

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 IMPROVED SAMPLE COMPLEXITY FOR DIFFUSION MODEL TRAINING WITHOUT EMPIRICAL RISK MIN- IMIZER ACCESS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Diffusion models have demonstrated state-of-the-art performance across vision, language, and scientific domains. Despite their empirical success, prior theoretical analyses of the sample complexity suffer from poor scaling with input data dimension or rely on unrealistic assumptions such as access to exact empirical risk minimizers. In this work, we provide a principled analysis of score estimation, establishing a sample complexity bound of  $\mathcal{O}(\epsilon^{-4})$ . Our approach leverages a structured decomposition of the score estimation error into statistical, approximation, and optimization errors, enabling us to eliminate the exponential dependence on neural network parameters that arises in prior analyses. It is the first such result that achieves sample complexity bounds without assuming access to the empirical risk minimizer of score function estimation loss.

## 1 INTRODUCTION

Diffusion models have emerged as a powerful class of generative models, achieving impressive performance across tasks such as image synthesis, molecular design, and audio generation. Central to the training of these models is the estimation of the *score function*, which characterizes the reverse-time dynamics in the diffusion process. Diffusion models are widely adopted in computer vision and audio generation tasks (Ulhaq & Akhtar, 2022; Bansal et al., 2023), text generation (Li et al., 2022), sequential data modeling (Tashiro et al., 2021), reinforcement learning and control (Zhu et al., 2023), and life sciences (Jing et al., 2022; Malusare & Aggarwal, 2024). For a more comprehensive exposition of applications, we refer readers to survey paper (Chen et al., 2024).

While diffusion models exhibit strong empirical performance, understanding their sample complexity is essential to guarantee their efficiency, generalization, and scalability, enabling high-quality generation with minimal data in real-world, resource-constrained scenarios. Some of the key works studying the sample complexity are summarized in Table 1. A key limitation of sample complexity analyses of diffusion models done thus far is the lack of the presence of finite-time sample complexity results under reasonable assumptions. This makes the theoretical analysis of diffusion models fall short of other machine learning areas such as reinforcement Learning (Kumar et al., 2023; Gaur et al., 2024), bi-level-optimization (Grazzi et al., 2023; Gaur et al., 2025) and graphical models (Fattah et al., 2019; Tran et al., 2019). In this work we aim to bridge that gap and obtain a sample complexity results on the same footing as results from the aforementioned areas. The iteration complexity or convergence has been studied in Li et al. (2024b); Benton et al. (2024); Li & Yan (2024); Huang et al. (2024); Dou et al. (2024); Liang et al. (2025a;b), while they assume bounded score estimates thus not providing the sample complexity which requires estimating the score function.

We note that works such as Zhang et al. (2024); Wibisono et al. (2024); Oko et al. (2023); Chen et al. (2023) have sample complexity results that depend exponentially on the data dimension, making the result less useful in high-dimensional settings. Recently, Gupta et al. (2024) improved upon this by obtaining  $\tilde{\mathcal{O}}(\epsilon^{-5})$  sample complexity without exponential dependence on data dimension.<sup>1</sup>.

<sup>1</sup>We note that Gupta et al. (2024) claimed a sample complexity of  $\tilde{\mathcal{O}}(\epsilon^{-3})$ , while this claim does not account for the accumulation of errors across discretization steps. Specifically, their bound at each step depends on  $\tilde{\mathcal{O}}(1/\epsilon^2)$  samples, and applying a union bound over  $\tilde{\mathcal{O}}(1/\epsilon^2)$  steps yields a total sample complexity of  $\tilde{\mathcal{O}}(\epsilon^{-5})$ . For more details, see Appendix B.

Reference	Sample Complexity	Empirical Risk Minimizer Assumption
Zhang et al. (2024)	$\tilde{O}(\epsilon^{-d})$	Yes
Wibisono et al. (2024)	$\tilde{O}(\epsilon^{-(d)})$	Yes
Oko et al. (2023)	$\tilde{O}(\epsilon^{-O(d)})$	Yes
Chen et al. (2023)	$\tilde{O}(\epsilon^{-O(d)})$	Yes
Gupta et al. (2024) <sup>1</sup>	$\tilde{O}(\epsilon^{-5})$	Yes
<b>This work</b>	$\tilde{O}(\epsilon^{-4})$	No

Table 1: Summary of sample complexity results for diffusion models, assuming no upper bound on score estimation error. For further details on how the sample complexity bounds are derived for Gupta et al. (2024), see Appendix B.

However, this work assumes access to the empirical risk minimizer (ERM) of the score estimation loss, a significant restriction that was explicitly highlighted as an open problem in Gupta et al. (2024) itself. While this assumption is present in all prior works, it is an unrealistic assumption regardless.

In this paper, we do not make this assumptions and establish an improved state-of-the-art sample complexity bound of  $\tilde{O}(\epsilon^{-4})$ . This represents a key step toward bridging the gap between the theory and practice of diffusion models. Specifically, we address the following fundamental question.

*How many samples are required for a sufficiently expressive neural network to estimate the score function well enough to generate high-quality samples using a DDPM algorithm?*

Our analysis directly connects the quality of the learned score function to the total variation distance between the generated and target distributions, offering more interpretable and practically relevant guarantees. Additionally, our work accounts for the unavailability of the empirical risk minimizer. Using our novel analysis of the score estimation error, we obtain the sample complexity bounds without exponential dependence on the data-dimension. Our principled analysis accounts for the statistical and optimization errors while not assuming access to the empirical risk minimizer of the score estimation loss, and achieves state-of-the-art sample complexity bounds.

The statistical error occurs due to the finite sample size used to obtain the score estimate. Existing methods used to bound statistical errors assume bounded loss functions, which is not true in the case of diffusion models. We thus use a novel analysis that uses the conditional normality of the score function to obtain upper bounds for the statistical error.

Finally, the optimization error occurs due to a finite number of SGD steps during the estimation of the score function. It is precisely the error that was not accounted for in the previous works due to the assumption that they have access to the empirical risk minimizer. We use the quadratic growth property implied by the Polyak-Łojasiewicz (PL) condition assumed in Assumption 2 and a novel recursive analysis of the error at each stochastic gradient descent (SGD) step to upper bound this error.

The main contributions of our work are summarized as:

- **Finite time sample complexity bounds.** We derive state-of-the-art sample complexity bound of  $\tilde{O}(\epsilon^{-4})$  for score-based diffusion models, without exponential dependence on the data dimension or neural network parameters. Our analysis avoids the unrealistic assumptions used in prior works such as access to an empirical loss minimizer.
- **Principled error decomposition.** We propose a structured decomposition of the score estimation error into approximation, statistical, and optimization components, enabling the characterization of how each factor contributes to sample complexity.

**Unconditional and Conditional Diffusion Models.** Diffusion models have emerged as leading frameworks across vision, audio, and scientific domains. Foundational works such as Sohl-Dickstein et al. (2015) and Ho et al. (2020) introduced and refined Denoising Diffusion Probabilistic Models (DDPMs), enabling high-quality sample generation. Subsequent advances include improved noise

108 schedules (Nichol & Dhariwal, 2021), score-based SDE formulations (Song et al., 2021), and ef-  
 109 ficient latent-space generation via Latent Diffusion Models (LDMs) (Rombach et al., 2022). Con-  
 110 ditional diffusion models extend these techniques for guided generation tasks, with applica-  
 111 tions in time-series Tashiro et al. (2021), speech Huang et al. (2022), and medical imaging Dorjsembe  
 112 et al. (2023). Conditioning mechanisms range from classifier-based Dhariwal & Nichol (2021) to  
 113 classifier-free guidance Ho & Salimans (2022), which enabled text-to-image models like Imagen  
 114 Saharia et al. (2022) and Stable Diffusion Rombach et al. (2022). Recent innovations focus on adap-  
 115 tive control Castillo et al. (2025), compositionality Liu et al. (2023), and multi-modal conditioning  
 116 Avrahami et al. (2022).

117 **Related Works:** Despite the empirical success of diffusion models, theoretical understanding  
 118 regarding the sample complexity remains limited. Assuming accurate score estimates, authors in  
 119 Chen et al. (2022) showed that score-based generative models can efficiently sample from a sub-  
 120 Gaussian data distribution. Assuming a bounded score function, iteration complexity bounds have  
 121 been extensively studied in recent works Li et al. (2024b); Benton et al. (2024); Li & Yan (2024);  
 122 Huang et al. (2024); Dou et al. (2024); Liang et al. (2025a;b). Some works, such as Zhang & Pilanci  
 123 (2024), establish iteration complexity for score matching. In particular, Benton et al. (2024); Li et al.  
 124 (2024b) establishes iteration complexity guarantees for DDPM algorithms. Several studies propose  
 125 accelerated denoising schedules to improve these rates Li et al. (2024a); Liang et al. (2025b); Dou  
 126 et al. (2024). Additionally, improved convergence rates under low-dimensional data assumptions  
 127 are demonstrated in Li & Yan (2024); Huang et al. (2024); Liang et al. (2025a). In contrast to these  
 128 works, our analysis addresses the sample complexity of score-based generative models, where the  
 129 errors introduced by the neural network approximation, data sampling, and optimization are also  
 130 accounted for.

131 Sample complexity bounds for diffusion models have been studied via diffusion SDEs under  
 132 smoothness and spectral assumptions in Chen et al. (2023); Zhang et al. (2024); Wibisono et al.  
 133 (2024); Oko et al. (2023). However, these bounds are exponential in the data dimension. Further,  
 134 authors of Gupta et al. (2024) use the quantile-based approach to get the sample complexity bounds.  
 135 In this work, we further improve on these guarantees. The detailed comparison of sample complex-  
 136 ities with the key approaches mentioned above is provided in Table 1.

## 137 2 PRELIMINARIES AND PROBLEM FORMULATION

139 We begin by outlining the theoretical basis of score-based diffusion models. In particular, we adopt  
 140 the continuous-time stochastic differential equation (SDE) framework, which provides a principled  
 141 basis for modeling the generative process. We then outline its practical discretization and formally  
 142 define our problem.

143 Score-based generative models enable sampling from a complex distribution  $p_0$  by learning to re-  
 144 verse a noise-adding diffusion process. This approach introduces a continuous-time stochastic pro-  
 145 cess that incrementally perturbs the data distribution into a tractable distribution (typically Gaus-  
 146 sian), and then seeks to reverse that transformation.

147 A canonical forward process used in diffusion models is the *Ornstein–Uhlenbeck (OU) process* (Øksendal, 2003), defined by the following SDE

$$150 \quad dx_t = -x_t dt + \sqrt{2} dB_t, \quad x_0 \sim p_0, x \subset \mathbb{R}^d, \quad (1)$$

151 where  $B_t$  denotes standard Brownian motion. The solution of this SDE in closed form is given by

$$153 \quad x_t \sim e^{-t} x_0 + \sqrt{1 - e^{-2t}} \epsilon, \quad \text{with } \epsilon \sim \mathcal{N}(0, I). \quad (2)$$

154 As  $t \rightarrow \infty$ , the process converges to the stationary distribution  $\mathcal{N}(0, I)$ . Let  $p_t$  denote the marginal  
 155 distribution of  $x_t$ . This defines a continuous-time smoothing of the data distribution, where  $p_t$   
 156 becomes increasingly Gaussian over time.

157 The reverse process is typically achieved using stochastic time-reversal theory (Anderson, 1982),  
 158 which yields a corresponding reverse-time SDE as follows.

$$159 \quad dx_{T-t} = (x_{T-t} + 2\nabla \log p_{T-t}(x_{T-t})) dt + \sqrt{2} dB_t, \quad (3)$$

160 where  $\nabla \log p_t(x)$  is known as the *score function* of the distribution  $p_t$ . Simulating this reverse  
 161 process starting from  $x_T \sim p_T \approx \mathcal{N}(0, I)$  yields approximate samples from the original distribution

*p*<sub>0</sub>. This motivates a sampling strategy where we begin from  $x_T \sim \mathcal{N}(0, I)$  for sufficiently large  $T$ , and then integrate the reverse SDE backward to  $t = 0$  using estimated score functions. In practice, the backward process is run up to a fixed time point  $t_0$  known as the *early stopping time* and not  $t = 0$ . This is done in order to improve performance and training speed (Lyu et al., 2022; Favero et al., 2025).

The continuous-time reverse SDE (Equation 3) is discretized over a finite sequence of times  $0 < t_0 < t_1 < \dots, t_k, \dots < t_K = (T - \kappa) < T$ . The score function  $s_t(x) := \nabla \log p_t(x)$  is approximated at these discrete points using a learned estimator  $\hat{s}_{t_k}$ . This discretization underlies the DDPM framework (Ho et al., 2020), where the reverse process is implemented by iteratively denoising the sample using the estimated scores at each time step. The detailed procedure is provided in Algorithm 1 in the Appendix C. We employ stochastic gradient descent (SGD) to learn the score function at each  $t_k$ , using either a constant learning rate, as justified in our analysis later.

