). Learning mixtures of gaussians using the ddpm objective. *Advances in Neural Information Processing Systems*, 36:19636–19649. [1.2](#)
- Song, J., Meng, C., and Ermon, S. (2020a). Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*. [1.1](#)
- Song, Y., Garg, S., Shi, J., and Ermon, S. (2020b). Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR. [1.3](#)
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. (2020c). Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*. [1](#)
- Taheri, M. and Lederer, J. (2025). Regularization can make diffusion models more efficient. *arXiv preprint arXiv:2502.09151*. [1](#)
- Tang, R. and Yang, Y. (2024). Adaptivity of diffusion models to manifold structures. In *International Conference on Artificial Intelligence and Statistics*, pages 1648–1656. PMLR. [1](#)
- Tang, W. and Zhao, H. (2024). Score-based diffusion models via stochastic differential equations—a technical tutorial. *arXiv preprint arXiv:2402.07487*. [1.1](#)
- Vempala, S. and Wang, G. (2004). A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841–860. [1.3](#)
- Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674. [2](#)
- Wang, P., Zhang, H., Zhang, Z., Chen, S., Ma, Y., and Qu, Q. (2024). Diffusion models learn low-dimensional distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*. [1.2](#)

- Wibisono, A., Wu, Y., and Yang, K. Y. (2024). Optimal score estimation via empirical bayes smoothing. *arXiv preprint arXiv:2402.07747*. [1.3](#)
- Wu, Y., Chen, M., Li, Z., Wang, M., and Wei, Y. (2024a). Theoretical insights for diffusion guidance: A case study for gaussian mixture models. *arXiv preprint arXiv:2403.01639*. [1.2](#)
- Wu, Y., Chen, Y., and Wei, Y. (2024b). Stochastic runge-kutta methods: Provable acceleration of diffusion models. *arXiv preprint arXiv:2410.04760*. [1](#)
- Wu, Y. and Yang, P. (2020). Optimal estimation of gaussian mixtures via denoised method of moments. *The Annals of Statistics*, 48(4):1981–2007. [1.3](#)
- Zhang, C.-H. (2009). Generalized maximum likelihood estimation of normal mixture densities. *Statistica Sinica*, pages 1297–1318. [1.3](#)
- Zhang, K., Yin, H., Liang, F., and Liu, J. (2024). Minimax optimality of score-based diffusion models: Beyond the density lower bound assumptions. *arXiv preprint arXiv:2402.15602*. [1.3](#)