**Problem formulation:** Let the score function be approximated using a parameterized family of neural networks  $\mathcal{F}_\Theta = \{s_\theta : \theta \in \Theta\}$ , where each  $s_\theta : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$  is represented by a neural network of depth  $D$  and width  $W$  with smooth activation functions. Given  $n_k$  i.i.d. samples  $\{x_i\}_{i=1}^n$  from the data distribution  $p_{t_k}$ , the score network is trained by minimizing the following time-indexed loss:

$$\mathcal{L}_k(\theta) := \mathbb{E}_{x \sim p_{t_k}} [\|s_\theta(x, t_k) - \nabla \log p_{t_k}(x)\|^2]. \quad (4)$$

**Objective.** Our goal is to quantify how well the learned generative model  $\hat{p}_{t_0}$  approximates the true data distribution  $p$  in terms of total variation (TV) distance. Specifically, we aim to show the number of samples needed so that with high probability, the TV distance  $\text{TV}(p_{t_0}, \hat{p}_{t_0})$  is bounded by  $\mathcal{O}(\epsilon)$ , where  $\epsilon$  is the  $L^2$  estimation error of the score function. This reduces the generative performance analysis to establishing tight sample complexity bounds on the score estimation error. We additionally define the following probability distributions:

$$\begin{aligned} p_{t_0} &:= \text{Distribution obtained after backward process till time } t_0 \text{ steps starting from } p_T \\ p_{t_0}^{dis} &:= \text{Distribution obtained by backward process till time } t_0 \text{ starting from } p_T \\ &\quad \text{at discretized time steps} \\ \tilde{p}_{t_0} &:= \text{Distribution obtained by backward process till time } t_0 \text{ starting from } p_T \\ &\quad \text{at discretized time steps using the estimated score functions} \\ \hat{p}_{t_0} &:= \text{Distribution obtained by backward process till time } t_0 \text{ starting from } \mathcal{N}(0, I) \\ &\quad \text{at discretized time steps using the estimated score functions} \end{aligned}$$

where  $t_0$  denotes the early stopping time.

### 3 SAMPLE COMPLEXITY OF DIFFUSION MODELS

In this section, we derive explicit sample complexity bounds for diffusion-based generative models. By leveraging tools from stochastic optimization and statistical learning theory, we provide bounds on the number of data samples required to accurately estimate the time-dependent score function  $s_t(x) := \nabla \log p_t(x)$  across the forward diffusion process. Note that accurate score estimation is critical for ensuring high-quality generation while sampling through the reverse-time SDE.

We first state the assumptions required throughout this work.

**Assumption 1** (Bounded Second Moment Data Distribution.). *The data distribution  $p_0$  of the data variable  $x_0$  has an absolutely continuous CDF, is supported on a continuous set  $\Gamma \in \mathbb{R}^d$ , and there exists a constant  $0 < C_1 < \infty$  such that  $\mathbb{E}(\|x_0\|^2) \leq C_1$ .*

Some works that analyze the convergence of score-based diffusion models, such as Chen et al. (2022); Oko et al. (2023), assume that the data distribution is supported on a bounded set, thereby excluding commonly encountered distributions such as Gaussian and sub-Gaussian families. In contrast, our analysis only requires the data distribution to have a finite second moment, making our results applicable to a significantly broader class of distributions.

216 **Assumption 2** (Polyak - Łojasiewicz (PL) condition.). *The loss  $\mathcal{L}_k(\theta)$  for all  $k \in [0, K]$  satisfies*  
 217 *the Polyak–Łojasiewicz condition, i.e., there exists a constant  $\mu > 0$  such that*  
 218

$$219 \quad \frac{1}{2} \|\nabla \mathcal{L}_k(\theta)\|^2 \geq \mu (\mathcal{L}_k(\theta) - \mathcal{L}_k(\theta^*)) , \quad \forall \theta \in \Theta, \quad (5)$$

221 where  $\theta^* = \arg \min_{\theta \in \Theta} \mathcal{L}_k(\theta)$  denotes the global minimizer of the population loss.  
 222

223 The Polyak–Łojasiewicz (PL) condition is significantly weaker than strong convexity and is known to  
 224 hold in many non-convex settings, including overparameterized neural networks trained with mean  
 225 squared error losses (Liu et al., 2022). Prior works such as Gupta et al. (2024) and Block et al. (2020)  
 226 implicitly assume access to an exact empirical risk minimizer (ERM) for score function estimation,  
 227 as reflected in their sample complexity analyses (see Assumption A2 in Gupta et al. (2024) and the  
 228 definition of  $\hat{f}$  in Theorem 13 of Block et al. (2020)). This assumption, however, introduces a major  
 229 limitation for practical implementations, where exact ERM is not attainable.

230 In contrast, the PL condition allows us to derive sample complexity bounds under realistic optimiza-  
 231 tion dynamics, without requiring exact ERM solutions. To our knowledge, this is the first theoretical  
 232 analysis of score-based generative models that explicitly accounts for inexact optimization, address-  
 233 ing a key gap in existing literature. Additionally, we establish convergence guarantees with both  
 234 constant and decreasing step sizes.

235 **Assumption 3** (Approximation error of the Class of Neural Networks). *For all  $t \in [0, T]$ , there*  
 236 *exists a neural network parameter  $\theta \in \Theta$  such that*

$$237 \quad \mathbb{E}_{x \sim p_t} \|s_\theta(x, t) - \nabla \log p_t(x)\|^2 \leq \epsilon_{approx} \quad (6)$$

239 This error is independent of the sampling algorithm, and describes the error due to neural network  
 240 parametrization. In learning theory, it is common to treat the *approximation error* of a model class  
 241 as a constant so that analyses can focus on the estimation/ optimization terms dependent on the sam-  
 242 ple. This convention appears in standard excess-risk decompositions for fixed hypothesis classes  
 243 (Shalev-Shwartz & Ben-David, 2014). In PAC-Bayesian analyses, approximation errors are denoted  
 244 by a constant once the class is fixed (Mai, 2025). In (NTK/RKHS) analyses of neural networks,  
 245 where it is assumed the target function lies in, or is well approximated by the specified function  
 246 class, the misspecification error is represented as a constant term (Bing et al., 2025). In reinforce-  
 247 ment learning algorithm analysis such as policy gradient, a task-dependent “inherent Bellman” or  
 248 function-approximation error that remains constant while deriving performance rates (Mondal &  
 249 Aggarwal, 2024; Fu et al., 2021; Gaur et al., 2024; Ganesh et al., 2025). Note that Gupta et al. (2024)  
 250 also make the same assumption implicitly, but assume this constant to be zero (in Assumption A2).  
 251 Note that in Gupta et al. (2024), Assumption A.2 states that the error in estimating the loss function  
 252 is ‘sufficiently small’. In practice, this assumption is used to make the score function estimation  
 253 error arbitrarily small, as is done in Theorem C.3, where it is stated that the there exists a neural net-  
 254 work such that the error in estimating the loss function is  $\mathcal{O}(\epsilon^3)$ . Thus, this is a stronger assumption  
 255 as compared to our Assumption 3.

256 Note that in certain works, such as (Jiao et al., 2023), it is shown that the network size has to be  
 257 exponential in data dimension in order to achieve a small approximation error. However, in practice,  
 258 that would require an impractically large neural network size. In practice neural network size is of  
 259 the same order as the data dimension. Thus for a fixed neural network size that we assume in this  
 260 work, it makes sense to assume the approximation error as a constant.

261 **Assumption 4** (Smoothness and bounded gradient variance of the score loss.). *For all  $k \in [0, K]$ ,*  
 262 *the population loss  $\mathcal{L}_k(\theta)$  is  $\kappa$ -smooth with respect to the parameters  $\theta$ , i.e., for all  $\theta, \theta' \in \Theta$*

$$263 \quad \|\nabla \mathcal{L}_k(\theta) - \nabla \mathcal{L}_k(\theta')\| \leq \kappa \|\theta - \theta'\|. \quad (7)$$

264 We assume that the estimators of the gradients  $\nabla \mathcal{L}_k(\theta)$  have bounded variance.  
 265

$$266 \quad \mathbb{E} \|\nabla \hat{\mathcal{L}}_k(\theta) - \nabla \mathcal{L}_k(\theta)\|^2 \leq \sigma^2. \quad (8)$$

267 Together, these assumptions form a minimal yet sufficient foundation for analyzing score estimation  
 268 in practice. *Smoothness and bounded gradient variance* implied by the sub-Gaussian assumption  
 269 are mild and generally satisfied for standard neural architectures such as GELU activations. The *PL*

270 *condition* has been shown to emerge in over-parameterized networks or under lazy training regimes,  
 271 where the function class is expressive enough to approximate the ground-truth score function (Liu  
 272 et al., 2022). Notably, these conditions are not only specific to our setting they have been widely  
 273 adopted in recent works studying the optimization landscape of deep diffusion models (Salimans &  
 274 Ho, 2022; Liu et al., 2022). Note that in no prior works were such assumptions stated since they  
 275 assumed access to the empirical risk minimizer.

276 **Theorem 1** (Total Variation Distance Bound). *Let  $p_{t_0}$  denote the distribution obtained by the back-  
 277 ward process till time  $t_0$  starting from  $p_T$ , and  $\hat{p}_{t_k}(x)$  be the distribution generated by the backward  
 278 process at discretized time steps  $\{t_k\}$ , starting from  $\mathcal{N}(0, I)$  using the estimated score functions  
 279  $\hat{s}_{t_k}(x)$  where  $k \in [0, K]$ . Let  $d$  be the data dimension, and  $n_k$  be the number of samples for score  
 280 estimation at time step  $t_k$ .*

281 *Assume that the data distribution satisfies Assumption 1, the loss function  $\mathcal{L}_k(\theta)$  satisfies Assump-  
 282 tions 2, 3,4 for all  $k \in [0, K]$  and the learning rate for estimating  $\mathcal{L}_k(\theta)$  using SGD satisfies  
 283  $0 \leq \eta \leq \frac{1}{\kappa}$  for all  $k \in [0, K]$ . Further assume*

$$285 \quad n_k = \Omega \left( W^{2D} \cdot d^2 \cdot \log \left( \frac{4K}{\delta} \right) \left( \frac{\epsilon^{-4}}{\sigma_k^{-4}} \right) \right), \quad (9)$$

287 *Then, with probability at least  $1 - \delta$ , the total variation distance between the  $p_{t_0}$  and  $\hat{p}_{t_0}$  satisfies*

$$289 \quad TV(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}(\exp^{-T}) + \mathcal{O} \left( \frac{1}{\sqrt{K}} \right) + \mathcal{O} \left( \epsilon \cdot \sqrt{\left( T + \log \frac{1}{\kappa} \right)} \right) + \epsilon_{approx} \quad (10)$$

292 *Furthermore, by setting  $T = \Omega(\log(\frac{1}{\epsilon}))$ ,  $\kappa = \Omega(\epsilon)$  and  $K = \Omega(\epsilon^{-2})$ , we obtain*

$$293 \quad TV(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}(\epsilon) + \epsilon_{approx}, \quad (11)$$

295 *with probability at least  $1 - \delta$ .*

296 Theorem 1 establishes that the total variation distance between the true data distribution and the  
 297 diffusion model's output can be made arbitrarily small specifically,  $\tilde{\mathcal{O}}(\epsilon)$  by properly scaling model  
 298 capacity and algorithmic parameters. To the best of our knowledge, these are the only known sample  
 299 complexity bounds for score-based diffusion models, improving upon the prior results as discussed  
 300 in the introduction without assuming access to empirical risk minimizer for the score estimation  
 301 loss.

302 **Usage of  $p_{t_0}$  instead of  $p_0$  in Theorem 1:** We have shown that the estimated distribution  $\hat{p}_{t_0}$  is  $\mathcal{O}(\epsilon)$ -  
 303 close in total variation (TV) to  $p_{t_0}$ , where  $p_{t_0}$  denotes the data distribution  $p_0$  pushed forward by  $t_0$   
 304 steps of the forward process. We do not claim that  $\hat{p}_{t_0}$  is  $\mathcal{O}(\epsilon)$ -close in TV to the true data distribution  
 305  $p_0$  (i.e., we do not bound  $TV(p_0, \hat{p}_{t_0})$ ), because doing so would require additional assumptions on  
 306  $p_0$ . For example, Fu et al. (2024) (in Lemma D.5) assumes a sub-Gaussian data distribution to show  
 307 that  $TV(p_0, p_{t_0}) \leq \mathcal{O}(\sqrt{t_0} \log(1/t_0))$ . We also note that all other works listed in Table 1 similarly  
 308 provide upper bounds on  $TV(p_{t_0}, \hat{p}_{t_0})$ , not on  $TV(p_0, \hat{p}_{t_0})$ .

309 However, it is to be noted that using the sub-Gaussian assumption, our analysis can be extended to a  
 310 bound  $TV(p_0, \hat{p}_{t_0})$  via the triangle inequality:

$$312 \quad TV(p_0, \hat{p}_{t_0}) \leq TV(p_0, p_{t_0}) + TV(p_{t_0}, \hat{p}_{t_0}).$$

313 We formally present the data assumption and the resulting theorem as follows

314 **Assumption 5** (Sub-Gaussian Data Distribution.). *The data distribution  $p_0$  of the data variable  $x_0$   
 315 has an absolutely continuous CDF, is supported on a continuous set  $\Gamma \in \mathbb{R}^d$ , and there exists a  
 316 constant  $0 < C_2 < \infty$  such that for every  $t \geq 0$  we have  $P(|x_0| \geq t) \leq 2 \cdot \exp^{-\frac{t^2}{C_2^2}}$ .*

318 **Theorem 2** (Total Variation Distance Bound Under Sub-Gaussian Assumption). *Assume that the  
 319 data distribution satisfies Assumption 5, the loss function  $\mathcal{L}_k(\theta)$  satisfies Assumptions 2 3,4 for all  
 320  $k \in [0, K]$  and the learning rate satisfies for estimating  $\mathcal{L}_k(\theta)$  using SGD satisfies  $0 \leq \eta \leq \frac{1}{\kappa}$  for  
 321 all  $k \in [0, K]$ . Further assume*

$$322 \quad n_k = \Omega \left( W^{2D} \cdot d^2 \cdot \log \left( \frac{4K}{\delta} \right) \left( \frac{\epsilon^{-4}}{\sigma_k^{-4}} \right) \right), \quad (12)$$

324 Then, with probability at least  $1 - \delta$ , the total variation distance between the  $p_0$  and  $\hat{p}_{t_0}$  satisfies  
 325

$$\begin{aligned} 326 \quad TV(p_0, \hat{p}_{t_0}) &\leq \mathcal{O}(\sqrt{t_0} \log(1/t_0)) + \mathcal{O}(\exp^{-T}) + \mathcal{O}\left(\frac{1}{\sqrt{K}}\right) \\ 327 \quad &+ \mathcal{O}\left(\epsilon \cdot \sqrt{\left(T + \log \frac{1}{\kappa}\right)}\right) + \epsilon_{approx} \end{aligned} \quad (13)$$

328 Furthermore, by setting  $t_0 = \Omega(\epsilon^2)$ ,  $T = \Omega(\log(\frac{1}{\epsilon}))$ ,  $\kappa = \Omega(\epsilon)$  and  $K = \Omega(\epsilon^{-2})$ , we obtain  
 329

$$330 \quad TV(p_0, \hat{p}_{t_0}) \leq \mathcal{O}(\epsilon) + \epsilon_{approx}, \quad (14)$$

331 with probability at least  $1 - \delta$ .  
 332

### 333 Proof of Theorem 1.

334 Recall that  $\hat{p}_{t_0}$  is derived via score-based sampling, so using the triangle inequality repeatedly to  
 335 decompose the TV distance between the true distribution  $p_{t_0}$  and  $\hat{p}_{t_0}$  we obtain  
 336

$$342 \quad TV(p_{t_0}, \hat{p}_{t_0}) \leq TV(p_{t_0}, p_{t_0}^{dis}) + TV(p_{t_0}^{dis}, \tilde{p}_{t_0}) + TV(\tilde{p}_{t_0}, \hat{p}_{t_0}) \quad (15)$$

343 The bounds on  $TV(p_{t_0}, p_{t_0}^{dis})$  and  $TV(\tilde{p}_{t_0}, \hat{p}_{t_0})$  follow from Lemma B.4 of Gupta et al. (2024) and  
 344 Proposition 4 of Benton et al. (2024), respectively to get  
 345

$$348 \quad TV(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right) + TV(p_{t_0}^{dis}, \tilde{p}_{t_0}) + \mathcal{O}(\exp(-T)) \quad (16)$$

351 Note that we have used results from Gupta et al. (2024) and Benton et al. (2024) which assume a  
 352 bounded second moment for the data distribution. This is satisfied by Assumption 1. Now from  
 353 lemma 4,  $TV(p_{t_0}^{dis}, \tilde{p}_{t_0})$  is upper bounded as follows  
 354

$$356 \quad TV(p_{t_0}^{dis}, \tilde{p}_{t_0}) \leq \frac{1}{2} \sqrt{\sum_{k=0}^K E_{x \sim p_{t_k}} \|\hat{s}_{t_k}(x, t_k) - \nabla \log p_{t_k}(x)\|^2 (t_{k+1} - t_k)} \quad (17)$$

360 In order to upper bound  $TV(p_{t_0}^{dis}, \tilde{p}_{t_0})$ , we denote  $A(k)$  as  
 361

$$363 \quad A(k) := E_{x \sim p_{t_k}} \|\hat{s}_{t_k}(x, t_k) - \nabla \log p_{t_k}(x)\|^2 dx \quad (18)$$

366 Therefore, bounding the TV distance between  $p_{t_0}$  and  $\hat{p}_{t_0}$  translates to bounding the cumulative  
 367 error in estimating the score function at different time steps. We now focus on bounding this term.  
 368 Specifically, we have that  
 369

$$371 \quad TV(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right) + \frac{1}{2} \sqrt{\sum_{k=0}^K A_k \cdot (t_{k+1} - t_k)(t_{k+1} - t_k)} + \mathcal{O}(\exp(-T)) \quad (19)$$

375 Now, for each time step  $k$ , we decompose the total score estimation error, denoted by  $A(k)$ , into  
 376 three primary components: approximation error, statistical error, and optimization error. Each of  
 377 these error corresponds to a distinct aspect of learning the reverse-time score function in a diffusion  
 378 model as described below.

$$\begin{aligned}
& \mathbb{E}_{x \sim p_{t_k}} \left[ \|\hat{s}_{t_k}(x, t_k) - \nabla \log p_{t_k}(x)\|^2 \right] \leq 4 \underbrace{\mathbb{E}_{x \sim p_{t_k}(x)} \left[ \|s_{t_k}^a(x, t_k) - \nabla \log p_{t_k}(x)\|^2 \right]}_{\mathcal{E}_k^{\text{approx}}} \\
& + 4 \underbrace{\mathbb{E}_{x \sim p_{t_k}(x)} \left[ \|s_{t_k}^a(x, t_k) - s_{t_k}^b(x, t_k)\|^2 \right]}_{\mathcal{E}_k^{\text{stat}}} \\
& + 4 \underbrace{\mathbb{E}_{x \sim p_{t_k}(x)} \left[ \|\hat{s}_{t_k}(x, t_k) - s_{t_k}^b(x, t_k)\|^2 \right]}_{\mathcal{E}_k^{\text{opt}}}, \quad (20)
\end{aligned}$$

where, we define the parameters

$$\theta_k^a = \arg \min_{\theta \in \Theta} \mathbb{E}_{x \sim p_{t_k}} \left[ \|s_\theta(x, t_k) - \nabla \log p_t(x, t_k)\|^2 \right], \quad (21)$$

$$\theta_k^b = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \|s_\theta(x_i, t_k) - \nabla \log p_t(x_i, t_k)\|^2 \quad (22)$$

and denote  $s_{t_k}^a$  and  $s_{t_k}^b$  as the estimated score functions associated with the parameters  $\theta_k^a$  and  $\theta_k^b$  respectively. *Approximation error*  $\mathcal{E}_k^{\text{approx}}$  captures the error due to the limited expressiveness of the function class  $\{s_\theta\}_{\theta \in \Theta}$ . The *statistical error*  $\mathcal{E}_k^{\text{stat}}$  is the error from using a finite sample size. Finally, the *optimization error*  $\mathcal{E}_k^{\text{opt}}$  is due to not reaching the global minimum during training.

One of our key contributions lies in rigorously bounding each of these error components and showing how their interplay governs the overall generative error. In particular, we derive novel bounds that explicitly capture the dependencies on sample size, neural network capacity, and optimization parameters, without any assumption on the access to the empirical risk minimizer of the score estimation loss. We formalize these results in the following lemmas. Detailed proofs are deferred to Appendices D.1, and D.2, respectively.

**Lemma 1** (Approximation Error).  $\mathcal{E}_k^{\text{approx}}$  is defined as follows

$$\mathcal{E}_k^{\text{approx}} = \min_{\theta \in \Theta} \mathbb{E}_{x \sim p_{t_k}} \left[ \|s_\theta(x, t_k) - \nabla \log p_t(x, t_k)\|^2 \right] \quad (23)$$

Then, under Assumption 3 for all  $k \in [0, K]$ , we have

$$\mathcal{E}_k^{\text{approx}} \leq \epsilon_{\text{approx}} \quad (24)$$

This result directly follows from Assumption 3 and the definition of  $\mathcal{E}_k^{\text{approx}}$ .

**Lemma 2** (Statistical Error). Let  $n_k$  denote the number of samples used to estimate the score function at time step  $t_k$ . If the data distribution satisfies the Assumption 1 and the loss function  $\mathcal{L}_k(\theta)$  satisfies Assumptions 2 for all  $k \in [0, K]$ , then with probability at least  $1 - \delta$ , we have

$$\mathcal{E}_k^{\text{stat}} \leq \mathcal{O} \left( W^D \cdot d \cdot \sqrt{\frac{\log \left( \frac{2}{\delta} \right)}{n_k}} \right) \quad (25)$$

This is the component of the error that accounts for the fact that we have a finite sample size and thus we solve an empirical loss function given in equation 21. The proof of this lemma follows from utilizing the definitions of  $s_t^a$  and  $s_t^b$ . Existing analyses of statistical errors, such as those given in Shalev-Shwartz & Ben-David (2014), only work when the loss function is bounded. This is not the case for diffusion models. Thus, we use a novel analysis that uses the conditional normality of the score function as well as the bounded second moment property of the data variable in Assumption 1 to obtain the upper bound on the statistical error. The details of the analysis are given in Appendix D.1.

**Lemma 3** (Optimization Error). Let  $n_k$  be the number of samples used to estimate the score function at time step  $t_k$ . Assume that the score loss function  $\mathcal{L}_k(\theta)$  satisfies the Assumptions 2 and 4, for all

432  $k \in [0, K]$ , and the learning rate for estimating  $\mathcal{L}_k$  using SGD satisfies  $0 \leq \eta \leq \frac{1}{\kappa}$ , then with  
 433 probability at least  $1 - \delta$

$$434 \quad 435 \quad \mathcal{E}_k^{\text{opt}} \leq \mathcal{O} \left( W^D \cdot d \cdot \sqrt{\frac{\log \left( \frac{2}{\delta} \right)}{n_k}} \right). \quad 436 \quad 437 \quad (26)$$

438 This is the component of the error that accounts for the fact that we do not have access to the  
 439 empirical risk minimizer. We leverage assumptions 2, 4, alongside our unique recursive at each  
 440 stochastic gradient descent (SGD) step, which captures the error introduced by the finite number of  
 441 SGD steps in estimating the score function. This is the first analysis of diffusion models to explicitly  
 442 account for this error. All prior works assumed no such error, treating the empirical loss minimizer  
 443 as if it were known exactly. The details of the analysis are given in Appendix D.2.

444 Combining the decomposition in equation 20 along with Lemmas 1–3, we obtain the following  
 445 bound on  $A(k)$  (equation 18) with probability at least  $1 - \delta$

$$446 \quad 447 \quad A(k) \leq \mathcal{O} \left( W^D \cdot d \cdot \sqrt{\frac{\log \left( \frac{2}{\delta} \right)}{n_k}} \right) + \mathcal{O} \left( W^D \cdot d \cdot \sqrt{\frac{\log \left( \frac{2}{\delta} \right)}{n_k}} \right) + \epsilon_{\text{approx}} \quad 448 \quad 449 \quad (27)$$

$$450 \quad 451 \quad \leq \mathcal{O} \left( W^D \cdot d \cdot \sqrt{\frac{\log \left( \frac{2}{\delta} \right)}{n_k}} \right) + \epsilon_{\text{approx}}, \quad 452 \quad 453 \quad (28)$$

454 where in the second inequality we combine the first two terms appropriately. Setting the sample size

$$456 \quad 457 \quad n_k = \Omega \left( W^{2D} \cdot d^2 \cdot \log^2 \left( \frac{4K}{\delta} \right) \left( \frac{\epsilon^{-4}}{\sigma_k^{-4}} \right) \right), \quad 458 \quad (29)$$

459 we ensure that  $A(k) \leq \frac{\epsilon^2}{\sigma_k^2} = \frac{\epsilon^2}{1 - e^{-2(T-t_k)}}$  for all  $k \in \{0, \dots, K\}$ . Summing over all time steps, we  
 460 obtain with probability at least  $1 - \delta$

$$463 \quad 464 \quad \sum_{k=0}^K A(k)(t_{k+1} - t_k) \leq \sum_{k=0}^K \frac{\epsilon^2}{1 - e^{-2(T-t_k)}} (t_{k+1} - t_k) \quad 465 \quad 466 \quad (30)$$

$$467 \quad \leq \int_0^{T-\kappa} \frac{\epsilon^2}{1 - e^{-2(T-t)}} dt \leq \epsilon^2 \left( T + \log \frac{1}{\kappa} \right). \quad 468 \quad (31)$$

469 Note that the term  $(\log^2 \left( \frac{4K}{\delta} \right))$  appears in the upper bound for  $n_k$  in equation 29 since we have to  
 470 take a union bound for Lemma 2 and Lemma 3 and then take a union bound over  $K$  discretization  
 471 steps. Substituting this bound into equation 18, and then substituting the result into equation 16, we  
 472 obtain that with probability at least  $1 - \delta$ .

$$474 \quad 475 \quad \text{TV}(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}(\exp^{-T}) + \mathcal{O} \left( \frac{1}{\sqrt{K}} \right) + \mathcal{O} \left( \epsilon \cdot \sqrt{\left( T + \log \frac{1}{\kappa} \right)} \right) + \epsilon_{\text{approx}} \quad 476 \quad (32)$$

477 Finally, by choosing  $T = \Omega \left( \log \left( \frac{1}{\epsilon} \right) \right)$ ,  $\kappa = \Omega(\epsilon)$  and  $K = \Omega(\epsilon^{-2})$ , we conclude that with proba-  
 478 bility at least  $1 - \delta$

$$480 \quad \text{TV}(p_{t_0}, \hat{p}_{t_0}) \leq \mathcal{O}(\epsilon) + \epsilon_{\text{approx}}, \quad 481 \quad 482 \quad (33)$$

483 completing the proof of Theorem 1.  $\square$

484 In summary, our work provides a principled decomposition of the errors in score-based generative  
 485 models, highlighting how each component contributes to the overall sample complexity. This leads  
 486 to the first finite sample complexity bound of  $\tilde{\mathcal{O}}(\epsilon^{-4})$  for diffusion models without assuming access  
 487 to the empirical minimizer of the score estimation function.

486 4 CONCLUSION AND FUTURE WORK  
487488 In this work, we investigate the sample complexity of training diffusion models via score estimation  
489 using neural networks. We derive a sample complexity bound of  $\tilde{\mathcal{O}}(\epsilon^{-4})$ , which, to our knowledge,  
490 is the first such result that does not assume access to an empirical risk minimizer of the score esti-  
491 mation loss. Notably, our bound does not depend exponentially on the number of neural network  
492 parameters. For comparison, the best-known existing result achieves a bound of  $\tilde{\mathcal{O}}(\epsilon^{-5})$ , but it  
493 crucially assumes access to an ERM. All prior results establishing sample complexity bounds for  
494 diffusion models have made this assumption. Our contribution is the first to establish a sample com-  
495 plexity bound for diffusion models under the more realistic setting where exact access to empirical  
496 risk minimizer of the score estimation loss is not available.497 While our analysis focuses on unconditional distributions, extending these guarantees to conditional  
498 settings remains an important direction for future work.  
499500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

540 REFERENCES  
541

542 Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.

544 Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of  
545 natural images. *CVPR*, 2022.

547 Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas  
548 Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.

550 Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly  $\$d\$$ -linear  
551 convergence bounds for diffusion models via stochastic localization. In *The Twelfth International  
552 Conference on Learning Representations*, 2024.

554 Xin Bing, Xin He, and Chao Wang. Kernel ridge regression with predicted feature inputs and  
555 applications to factor-based nonparametric regression. *arXiv preprint arXiv:2505.20022*, 2025.

556 Adam Block, Youssef Mroueh, and Alexander Rakhlin. Generative modeling with denoising auto-  
557 encoders and langevin sampling. *arXiv preprint arXiv:2002.00107*, 2020.

559 Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. Introduction to statistical learning the-  
560 ory. In *Summer school on machine learning*, pp. 169–207. Springer, 2003.

561 Angela Castillo, Jonas Kohler, Juan C Pérez, Juan Pablo Pérez, Albert Pumarola, Bernard Ghanem,  
562 Pablo Arbeláez, and Ali Thabet. Adaptive guidance: Training-free acceleration of conditional  
563 diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39,  
564 pp. 1962–1970, 2025.

566 Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and  
567 distribution recovery of diffusion models on low-dimensional data. In *International Conference  
568 on Machine Learning*, pp. 4672–4712. PMLR, 2023.

570 Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models: Ap-  
571 plications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*,  
571 2024.

572 Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as  
573 learning the score: theory for diffusion models with minimal data assumptions. In *NeurIPS 2022  
574 Workshop on Score-Based Methods*, 2022.

576 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Ad-  
577 vances in Neural Information Processing Systems*, volume 34, pp. 8780–8794, 2021.

578 Zolnamar Dorjsembe, Hsing-Kuo Pao, Sodtavilan Odonchimed, and Furen Xiao. Conditional diffu-  
579 sion models for semantic 3d medical image synthesis. *imaging*, 19:20, 2023.

581 Zehao Dou, Minshuo Chen, Mengdi Wang, and Zhuoran Yang. Theory of consistency diffusion  
582 models: Distribution estimation meets fast sampling. In *Forty-first International Conference on  
583 Machine Learning*, 2024.

584 Salar Fattah, Richard Y. Zhang, and Somayeh Sojoudi. Linear-time algorithm for learning large-  
585 scale sparse graphical models. *IEEE Access*, 7:12658–12672, 2019. doi: 10.1109/ACCESS.2018.  
586 2890583.

588 Alessandro Favero, Antonio Sclocchi, and Matthieu Wyart. Bigger isn’t always memorizing: Early  
589 stopping overparameterized diffusion models. *arXiv preprint arXiv:2505.16959*, 2025.

590 Hengyu Fu, Zhuoran Yang, Mengdi Wang, and Minshuo Chen. Unveil conditional diffusion models  
591 with classifier-free guidance: A sharp statistical theory. *arXiv preprint arXiv:2403.11968*, 2024.

593 Zuyue Fu, Zhuoran Yang, and Zhaoran Wang. Single-timescale actor-critic provably finds globally  
594 optimal policy. In *International Conference on Learning Representations*, 2021.

594 Swetha Ganesh, Washim Uddin Mondal, and Vaneet Aggarwal. A sharper global convergence anal-  
 595 ysis for average reward reinforcement learning via an actor-critic approach. In *Forty-second*  
 596 *International Conference on Machine Learning*, 2025.

597 Mudit Gaur, Amrit Bedi, Di Wang, and Vaneet Aggarwal. Closing the gap: Achieving global conver-  
 598 gence (last iterate) of actor-critic under markovian sampling with neural network parametrization.  
 599 In *International Conference on Machine Learning*, pp. 15153–15179. PMLR, 2024.

600 Mudit Gaur, Utsav Singh, Amrit Singh Bedi, Raghu Pasupathu, and Vaneet Aggarwal. On the  
 601 sample complexity bounds in bilevel reinforcement learning. *arXiv preprint arXiv:2503.17644*,  
 602 2025.

603 Riccardo Grazzi, Massimiliano Pontil, and Saverio Salzo. Bilevel optimization with a lower-level  
 604 contraction: Optimal sample complexity without warm-start. *Journal of Machine Learning Re-*  
 605 *search*, 24(167):1–37, 2023.

606 Shivam Gupta, Aditya Parulekar, Eric Price, and Zhiyang Xun. Improved sample complexity bounds  
 607 for diffusion model training. In *The Thirty-eighth Annual Conference on Neural Information*  
 608 *Processing Systems*. PMLR, 2024.

609 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
 610 *arXiv:2207.12598*, 2022.

611 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
 612 *neural information processing systems*, 33:6840–6851, 2020.

613 Rongjie Huang, Max WY Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. Fast-  
 614 diff: A fast conditional diffusion model for high-quality speech synthesis. *arXiv preprint*  
 615 *arXiv:2204.09934*, 2022.

616 Zhihan Huang, Yuting Wei, and Yuxin Chen. Denoising diffusion probabilistic models are optimally  
 617 adaptive to unknown low dimensionality, 2024. URL <https://arxiv.org/abs/2410.18784>.

618 Yuling Jiao, Guohao Shen, Yuanyuan Lin, and Jian Huang. Deep nonparametric regression on  
 619 approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *The Annals of*  
 620 *Statistics*, 51(2):691–716, 2023.

621 Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi Jaakkola. Torsional dif-  
 622 fusion for molecular conformer generation. *Advances in neural information processing systems*,  
 623 35:24240–24253, 2022.

624 Harshat Kumar, Alec Koppel, and Alejandro Ribeiro. On the sample complexity of actor-critic  
 625 method for reinforcement learning with function approximation. *Machine Learning*, 112(7):  
 626 2433–2467, 2023.

627 Gen Li and Yuling Yan. Adapting to unknown low-dimensional structures in score-based diffusion  
 628 models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,  
 629 2024.

630 Gen Li, Yu Huang, Timofey Efimov, Yuting Wei, Yuejie Chi, and Yuxin Chen. Accelerating conver-  
 631 gence of score-based diffusion models, provably. In Ruslan Salakhutdinov, Zico Kolter, Katherine  
 632 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings*  
 633 *of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Ma-  
 634 chine Learning Research*, pp. 27942–27954, 21–27 Jul 2024a.

635 Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. A sharp convergence theory for the probability  
 636 flow odes of diffusion models. *arXiv preprint arXiv:2408.02320*, 2024b.

637 Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-  
 638 lm improves controllable text generation. *Advances in neural information processing systems*, 35:  
 639 4328–4343, 2022.

648 Jiadong Liang, Zhihan Huang, and Yuxin Chen. Low-dimensional adaptation of diffusion models:  
 649 Convergence in total variation. *arXiv preprint arXiv:2501.12982*, 2025a.  
 650

651 Yuchen Liang, Peizhong Ju, Yingbin Liang, and Ness Shroff. Broadening target distributions for  
 652 accelerated diffusion models via a novel analysis approach. In *The Thirteenth International Con-*  
 653 *ference on Learning Representations*, 2025b.

654 Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-  
 655 parameterized non-linear systems and neural networks. *Applied and Computational Harmonic*  
 656 *Analysis*, 59:85–116, 2022.

657

658 Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu,  
 659 Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with  
 660 semantic diffusion guidance. In *Proceedings of the IEEE/CVF winter conference on applications*  
 661 *of computer vision*, pp. 289–299, 2023.

662

663 Zhaoyang Lyu, Xudong Xu, Ceyuan Yang, Dahua Lin, and Bo Dai. Accelerating diffusion models  
 664 via early stop of the diffusion process. *arXiv preprint arXiv:2205.12524*, 2022.

665

666 The Tien Mai. Pac-bayesian risk bounds for fully connected deep neural network with gaussian  
 667 priors. *arXiv preprint arXiv:2505.04341*, 2025.

668

669 Aditya Malusare and Vaneet Aggarwal. Improving molecule generation and drug discovery with a  
 670 knowledge-enhanced generative model. *IEEE/ACM Transactions on Computational Biology and*  
 671 *Bioinformatics*, 2024.

672

673 Washim U Mondal and Vaneet Aggarwal. Improved sample complexity analysis of natural policy  
 674 gradient algorithm with general parameterization for infinite horizon discounted reward markov  
 675 decision processes. In *International Conference on Artificial Intelligence and Statistics*, pp. 3097–  
 676 3105. PMLR, 2024.

677

678 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
 679 In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.

680

681 Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distri-  
 682 bution estimators. In *International Conference on Machine Learning*, pp. 26517–26582. PMLR,  
 683 2023.

684

685 Bernt Øksendal. *Stochastic differential equations*. Springer, 2003.

686

687 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
 688 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-*  
 689 *ference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, 2022.

690

691 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
 692 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
 693 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*  
 694 *tion processing systems*, 35:36479–36494, 2022.

695

696 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In  
 697 *International Conference on Learning Representations*, 2022.

698

699 Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algo-*  
 700 *rithms*. Cambridge university press, 2014.

701

Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-  
 702 vised learning using nonequilibrium thermodynamics. In *International Conference on Machine*  
 703 *Learning*, pp. 2256–2265. PMLR, 2015.

704

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
 705 Poole. Score-based generative modeling through stochastic differential equations. In *Interna-*  
 706 *tional Conference on Learning Representations*, 2021.

702 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csd: Conditional score-based dif-  
703 fusion models for probabilistic time series imputation. *Advances in neural information processing*  
704 systems, 34:24804–24816, 2021.

705  
706 Nguyen Tran, Oleksii Abramenco, and Alexander Jung. On the sample complexity of graphical  
707 model selection from non-stationary samples. *IEEE Transactions on Signal Processing*, 68:17–  
708 32, 2019.

709 Anwaar Ulhaq and Naveed Akhtar. Efficient diffusion models for vision: A survey. *arXiv preprint*  
710 *arXiv:2210.09292*, 2022.

711 Andre Wibisono, Yihong Wu, and Kaylee Yingxi Yang. Optimal score estimation via empirical  
712 bayes smoothing. In *The Thirty Seventh Annual Conference on Learning Theory*, pp. 4958–4991.  
713 PMLR, 2024.

714  
715 Fangzhao Zhang and Mert Pilanci. Analyzing neural network-based generative diffusion models  
716 through convex optimization, 2024. URL <https://arxiv.org/abs/2402.01965>.

717  
718 Kaihong Zhang, Heqi Yin, Feng Liang, and Jingbo Liu. Minimax optimality of score-based diffusion  
719 models: Beyond the density lower bound assumptions. In *Forty-first International Conference on*  
720 *Machine Learning*, 2024.

721 Zhengbang Zhu, Hanye Zhao, Haoran He, Yichao Zhong, Shenyu Zhang, Haoquan Guo, Tingting  
722 Chen, and Weinan Zhang. Diffusion models for reinforcement learning: A survey. *arXiv preprint*  
723 *arXiv:2311.01223*, 2023.

724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756 A APPENDIX  
757758 B COMPARISON WITH PRIOR WORKS  
759760 In this section, we provide a detailed comparison of our results with prior work. Specifically, we  
761 analyze the sample complexity bounds presented in Gupta et al. (2024), and show how combining  
762 their results with those of Block et al. (2020) leads to an alternative bound.  
763764 B.1 SAMPLE COMPLEXITY OF GUPTA ET AL. (2024)  
765766 We begin by examining the sample complexity claim of  $\mathcal{O}(1/\epsilon^3)$  reported in Gupta et al. (2024).  
767 A closer analysis reveals that the actual sample complexity is  $\tilde{\mathcal{O}}(1/\epsilon^5)$ , once the error over all the  
768 discretization steps is properly accounted for and a union bound is applied.  
769770 The main result regarding the sample complexity of estimating the score function as given in Theorem  
771 C.2 of Gupta et al. (2024) is as follows. We re-iterate this Theorem here.  
772773 **Theorem C.2** Gupta et al. (2024). *Let  $q$  be a distribution of  $\mathbb{R}^d$  with second moment  $m_2^2$ . Let  
774  $\phi_\theta(\cdot)$  be the fully connected neural network with ReLU activations parameterized by  $\theta$ , with  $P$  total  
775 parameters and depth  $D$ . Let  $\Theta > 1$ . For any  $\gamma > 0$ , there exist  $K = \tilde{\mathcal{O}}\left(\frac{d}{\epsilon^2 + \delta^2} \log^2 \frac{m_2 + 1/m_2}{\gamma}\right)$   
776 discretization times  $0 = t_0 < \dots < t_K < T$  such that if for each  $t_k$ , there exists some score function  
777  $\hat{s}_\theta$  with  $\|\theta^*\|_\infty \leq \Theta$  such that*  
778

779 
$$\mathbb{E}_{x \sim p_{t_k}} \left[ \|s_\theta(x) - s_{t_k}(x)\|_2^2 \right] \leq \frac{\delta \cdot \epsilon^3}{CK^2 \sigma_{T-t_k}^2} \cdot \frac{1}{\log \frac{d+m_2+1/m_2}{\gamma}} \quad (34)$$
  
780

781 for sufficiently large constant  $C$ , then consider the score functions trained from  
782

783 
$$m > \tilde{\mathcal{O}}\left(\frac{K(d + \log \frac{1}{\delta}) \cdot PD}{\epsilon^3} \cdot \log\left(\frac{\max(m_2, 1) \cdot \Theta}{\delta}\right) \cdot \log\left(\frac{m_2 + 1/m_2}{\gamma}\right)\right). \quad (35)$$
  
784

785 *i.i.d. samples of  $q$ , with  $1 - \delta$  probability, DDPM can sample from a distribution  $\epsilon$ -close in TV to a  
786 distribution  $\gamma m_2$ -close in 2-Wasserstein to  $q$  in  $N$  steps.*  
787788 Note that Lemma B.6 of Gupta et al. (2024) states that  
789

790 
$$\text{TV}(p_{t_0}, \hat{p}_{t_0}) \leq \delta + \mathcal{O}\left(\frac{1}{\sqrt{N}}\right) + \epsilon \cdot \sqrt{T} + \mathcal{O}(\exp^{-T}) \quad (36)$$
  
791

792 Here  $p$  and  $\hat{p}$  are the true and learned data distributions, respectively.  
793794 In order to achieve  $\text{TV}(p, \hat{p}) \leq \epsilon$  we have to set  $\delta = \epsilon$ . This would imply  $N = \mathcal{O}(\epsilon^{-2})$  and  
795  $T = \mathcal{O}\left(\frac{1}{\epsilon}\right)$ . Putting this value of  $N$  in Equation equation 35 we obtain that for  
796

797 
$$m > \tilde{\mathcal{O}}\left(\frac{(d + \log(1/\delta)) \cdot PD}{\epsilon^5} \cdot \log\left(\frac{\Theta}{\epsilon}\right)\right). \quad (37)$$
  
798

800 we have with probability at least  $1 - \epsilon$   
801

802 
$$\text{TV}(p_{t_0}, \hat{p}_{t_0}) \leq \epsilon \quad (38)$$

803 This reveals a discrepancy between the reported sample complexity and the actual bound derived  
804 above, highlighting that the true complexity is significantly higher than what was originally reported.  
805806 In contrast, our analysis reduces the overall complexity by a factor of  $\tilde{\mathcal{O}}(1/\epsilon)$ , yielding the tightest  
807 known bounds for neural score estimation in diffusion models, i.e.,  $\tilde{\mathcal{O}}(1/\epsilon^4)$ . Further, unlike Gupta  
808 et al. (2024), our analysis avoids using the  $1 - \delta$ -quantile bound on the score norm and instead  
809 directly bounds the global  $L^2$  score estimation error thus avoids applying a union bound across time  
steps, and finally achieve tighter sample complexity guarantees.  
810

---

810    **C SCORE ESTIMATION ALGORITHM**  
811

812    In this section, we provide a detailed description of the algorithm used for estimating the score  
813    function in diffusion models.  
814

815    **Algorithm 1** Denoising Diffusion Probabilistic Model (DDPM)  
816

817    1: **Input:** Dataset  $\mathcal{D}$ , timesteps  $T$ , stop time  $t_{\text{stop}}$ , schedule  $\{\beta_t\}_{t=1}^T$ , network  $\epsilon_\theta$ , learning rate  $\eta$ ,  
818    iterations  $K$   
819    2: Precompute:  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$   
820       **Training (Score Estimation)**  
821    3: **for**  $i = 1$  to  $N$  **do**  
822    4:    Sample  $x_j \sim \mathcal{D}$ ,  $k \sim \text{Uniform}([1, T])$ ,  $\epsilon_k \sim \mathcal{N}(0, I)$  for  $i = 1, \dots, n$   
823    5:     $x_{t_i} = e^{-t_k} x_i + \sqrt{1 - e^{-2t_k}} \epsilon_i$   
824    6:    Compute loss:  $\hat{L}(\theta) = \|\epsilon_i - \epsilon_\theta(x_{t_i}, t_i)\|^2$   
825    7:    Update  $\theta \leftarrow \theta - \eta_k \cdot \nabla_\theta \hat{L}(\theta)$   
826    8: **end for**  
827       **Sampling**  
828    9: Sample  $x_T \sim \mathcal{N}(0, I)$   
829    10: **for**  $t = T$  down to  $t_{\text{stop}} + 1$  **do**  
830    11:     $z \sim \mathcal{N}(0, I)$  if  $t > 1$  else  $z = 0$   
831    12:     $\hat{\epsilon} = \epsilon_\theta(x_t, t)$   
832    13:     $\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\epsilon} \right)$   
833    14:     $x_{t-1} = \tilde{\mu}_t + \sqrt{\beta_t} \cdot z$   
834    15: **end for**  
835    16: **Return**  $x_{t_{\text{stop}}}$   
836  
837

838    **D PROOFS OF INTERMEDIATE LEMMAS**  
839

840    In this section, we present the proofs of intermediate lemmas used to bound the statistical error and  
841    optimization error in our analysis.  
842

843    **D.1 BOUNDING THE STATISTICAL ERROR**  
844

845    *Proof.* Let us define the population loss at time  $t_k$  for  $k \in [0, K]$  as  
846

847    
$$\mathcal{L}_k(\theta) = \mathbb{E}_{x \sim p_{t_k}} \|s_\theta(x, t_k) - \nabla \log p_{t_k}(x)\|^2, \quad (39)$$

848    where  $s_\theta$  denotes the score function estimated by a neural network parameterized by  $\theta$ , and  $x$  denotes  
849    samples at time  $t$  used in Algorithm 1. The corresponding empirical loss is defined as:  
850

851    
$$\hat{\mathcal{L}}_k(\theta) = \frac{1}{n} \sum_{i=1}^n \|s_\theta(x_i, t_k) - \nabla \log p_{t_k}(x_i)\|^2. \quad (40)$$

852    Let  $\theta_k^a$  and  $\theta_k^b$  be the minimizers of  $\mathcal{L}_k(\theta)$  and  $\hat{\mathcal{L}}_k(\theta)$ , respectively, corresponding to score functions  
853     $s_{t_k}^a$  and  $s_{t_k}^b$ . By the definitions of minimizers, we can write  
854

855    
$$\mathcal{L}_k(\theta_k^b) - \mathcal{L}_k(\theta_k^a) \leq \mathcal{L}_k(\theta_k^b) - \mathcal{L}_k(\theta_k^a) + \hat{\mathcal{L}}_k(\theta_k^a) - \hat{\mathcal{L}}_k(\theta_k^b) \quad (41)$$

856    
$$\leq \underbrace{|\mathcal{L}_k(\theta_k^b) - \hat{\mathcal{L}}_k(\theta_k^b)|}_{(I)} + \underbrace{|\mathcal{L}_k(\theta_k^a) - \hat{\mathcal{L}}_k(\theta_k^a)|}_{(II)}. \quad (42)$$

857    Note that the right-hand side of equation 41 is greater than the left-hand side since we  
858    have added the quantity  $\hat{\mathcal{L}}_k(\theta_k^a) - \hat{\mathcal{L}}_k(\theta_k^b)$  which is strictly positive since  $\theta_K^b$  is the minimizer of the  
859    function  $\hat{\mathcal{L}}_k(\theta)$  by definition. We then take the absolute value on both sides of the equation 42 to get  
860

We now bound terms (I) and (II) using generalization results. From Lemma 5 (Theorem 26.5 of Shalev-Shwartz & Ben-David (2014)), if the loss function  $\widehat{\mathcal{L}}(\theta)$  is uniformly bounded over the parameter space  $\Theta'' = \{\theta_k^a, \theta_k^b\}$ , then with probability at least  $1 - \delta$ , we have

$$|\mathcal{L}_k(\theta) - \widehat{\mathcal{L}}_k(\theta)| \leq \widehat{R}(\Theta'') + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right), \quad \forall \theta \in \Theta'' \quad (43)$$

where  $\widehat{R}(\Theta'')$  denotes the empirical Rademacher complexity of the function class restricted to  $\Theta''$ . Now since  $x$  is not bounded, this result does not hold. We then define the following two functions

$$\mathcal{L}'_k(\theta) = \mathbb{E}_{x \sim \mu_{t_k}} \|v_\theta(x, t_k) - v_{t_k}(x)\|^2, \quad (44)$$

and

$$\widehat{\mathcal{L}}'_k(\theta) = \frac{1}{n} \sum_{i=1}^n \|v_\theta(x_i, t_k) - v_{t_k}(x_i)\|^2. \quad (45)$$

where we define the functions

$$(v_t(x))_j = \begin{cases} (\nabla \log p_t(x))_j & \text{if } |\frac{x - e^{-t} x_0}{\sigma_t^2}|_j \leq \kappa \\ 0 & \text{if } |\frac{x - e^{-t} x_0}{\sigma_t^2}|_j \geq \kappa \end{cases} \quad (46)$$

and

$$(v_\theta(x, t))_j = \begin{cases} (s_\theta(x, t))_j & \text{if } |\frac{x - e^{-t} x_0}{\sigma_t^2}|_j \leq \kappa \\ 0 & \text{if } |\frac{x - e^{-t} x_0}{\sigma_t^2}|_j \geq \kappa \end{cases} \quad (47)$$

Here  $(v_t(x))_j, (\nabla \log p_t(x))_j, (v_\theta(x, t))_j$  and  $(s_\theta(x, t))_j$  denote the  $j^{th}$  co-ordinate of  $v_t(x)$ ,  $(\nabla \log p_t(x))$ ,  $v_\theta(x, t)$  and  $s_\theta(x, t)$  respectively. Further,  $|\frac{x - e^{-t} x_0}{\sigma_t^2}|_j$  denotes the  $j^{th}$  co-ordinate of the  $i^{th}$  sample of the score function in  $\widehat{\mathcal{L}}_k(\theta)$  which is given by  $\log p_t(x) = |\frac{x - e^{-t} x_0}{\sigma_t^2}|$ .

Note that the functions  $v_t(x)$  and  $v_\theta(x, t)$  are uniformly bounded. Thus using Theorem 26.5 of Shalev-Shwartz & Ben-David (2014) we have with probability at least  $1 - \delta$ ,

$$|\mathcal{L}'_k(\theta) - \widehat{\mathcal{L}}'_k(\theta)| \leq \widehat{R}(\theta) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right), \quad \forall \theta \in \Theta''. \quad (48)$$

Since  $\Theta'' = \{\theta_a, \theta_b\}$  is a finite class (just two functions). We can apply Lemma E.5 to bound the empirical Rademacher complexity  $\widehat{R}(\theta)$  in terms of the Rademacher complexity  $R(\theta)$  of the function class  $\Theta''$ . Since  $\widehat{R}(\theta) = \frac{1}{m} \mathbb{E}_\sigma [\max_{\theta \in \Theta''} \sum_{i=1}^n f(\theta) \sigma_i]$ , applying Lemma E.5, we have with probability at least  $1 - 2\delta$

$$|\mathcal{L}'_k(\theta) - \widehat{\mathcal{L}}'_k(\theta)| \leq \mathcal{O}\left(\frac{d \cdot W^D}{n}\right) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right), \quad \forall \theta \in \Theta''. \quad (49)$$

This yields that with probability at least  $1 - \delta$  we have

$$|\mathcal{L}'_k(\theta) - \widehat{\mathcal{L}}'_k(\theta)| \leq \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right), \quad \forall \theta \in \Theta'' \quad (50)$$

From this we have

$$|\mathcal{L}'_k(\theta) - \widehat{\mathcal{L}}_k(\theta)| \leq \mathcal{O}\left(d \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) \quad (51)$$

918 Now consider the probability of the event  
 919

920  
 921 
$$A_{i,j} = \left\{ \left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa \right\} \quad (52)$$
  
 922  
 923

924  
 925 Where  $\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j$  denotes the  $k^{th}$  co-ordinate of of the  $i^{th}$  sample of the score function given  
 926  
 927 by  $\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|$  We have the probability of this event upper bounded as  
 928  
 929

930  
 931 
$$P \left( \left| \frac{x_i - e^{-t}x_0}{\sigma_t^2} \right|_j \geq \kappa \right) = \mathbb{E}_z P \left( \left| \frac{x_i - e^{-t}x_0}{\sigma_t^2} \right|_j \geq \kappa \middle| x_0 \right) \quad (53)$$
  
 932  
 933

934 
$$\leq \exp(-\kappa^2(1 - e^{-t})) \quad (54)$$
  
 935

936 
$$\leq \exp(-\kappa^2) \quad (55)$$
  
 937

938 We get equation 54 from equation 53 since the score variable is conditionally normal given  $x_0$ .  
 939

940 Setting  $\kappa = \log(\frac{dn}{\delta})$ , we have  
 941

942  
 943 
$$P \left( \left| \frac{x_i - e^{-t}x_0}{\sigma_t^2} \right|_j \geq \kappa \right) \leq \frac{\delta}{dn} \quad (56)$$
  
 944

945 If we denote the event  $A = \{\hat{L}'(\theta) = \hat{L}(\theta)\}$ , then by union bound we have  $P(A) = P(\cup_{i,j} A_{i,j}) \leq$   
 946  $\sum_{i,j} P(A_{i,j}) \leq \delta$ . Let event  $B$  denote the failure of the generalization bound, i.e.,  
 947

948  
 949 
$$B := \left\{ \left| \mathcal{L}'_t(\theta) - \hat{\mathcal{L}}'_t(\theta) \right| > \hat{R}(\Theta'') + \mathcal{O} \left( \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right) \right\}. \quad (57)$$
  
 950  
 951

952 From above, we know  $\mathbb{P}(B) \leq \delta$  under the boundedness condition. Therefore, by the union bound,  
 953 we have  
 954

955 
$$\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B) \leq 2\delta, \quad (58)$$

956 
$$\implies \mathbb{P}(A^c \cap B^c) = 1 - P(A \cup B) \geq 1 - 2\delta. \quad (59)$$

957  
 958 On this event  $(A^c \cap B^c)$ , we have  $\hat{\mathcal{L}}'(\theta) = \hat{\mathcal{L}}(\theta)$ . Hence, with probability at least  $1 - 2\delta$ , we have  
 959

960  
 961 
$$|\mathcal{L}_k(\theta_k^b) - \mathcal{L}_k(\theta_k^a)| \leq \left| \mathcal{L}_k(\theta_k^b) - \hat{\mathcal{L}}_t(\theta_k^b) \right| + \left| \mathcal{L}_k(\theta_k^a) - \hat{\mathcal{L}}_t(\theta_k^a) \right|. \quad (60)$$
  
 962

963 
$$\leq \left| \mathcal{L}_k(\theta_k^b) - \hat{\mathcal{L}}'_t(\theta_k^b) \right| + \left| \mathcal{L}_k(\theta_k^a) - \hat{\mathcal{L}}'_t(\theta_k^a) \right|. \quad (61)$$
  
 964

965 
$$= \left| \mathcal{L}_k(\theta_k^b) - \mathcal{L}'_t(\theta_k^b) \right| + \left| \mathcal{L}_k(\theta_k^a) - \mathcal{L}'_t(\theta_k^a) \right|.$$
  
 966  
 967 
$$+ \left| \mathcal{L}'_t(\theta_k^b) - \hat{\mathcal{L}}'_t(\theta_k^b) \right| + \left| \mathcal{L}_k(\theta_k^a) - \hat{\mathcal{L}}'_t(\theta_k^a) \right|. \quad (62)$$
  
 968

969  
 970 
$$\leq |\mathcal{L}'_t(\theta_k^a) - \mathcal{L}_t(\theta_k^a)| + |\mathcal{L}'_t(\theta_k^b) - \mathcal{L}_t(\theta_k^b)| + \mathcal{O} \left( d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right) \quad (63)$$
  
 971

In order to bound  $|\mathcal{L}_k(\theta) - \mathcal{L}'_t(\theta)|$  we have the following

$$|\mathcal{L}_k(\theta) - \mathcal{L}'_t(\theta)| \leq \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} |(s_\theta(x, t_k)) - (\nabla \log p_{t_k}(x))|_j^2 - \mathbb{E}_{x \sim u_t} |(v_t(x))_k - (v_\theta(x, t))|_j^2 \quad (64)$$

$$\leq \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( |(\nabla \log p_{t_k}) - (s_\theta(x, t_k))|_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)_i}{\sigma_t^2} \right|_i \geq \kappa} \right) \quad (65)$$

$$\leq \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})} \left( \left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} - (s_\theta(x, t))_k \right|_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j \geq \kappa} \right) \quad (66)$$

$$\leq 2 \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( \left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j \geq \kappa} \right)$$

$$+ \sum_{j=1}^d 2\mathbb{E}_{x_j \sim (u_{t_k})_j} \left( (s_\theta(x, t))_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right| \geq \kappa} \right) \quad (67)$$

$$\leq 2 \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( \left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j \geq \kappa} \right)$$

$$+ \sum_{j=1}^d C_{\Phi''} \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( \left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right| \geq \kappa} \right) \quad (68)$$

$$\leq \left( \frac{4 + 2C_{\Phi''}}{\sigma_t^2} \right) \sum_{k=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \underbrace{\left( |x_{ij}^2 \mathbf{1}|_{\frac{x - e^{-t} (x_0)}{\sigma_t^2}} \right)_j}_{\geq \kappa}$$

$$+ \underbrace{\frac{2}{\sigma_t^2} \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( |x_0|_j^2 \mathbf{1}_{\left| \frac{x - e^{-t}(x_0)}{\sigma_t^2} \right|_j \geq \kappa} \right)}_I \quad (69)$$

We get Equation 67 from Equation 66 by using the identity  $(a - b)^2 \leq 2|a|^2 + 2|b|^2$ . We get Equation 68 from Equation 67 by using Lemma 8. We get Equation 69 from Equation 68 by using the identity  $(a - b)^2 \leq 2|a|^2 + 2|b|^2$  again. Now

1026 we separately obtain upper bounds for the terms  $I$  and  $II$  as follows.  
 1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

$$\sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( |x_j|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \quad (70)$$

$$= \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( |x_j|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \quad (71)$$

$$\leq \sum_{j=1}^d \mathbb{E}_{x_0} \mathbb{E}_{x_j \sim (u_{t_k})_j | x_0} \left( |x_j|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \quad (72)$$

$$\leq \sum_{j=1}^d \mathbb{E}_{x_0} \mathbb{E}_{x_j \sim (u_t)_k | x_0} \left( |x_j|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) P \left( \left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right| \geq \kappa \middle| x_0 \right) \quad (73)$$

$$\leq \exp(-\kappa^2) \sum_{j=1}^d \mathbb{E}_{x_0} \mathbb{E}_{x_0 \sim (u_t)_j | x_0} \left( |x_k|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \quad (74)$$

$$\leq \exp(-\kappa^2) \sum_{k=1}^d \mathbb{E}_{x_0} \left( \sigma_t^2 + \sigma_t^2 \cdot \kappa \cdot \sigma_t^2 \cdot \frac{\phi(\kappa \cdot \sigma_t^2)}{1 - \Phi((\kappa \cdot \sigma_t^2))} \right) \quad (75)$$

$$\leq \exp(-\kappa^2) \sum_{k=1}^d \mathbb{E}_z (2 \cdot \sigma_t^2) \quad (76)$$

$$\leq \mathcal{O}(\exp(-\kappa^2)) \quad (77)$$

We get Equation equation 75 from Equation equation 74 by using Lemma 7. We get Equation equation 77 from Equation equation 76 from Assumption 1 and by using the upper bound on the Mill's ration which implies that  $\frac{\phi(\kappa)}{1 - \Phi(\kappa)} \leq \kappa + \frac{1}{\kappa}$ . We get Equation equation 77 from Equation equation 76 from Assumption 1, which implies that the second moment of  $z$  is bounded.

$$\begin{aligned}
& \sum_{j=1}^d \mathbb{E}_{x_j \sim (u_{t_k})_j} \left( |x_0|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \tag{78}
\end{aligned}$$

$$\begin{aligned}
& = \sum_{j=1}^d \mathbb{E}_{x_k \sim (u_t)_k} \left( |x_0|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \tag{79}
\end{aligned}$$

$$\begin{aligned}
& \leq \sum_{j=1}^d \mathbb{E}_{x_0} \mathbb{E}_{x_k \sim (u_t)_k | x_0} \left( |x_0|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) \tag{80}
\end{aligned}$$

$$\begin{aligned}
& \leq \sum_{j=1}^d \mathbb{E}_{x_0} \mathbb{E}_{x_k \sim (u_t)_k | x_0} \left( |x_0|^2 \mathbf{1}_{\left| \frac{x_i - e^{-t}(x_0)_i}{\sigma_t^2} \right|_j \geq \kappa} \right) P \left( \left| \frac{x_k - t z_k}{1-t} \right|_j \geq \kappa | x_0 \right) \tag{81}
\end{aligned}$$

$$\begin{aligned}
& \leq \exp(-\kappa^2) \sum_{j=1}^d \mathbb{E}_{x_0} |x_0|^2 \tag{82}
\end{aligned}$$

$$\begin{aligned}
& \leq \mathcal{O}(\exp(-\kappa^2)) \tag{83}
\end{aligned}$$

Setting  $\kappa = \log \frac{dn}{\delta}$  Plugging Equation equation 77, equation 83 into Equation equation 69. Then we have

$$|\mathcal{L}_k(\theta) - \mathcal{L}'_t(\theta)| \leq \mathcal{O}(\exp(-\kappa^2)), \quad \forall \theta = \{\theta_k^a, \theta_k^b\} \tag{84}$$

$$\leq \mathcal{O}\left(\frac{\delta}{dn}\right) \tag{85}$$

Now plugging Equation equation 84 into Equation equation 63 we get with probability at least  $1 - 2\delta$

$$|\mathcal{L}_k(\theta_k^b) - \mathcal{L}_k(\theta_k^a)| \leq \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) \tag{86}$$

Finally, using the Polyak-Łojasiewicz (PL) condition for  $\mathcal{L}_k(\theta)$ , from Assumption 2, we have from the quadratic growth condition of PL functions the following,

$$\|\theta_k^a - \theta_k^b\|^2 \leq \mu |\mathcal{L}_k(\theta_k^a) - \mathcal{L}_k(\theta_k^b)|, \tag{87}$$

and applying Lipschitz continuity of the velocity fields with respect to parameter  $x$

$$\|v^{\theta_k^a}(x, t_k) - v^{\theta_k^b}(x, t_k)\|^2 \leq L_t \cdot \|\theta_k^a - \theta_k^b\|^2 \tag{88}$$

$$\leq L_t \cdot \mu |\mathcal{L}_k(\theta_k^a) - \mathcal{L}_k(\theta_k^b)| \tag{89}$$

$$\leq \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right). \tag{90}$$

1134 Here  $L_t$  is the Lipschitz parameter of the neural networks. It is always possible to obtain this  
 1135 Lipschitz constant as the quantity  $\|v^a(x, t) - v^b(x, t)\|^2 \leq L_t$  is non-zero only over a finite domain  
 1136 of  $x$ . Taking expectation with respect to  $x$ , we obtain the following.

$$1138 \mathbb{E}_{x \sim u_{t_k}} \|s^{\theta_k^a}(x, t_k) - s^{\theta_k^b}(x, t_k)\|^2 \leq 2\mathbb{E}_{x \sim u_t} \|s^{\theta_k^a}(x, t_k) - s^{\theta_k^b}(x, t_k) - v_t^a(x) + v_t^a(x)\|^2 \\ 1139 + 2\mathbb{E}_{x \sim u_t} \|v_t^a(x) - v_t^b(x)\|^2 \quad (91)$$

$$1140 \leq 4\mathbb{E}_{x \sim u_t} \|v_t^a(x) - s^{\theta_k^a}(x, t_k)\|^2 \quad (92)$$

$$1141 + 4\mathbb{E}_{x \sim u_t} \|s^{\theta_k^b}(x, t_k) - v_t^b(x)\|^2 \\ 1142 + 4\mathbb{E}_{x \sim u_t} \|v_t^a(x) - v_t^b(x)\|^2 \quad (93)$$

$$1143 \leq \mathcal{O}(\kappa^{-2}) + \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right). \quad (94)$$

$$1144 \leq \mathcal{O}\left(\frac{\delta}{dn}\right) + \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right). \quad (95)$$

$$1145 \leq \mathcal{O}\left(d \cdot W^D \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}}\right). \quad (96)$$

1155 This completes the proof. Note that the quantities  $4\mathbb{E}_{x \sim u_t} \|u_t^a(x) - v_t^a(x)\|^2$  and  $4\mathbb{E}_{x \sim u_t} \|u_t^a(x) - v_t^b(x)\|^2$  are bounded in the same manner as is done in Equation equation 83.

□

## 1159 D.2 BOUNDING OPTIMIZATION ERROR

1160 The optimization error ( $\mathcal{E}_{\text{opt}}$ ) accounts for the fact that gradient-based optimization does not necessarily find the optimal parameters due to limited steps, local minima, or suboptimal learning rates. This can be bounded as follows.

1165 *Proof.* Let  $\mathcal{E}_k^{\text{opt}}$  denote the optimization error incurred when performing stochastic gradient descent (SGD), with the empirical loss defined by

$$1166 \mathcal{L}_k^i(\theta) = \|s_\theta(x_i, t_k) - \nabla \log p_t(x_i)\|^2. \quad (97)$$

1169 The corresponding population loss is

$$1170 \mathcal{L}_k(\theta) = \mathbb{E}_{x \sim p_{t_k}} \left[ \|s_\theta(x, t_k) - \nabla \log p_{t_k}(x)\|^2 \right], \quad (98)$$

1174 Thus,  $\mathcal{E}_t^{\text{opt}}$  captures the error incurred during the stochastic optimization at each fixed time step  $t$ .  
 1175 We now derive upper bounds on this error.

1176 From the smoothness of  $\mathcal{L}_k(\theta)$  through Assumption 4, we have

$$1177 \mathcal{L}_k(\theta_{i+1}) \leq \mathcal{L}_k(\theta_i) + \langle \nabla \mathcal{L}_k(\theta_i), \theta_{i+1} - \theta_i \rangle + \frac{\kappa}{2} \|\theta_{i+1} - \theta_i\|^2. \quad (99)$$

1180 Taking conditional expectation given  $\theta_i$ , and using the unbiasedness of the stochastic gradient  
 1181  $\nabla \widehat{\mathcal{L}}_k(\theta_i)$ , we get:

$$1182 \mathbb{E}[\mathcal{L}_k(\theta_{i+1}) \mid \theta_i] \leq \mathcal{L}_k(\theta_i) - \alpha_t \|\nabla \mathcal{L}_k(\theta_i)\|^2 + \frac{\kappa \alpha_t^2}{2} \mathbb{E}[\|\nabla \widehat{\mathcal{L}}_k(\theta_i)\|^2 \mid \theta_i]. \quad (100)$$

1186 Now using the variance bound on the stochastic gradients using Assumption 4, we have

$$1187 \mathbb{E}[\|\nabla \widehat{\mathcal{L}}_k(\theta_i)\|^2 \mid \theta_i] \leq \|\nabla \mathcal{L}_k(\theta_i)\|^2 + \sigma^2, \quad (101)$$

1188 Using this in the previous equation, we have that  
 1189

$$\mathbb{E}[\mathcal{L}_k(\theta_{t+1}) | \theta_i] \leq \mathcal{L}_k(\theta_i) - \eta \|\nabla \mathcal{L}(\theta_i)\|^2 + \frac{\kappa \eta^2}{2} (\|\nabla \mathcal{L}(\theta_i)\|^2 + \sigma^2) \quad (102)$$

$$= \mathcal{L}(\theta_i) - \left( \eta - \frac{\kappa \eta^2}{2} \right) \|\nabla \mathcal{L}(\theta_i)\|^2 + \frac{\kappa \eta^2 \sigma^2}{2}. \quad (103)$$

1194 Now applying the PL inequality (Assumption 2),  $\|\nabla L(\theta_i)\|^2 \geq 2\mu(L(\theta_i) - L^*)$ , we substitute in  
 1195 the above inequality to get

$$\mathbb{E}[\mathcal{L}(\theta_{t+1}) | \theta_i] - \mathcal{L}^* \leq \left( 1 - 2\mu \left( \eta - \frac{\kappa \eta^2}{2} \right) \right) (\mathcal{L}(\theta_i) - \mathcal{L}^*) + \frac{\kappa \eta^2 \sigma^2}{2}. \quad (104)$$

1199 Define the contraction factor  
 1200

$$\rho = 1 - 2\mu \left( \eta - \frac{\kappa \eta^2}{2} \right). \quad (105)$$

1203 Taking total expectation and defining  $\delta_t = \mathbb{E}[L(\theta_i) - L^*]$ , we get the recursion:  
 1204

$$\delta_{t+1} \leq \rho \cdot \delta_t + \frac{\kappa \eta^2 \sigma^2}{2}. \quad (106)$$

1207 When  $\eta \leq \frac{1}{\kappa}$ , we have

$$\eta - \frac{\kappa \eta^2}{2} \geq \frac{\eta}{2} \Rightarrow \rho \leq 1 - \mu \eta. \quad (107)$$

1211 Unrolling the recursion we have

$$\delta_t \leq (1 - \mu \eta)^t \delta_0 + \frac{\kappa \eta^2 \sigma^2}{2} \sum_{j=0}^{t-1} (1 - \mu \eta)^j. \quad (108)$$

1215 Using the geometric series bound:

$$\sum_{j=0}^{t-1} (1 - \mu \eta)^j \leq \frac{1}{\mu \eta}, \quad (109)$$

1219 we conclude that

$$\delta_t \leq (1 - \mu \eta)^t \delta_0 + \frac{\kappa \eta \sigma^2}{2\mu}. \quad (110)$$

1223 Hence, we have the convergence result

$$\mathbb{E}[\mathcal{L}_k(\theta_n) - \mathcal{L}^*] \leq (1 - \mu \eta)^n \delta_0 + \frac{\kappa \eta \sigma^2}{2\mu}. \quad (111)$$

$$\leq \exp(-\eta \cdot \mu \cdot n) \delta_0 + \frac{\kappa \eta \sigma^2}{2\mu} \quad (112)$$

$$\leq \mathcal{O}\left(\frac{1}{n}\right) \quad (113)$$

1232 We get Equation equation 112 from Equation equation 111 by the identity  $(1 - x) \leq e^{-x}$ . We get  
 1233 Equation equation 113 from Equation equation 112 by setting the step size  $\eta = \mathcal{O}(\frac{1}{n})$ .  
 1234

1235 Note that  $\hat{s}_{t_k}$  and  $\hat{\theta}_k$  denote our estimate of the loss function and associated parameter obtained  
 1236 from the SGD. Also note that  $\mathcal{L}^*$  is the loss function corresponding whose minimizer is the neu-  
 1237 ral network  $s_{t_k}^a$  and the neural parameter  $\theta_k^a$  is our estimated score parameter. Thus applying the  
 1238 quadratic growth inequality.

$$\|\hat{s}_{t_k}(x, t_k) - s_{t_k}^a(x, t_k)\|^2 \leq L \|\hat{\theta}_k - \theta_k^a\|^2 \leq \|\mathcal{L}(\theta_k) - \mathcal{L}^*\| \quad (114)$$

$$\leq \mathcal{O}\left(\frac{1}{n}\right) \quad (115)$$

1242 From lemma 2 we have with probability  $1 - \delta$  that  
 1243

$$\|s_{t_k}^a(x, t_k) - s_{t_k}^b(x, t_k)\|^2 \leq L \|\theta_t^a - \theta_t^b\|^2 \quad (116)$$

$$\leq L \cdot \mu |\mathcal{L}_k(\theta_t^a) - \mathcal{L}_k(\theta_t^b)| \quad (117)$$

$$\leq \mathcal{O} \left( d \cdot W^D \cdot \sqrt{\frac{\log \frac{2}{\delta}}{n}} \right). \quad (118)$$

1250 Thus we have with probability at least  $1 - \delta$   
 1251

$$\|s_t(x, t_k) - s_t^b(x, t_k)\|^2 \leq 2\|s_{t_k}(x, t_k) - s_{t_k}^a(x, t_k)\| + 2\|s_{t_k}^a(x, t_k) - s_{t_k}^b(x, t_k)\| \quad (119)$$

$$\leq \mathcal{O} \left( \log \left( \frac{1}{n} \right) \right) + \mathcal{O} \left( d \cdot \sqrt{\frac{\log \frac{2}{\delta}}{n}} \right). \quad (120)$$

$$\leq \mathcal{O} \left( d \cdot W^D \cdot \sqrt{\frac{\log \frac{2}{\delta}}{n}} \right). \quad (121)$$

$$(122)$$

1262 Taking expectation with respect to  $x \sim p_{t_k}$  on both sides completes the proof.  
 1263

□

## 1266 E INTERMEDIATE LEMMAS

1268 **Lemma 4** (TV bound via Girsanov for reverse diffusions). *Let  $X$  and  $\tilde{X}$  on  $[0, T]$  solve*

$$1270 dX_t = (f(X_t, t) - \sigma^2(t)s_*(X_t, t))dt + \sigma(t)d\bar{W}_t, \quad d\tilde{X}_t = (f(\tilde{X}_t, t) - \sigma^2(t)s_\theta(\tilde{X}_t, t))dt + \sigma(t)d\bar{W}_t,$$

1271 with the same nondegenerate diffusion  $\sigma(t) \in \mathbb{R}^{d \times d}$  (invertible for a.e.  $t$ ) and the same initial law  
 1272 at time  $T$ . Let  $\mathbb{P}$  and  $\mathbb{Q}$  be the path measures of  $X$  and  $\tilde{X}$  on  $C([0, T], \mathbb{R}^d)$ . Assume Novikov's  
 1273 condition

$$1274 \mathbb{E}_{\mathbb{Q}} \exp \left( \frac{1}{2} \int_0^T \|\sigma(t)(s_\theta(\tilde{X}_t, t) - s_*(\tilde{X}_t, t))\|_2^2 dt \right) < \infty.$$

1275 Then

$$1276 \text{TV}(\mathbb{P}, \mathbb{Q}) \leq \frac{1}{2} \left( \mathbb{E}_{\mathbb{Q}} \int_0^T \|\sigma(t)(s_\theta(\tilde{X}_t, t) - s_*(\tilde{X}_t, t))\|_2^2 dt \right)^{1/2}.$$

1281 *Proof.* Write the drift difference as

$$1282 \Delta b(x, t) = -\sigma^2(t)(s_*(x, t) - s_\theta(x, t)).$$

1283 By Girsanov's theorem (under the stated Novikov condition),  $\mathbb{P} \ll \mathbb{Q}$  and the Radon–Nikodym  
 1284 derivative is the exponential martingale driven by  $u_t = \sigma(t)^{-1}\Delta b(\tilde{X}_t, t) = \sigma(t)(s_\theta(\tilde{X}_t, t) -$   
 1285  $s_*(\tilde{X}_t, t))$ . The Cameron–Martin formula yields

$$1287 \text{KL}(\mathbb{P} \parallel \mathbb{Q}) = \frac{1}{2} \mathbb{E}_{\mathbb{Q}} \left[ \int_0^T \|u_t\|_2^2 dt \right] = \frac{1}{2} \mathbb{E}_{\mathbb{Q}} \left[ \int_0^T \|\sigma(t)(s_\theta(\tilde{X}_t, t) - s_*(\tilde{X}_t, t))\|_2^2 dt \right].$$

1288 Applying Pinsker's inequality  $\text{TV}(\mathbb{P}, \mathbb{Q}) \leq \sqrt{\text{KL}(\mathbb{P} \parallel \mathbb{Q})/2}$  gives

$$1289 \text{TV}(\mathbb{P}, \mathbb{Q}) \leq \frac{1}{2} \left( \mathbb{E}_{\mathbb{Q}} \int_0^T \|\sigma(t)(s_\theta - s_*)\|_2^2 dt \right)^{1/2}.$$

1290 Finally, the evaluation map  $C([0, T], \mathbb{R}^d) \rightarrow \mathbb{R}^d, \omega \mapsto \omega(0)$ , is measurable, so by data processing  
 1291 for  $f$ -divergences,  $\text{TV}(\mathcal{L}(X_0), \mathcal{L}(\tilde{X}_0)) \leq \text{TV}(\mathbb{P}, \mathbb{Q})$ . □

1296 Let  $\{x_k\}_{k=0}^N$  and  $\{\tilde{x}_k\}_{k=0}^N$  be Euler schemes with the same Gaussian noises,

1298  $x_{k-1} = x_k + (f_k - \sigma_k^2 s_*(x_k, t_k)) \Delta t_k + \sigma_k \sqrt{\Delta t_k} \xi_k, \quad \tilde{x}_{k-1} = \tilde{x}_k + (f_k - \sigma_k^2 s_\theta(\tilde{x}_k, t_k)) \Delta t_k + \sigma_k \sqrt{\Delta t_k} \xi_k,$   
1299  
1300  $\xi_k \sim \mathcal{N}(0, I)$  i.i.d. Then, with “traj” denoting trajectory measures,

1301  
1302  $\text{KL}(\text{traj}_* \| \text{traj}_\theta) = \frac{1}{2} \sum_{k=1}^N \mathbb{E} \left[ \left\| \sigma_k(s_\theta(x_k, t_k) - s_*(x_k, t_k)) \right\|_2^2 \Delta t_k \right],$   
1303

1304 and hence by Pinsker’s inequality we get,

1305  
1306  $\text{TV}(\text{traj}_*, \text{traj}_\theta) \leq \frac{1}{2} \left( \sum_{k=1}^N \mathbb{E} \left\| \sigma_k(s_\theta - s_*) \right\|_2^2 \Delta t_k \right)^{1/2}.$   
1307  
1308

1309 **Lemma 5** (Theorem 26.5 of Shalev-Shwartz & Ben-David (2014)). *Consider data  $z \in Z$ , the*  
1310 *parametrized hypothesis class  $h_\theta, \theta \in \Theta$ , and the loss function  $\ell(h, z) : \mathbb{R}^d \rightarrow \mathbb{R}$ , where  $|\ell(h, z)| \leq$*   
1311 *c. We also define the following terms*

1313  
1314  $L_D(h) = \mathbb{E} \ell(h, z) \tag{123}$   
1315  
1316

1317  
1318  $L_S(h) = \frac{1}{m} \sum_{z_i \in \mathcal{S}} \ell(h, z_i) \tag{124}$   
1319  
1320

1321 which denote the expected and empirical loss functions respectively.

1322 Then,

1323 With probability of at least  $1 - \delta$ , for all  $h \in \mathcal{H}$ ,

1324  
1325  $L_D(h) - L_S(h) \leq 2R(\ell \circ \Theta \circ S) + 4c \sqrt{\frac{2 \ln(4/\delta)}{m}}. \tag{125}$   
1326  
1327

1328 where  $2R(\ell \circ \Theta \circ S)$  denotes the empirical Radamacher complexity over the loss function  $\ell$ , hypothesis parameter set  $\Theta$  and the dataset  $\mathcal{S}$

1329 **Lemma 6** (Extewnson of Massart’s Lemma Bousquet et al. (2003)). *Let  $\Theta''$  be a finite function*  
1330 *class. Then, for any  $\theta \in \Theta''$ , we have*

1331  
1332  
1333  $\mathbb{E}_\sigma \left[ \max_{\theta \in \Theta''} \sum_{i=1}^n f(\theta) \sigma_i \right] \leq \|f(\theta)\|_2 \leq (BW)^L \left( d + \frac{L}{W} \right) \tag{126}$   
1334  
1335  
1336

1337 where  $\sigma_i$  are i.i.d random variables such that  $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = \frac{1}{2}$ . We get the second  
1338 inequality by denoting  $L$  as the number of layers in the neural network,  $W$  and  $B$  a constant such  
1339 all parameters of the neural network upper bounded by  $B$ .

1340 *Proof.* Let  $h_0 = x$ , and for  $\ell = 0, \dots, L-1$  define the layer recursion

1341  
1342  $h_{\ell+1} = \sigma(W_\ell h_\ell + b_\ell),$   
1343

1344 where  $W_\ell \in \mathbb{R}^{n_{\ell+1} \times n_\ell}$ ,  $b_\ell \in \mathbb{R}^{n_{\ell+1}}$ , and  $n_\ell \leq W$  for hidden layers. We work with the  $\ell_\infty$  operator  
1345 norm:

1346  
1347  $\|W_\ell\|_\infty = \max_i \sum_j |(W_\ell)_{ij}| \leq B n_\ell \leq BW = \alpha.$

1348 Since  $\sigma$  is 1-Lipschitz with  $\sigma(0) = 0$ , we have  $\|\sigma(u)\|_\infty \leq \|u\|_\infty$  and thus

1349  
 $\|h_{\ell+1}\|_\infty \leq \|W_\ell\|_\infty \|h_\ell\|_\infty + \|b_\ell\|_\infty \leq \alpha \|h_\ell\|_\infty + B.$

1350 With  $\|h_0\|_\infty \leq d$ , iterating this affine recursion yields the standard geometric-series bound  
 1351

$$1352 \quad \|h_L\|_\infty \leq \alpha^L d + B \sum_{i=0}^{L-1} \alpha^i = \alpha^L d + B \frac{\alpha^L - 1}{\alpha - 1} \quad (\alpha \neq 1),$$

1355 and for  $\alpha = 1$ ,  $\|h_L\|_\infty \leq d + BL$ . The scalar output  $f(x)$  is either a coordinate of  $h_L$  or obtained  
 1356 by applying the same 1-Lipschitz activation to a linear form of  $h_L$ ; in either case,  $|f(x)| \leq \|h_L\|_\infty$ ,  
 1357 giving the stated bound.

1358 For the  $\alpha \geq 1$  simplification, use  $\sum_{i=0}^{L-1} \alpha^i \leq L\alpha^{L-1}$  to obtain  
 1359

$$1360 \quad |f(x)| \leq \alpha^L d + BL\alpha^{L-1} = (BW)^L \left( d + \frac{L}{W} \right).$$

1362 For  $\alpha < 1$ , since  $\alpha^i \leq 1$ ,  $\sum_{i=0}^{L-1} \alpha^i \leq L$  and hence  $|f(x)| \leq d + BL$ . Finally, substituting  
 1363  $W = S/L$  gives the size-based form  
 1364

$$1365 \quad |f(x)| \leq (BS/L)^L \left( d + \frac{L^2}{S} \right).$$

□

1368  
 1369 **Lemma 7** (Second Moment of a Symmetrically Truncated Normal). *Let  $X \sim \mathcal{N}(\mu, \sigma^2)$ , and let  $a >$   
 1370  $0$ . Then the second moment of  $X$  conditioned on being outside the symmetric interval  $[\mu - a, \mu + a]$   
 1371 is given by*

$$1372 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \mu^2 + \sigma^2 + \sigma a \cdot \frac{\phi\left(\frac{a}{\sigma}\right)}{1 - \Phi\left(\frac{a}{\sigma}\right)},$$

1375 where  $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$  is the standard normal probability density function (PDF), and  $\Phi(z)$  is  
 1376 the standard normal cumulative distribution function (CDF).

1377 *Proof.* Let  $X \sim \mathcal{N}(\mu, \sigma^2)$ . We aim to compute the second moment of  $X$  conditioned on the event  
 1378 that it lies outside an interval centered at its mean

$$1381 \quad \mathbb{E}[X^2 \mid |X - \mu| > a]$$

1383 This represents the expected squared value of  $X$ , given that  $X$  is in the tails of the distribution (i.e.,  
 1384 more than  $a$  units away from the mean).

1385 By definition, the conditional expectation is  
 1386

$$1388 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \frac{\mathbb{E}[X^2 \cdot \mathbf{1}_{\{|X - \mu| > a\}}]}{\mathbb{P}(|X - \mu| > a)}$$

1390 The numerator integrates  $X^2$  over the tail regions  $(-\infty, \mu - a) \cup (\mu + a, \infty)$ , while the denominator  
 1391 is the probability mass in those same regions.

1392 To simplify the integrals, we standardize  $X$ . Define the standard normal variable  
 1393

$$1395 \quad Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1) \quad \Rightarrow \quad X = \mu + \sigma Z$$

1397 Define  $\alpha = \frac{a}{\sigma}$ . Then  
 1398

$$1399 \quad |X - \mu| > a \quad \Leftrightarrow \quad |Z| > \alpha$$

1401 Our conditional second moment becomes  
 1402

$$1403 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \mathbb{E}[(\mu + \sigma Z)^2 \mid |Z| > \alpha]$$

1404 Expanding the square inside the expectation  
 1405

$$1406 \quad (\mu + \sigma Z)^2 = \mu^2 + 2\mu\sigma Z + \sigma^2 Z^2$$

1408 Taking the conditional expectation  
 1409

$$1411 \quad \mathbb{E}[(\mu + \sigma Z)^2 \mid |Z| > \alpha] = \mu^2 + 2\mu\sigma\mathbb{E}[Z \mid |Z| > \alpha] + \sigma^2\mathbb{E}[Z^2 \mid |Z| > \alpha]$$

1412 Since the standard normal distribution is symmetric and the region  $|Z| > \alpha$  is also symmetric, we  
 1413 have  
 1414

$$1416 \quad \mathbb{E}[Z \mid |Z| > \alpha] = 0$$

1418 Thus, the expression simplifies to  
 1419

$$1420 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \mu^2 + \sigma^2\mathbb{E}[Z^2 \mid |Z| > \alpha]$$

1422 By definition  
 1423

$$1424 \quad \mathbb{E}[Z^2 \mid |Z| > \alpha] = \frac{\int_{|z|>\alpha} z^2\phi(z) dz}{\mathbb{P}(|Z| > \alpha)} = \frac{2\int_{\alpha}^{\infty} z^2\phi(z) dz}{2(1 - \Phi(\alpha))} = \frac{\int_{\alpha}^{\infty} z^2\phi(z) dz}{1 - \Phi(\alpha)}$$

1427 Using Intergration by Parts we get,  
 1428

$$1429 \quad \int_{\alpha}^{\infty} z^2\phi(z) dz = \phi(\alpha)\alpha + 1 - \Phi(\alpha)$$

1432 Therefore  
 1433

$$1434 \quad \mathbb{E}[Z^2 \mid |Z| > \alpha] = \frac{\phi(\alpha)\alpha + 1 - \Phi(\alpha)}{1 - \Phi(\alpha)} = 1 + \frac{\alpha\phi(\alpha)}{1 - \Phi(\alpha)}$$

1437 Substitute back into the expression for  $\mathbb{E}[X^2 \mid |X - \mu| > a]$   
 1438

$$1439 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \mu^2 + \sigma^2 \left( 1 + \frac{\alpha\phi(\alpha)}{1 - \Phi(\alpha)} \right)$$

1442 Recall that  $\alpha = \frac{a}{\sigma}$ , so the final expression becomes  
 1443

$$1444 \quad \mathbb{E}[X^2 \mid |X - \mu| > a] = \mu^2 + \sigma^2 + \sigma a \cdot \frac{\phi\left(\frac{a}{\sigma}\right)}{1 - \Phi\left(\frac{a}{\sigma}\right)}$$

1447  $\square$   
 1448

1449 **Lemma 8** (Linear Growth of Finite Neural Networks). *Let  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}$  be the output of a feedfor-  
 1450 ward neural network with a finite number of layers and parameters and  $\theta \in \Theta$  where  $\Theta$  has a finite  
 1451 number of elements. Suppose that each activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  satisfies the growth condition*

$$1453 \quad |\sigma(z)| \leq A + B|z|, \quad \text{for all } z \in \mathbb{R},$$

1454 for constants  $A, B \geq 0$ . Then there exists a constant  $C_{\Theta} > 0$  such that for all  $x \in \mathbb{R}^d$ ,

$$1455 \quad |f(x)| \leq C_{\Theta}(1 + \|x\|).$$

1457 *Proof.* We proceed by induction on the number of layers in the network.

1458 **Base case: One-layer network.** Let the network be a single-layer function  
 1459

$$1460 \quad 1461 \quad 1462 \quad f(x) = \sum_{i=1}^k a_i \sigma(w_i^\top x + b_i),$$

1463 where  $w_i \in \mathbb{R}^d$ ,  $b_i \in \mathbb{R}$ , and  $a_i \in \mathbb{R}$ . Then  
 1464

$$1465 \quad 1466 \quad 1467 \quad |f(x)| \leq \sum_{i=1}^k |a_i| \cdot |\sigma(w_i^\top x + b_i)|.$$

1468 Using the growth condition on  $\sigma$ , we get  
 1469

$$1470 \quad |\sigma(w_i^\top x + b_i)| \leq A + B|w_i^\top x + b_i| \leq A + B(\|w_i\|\|x\| + |b_i|).$$

1471 Hence

$$1472 \quad 1473 \quad 1474 \quad |f(x)| \leq \sum_{i=1}^k |a_i| (A + B(\|w_i\|\|x\| + |b_i|)) = C_0 + C_1\|x\|,$$

1475 where  $C_0, C_1$  are constants depending only on the network parameters. Therefore  
 1476

$$1477 \quad |f(x)| \leq C(1 + \|x\|) \quad \text{with } C = \max\{C_0, C_1\}.$$

1478 **Inductive step.** Assume the result holds for all networks with  $L$  layers, i.e., for any such network  
 1479  $f_L(x)$ ,

$$1480 \quad |f_L(x)| \leq C_L(1 + \|x\|).$$

1481 Now consider a network with  $L + 1$  layers, defined by  
 1482

$$1483 \quad 1484 \quad 1485 \quad f_{L+1}(x) = \sum_{j=1}^k a_j \sigma(f_L^{(j)}(x)),$$

1486 where each  $f_L^{(j)}(x)$  is an output of a depth- $L$  subnetwork. By the inductive hypothesis  
 1487

$$1488 \quad |f_L^{(j)}(x)| \leq C_j(1 + \|x\|).$$

1489 Applying the activation bound  
 1490

$$1491 \quad 1492 \quad |\sigma(f_L^{(j)}(x))| \leq A + B|f_L^{(j)}(x)| \leq A + BC_j(1 + \|x\|).$$

1493 Then

$$1494 \quad 1495 \quad 1496 \quad |f_{L+1}(x)| \leq \sum_{j=1}^k |a_j| \cdot |\sigma(f_L^{(j)}(x))| \leq \sum_{j=1}^k |a_j| (A + BC_j(1 + \|x\|)) = C_{L+1}(1 + \|x\|),$$

1497 for some constant  $C_{L+1} > 0$ . This completes the induction.  
 1498

## 1499 EXAMPLES OF VALID ACTIVATION FUNCTIONS

1500 The condition  $|\sigma(z)| \leq A + B|z|$  holds for most common activations  
 1501

- 1502 • **ReLU:**  $\sigma(z) = \max(0, z) \Rightarrow |\sigma(z)| \leq |z|$
- 1503 • **Leaky ReLU:** bounded by linear function of  $|z|$
- 1504 • **Tanh:** bounded by 1  $\Rightarrow A = 1, B = 0$
- 1505 • **Sigmoid:** bounded by 1

1506  
 1507  
 1508  
 1509  
 1510  
 1511

□