# STATISTICAL FOUNDATIONS OF CONDITIONAL DIFFU SION TRANSFORMERS

Anonymous authors

003 004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

## ABSTRACT

We explore the statistical foundations of conditional diffusion transformers (DiTs) with classifier-free guidance. Through a comprehensive analysis of "in-context" conditional DiTs under four data assumptions, we demonstrate that both conditional DiTs and their latent variants achieve minimax optimality for unconditional DiTs. By discretizing input domains into infinitesimal grids and performing term-by-term Taylor expansions on the conditional score function, we enable leveraging transformers' universal approximation capabilities through detailed piecewise constant approximations, resulting in tighter bounds. Extending our analysis to the latent setting under a linear latent subspace assumption, we show that latent conditional DiTs achieve lower bounds than their counterparts in both approximation and estimation. We also establish the minimax optimality of latent unconditional DiTs. Our findings provide statistical limits for conditional and unconditional DiTs and offer practical guidance for developing more efficient and accurate models.

## 1 INTRODUCTION

We investigate the approximation and estimation rates of conditional diffusion transformers (DiTs) with classifier-free guidance. Specifically, we derive score approximation, score estimation, and 025 distribution estimation guarantees for both conditional DiTs and their latent variants under various 026 data conditions. We also demonstrate that both conditional DiTs and their latent variants lead to 027 the minimax optimality of unconditional DiTs under identified settings. This analysis is not only 028 practical but also timely. Transformer-based conditional diffusion models are leading advancements 029 in generative AI due to their success as scalable and flexible frameworks for image (Wu et al., 2024) and video generation (Saharia et al., 2022). But our knowledge of the theory behind conditional DiTs 031 is still limited. While Hu et al. (2024b) analyze approximation and estimation rates using transformer universality, their results are not tight and only focus on unconditional diffusion. Meanwhile, existing 033 theoretical studies on conditional diffusion models have primarily examined ReLU networks (Fu et al., 2024a), model-free settings (Ye et al., 2024), or generative sampling processes (Dinh et al., 034 2023), without addressing transformer architectures. This work fills the gap by examining the statistical boundaries of conditional DiTs.

In this work, we provide a thorough analysis of conditional DiT and its latent variant under four standard data assumptions and establish their minimax optimality through tight distribution estimation bounds. Our approach employs two key techniques: discretizing input domains into infinitesimal grids and performing term-by-term Taylor expansions of the conditional diffusion score function under Hölder smoothness assumptions, motivated by the local diffused polynomial analysis (Fu et al., 2024a; Oko et al., 2023). These methods leverage the regularity of the score function, enabling efficient use of transformers' universal approximation capabilities through detailed piecewise approximations. Consequently, we achieve tighter bounds. We summarize the theoretical results in Table 1.

044 045 046

047

048

049

## 2 BACKGROUNDS AND PRELIMINARIES

**Conditional Diffusion Model.** The forward process adds noise to data  $x_0$  given condition y, resulting in a noisy distribution  $P_t(x_t|y) \sim N(\alpha_t x_0, \sigma_t^2 I_{d_x})$ . The backward process reverses this using the score function  $\nabla \log p_t(\cdot|y)$ .

**Classifier-Free Guidance.** This method approximates conditional and unconditional score functions using a neural network  $s_W$ . The loss function is:

$$\ell(x_0, y; s_W) = \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_t \sim N(\alpha_t x_0, \sigma_t^2 I_{d_x})} \left[ \|s_W(x_t, \tau y, t) - \nabla_{x_t} \log \phi_t(x_t | x_0) \|_2^2 \right] \mathrm{d}t.$$

054	Table 1: Summary of Theoretical Results. The initial data is $d_x$ -dimensional, and the condition is $d_y$
055	dimensional. For latent DiT, the latent variable is $d_0$ -dimensional. $\sigma_t^2 = 1 - e^{-t}$ is the denoising scheduler. The
056	sample size is n, and $0 < \epsilon < 1$ represents the score approximation error. While we report asymptotics for large
057	$d_x, d_0$ , we reintroduce the <i>n</i> dependence in the estimation results to emphasize sample complexity convergence

058 059	Assumption	Score Approximation	Score Estimation	Dist. Estimation (Total Variation)	Minimax Optimality
060	Generic Hölder Smooth	$\mathcal{O}\Big( \left( \log \left( \frac{1}{\epsilon}  ight)  ight)^{d_x} / \sigma_t^4 \Big)$	$   n^{-o(1/d_x)} \cdot (\log n)^{\mathcal{O}(d_x)} $	$ \mid n^{-o(1/d_x)} \cdot (\log n)^{\mathcal{O}(d_x)} $	×
062	Stronger Hölder Smooth	$\left(\log\left(\frac{1}{\epsilon}\right)\right)^{\mathcal{O}(1)}/\sigma_t^2$	$ \qquad n^{-o(1)} \cdot (\log n)^{\mathcal{O}(1)} $	$ \qquad \qquad n^{-o(1)} \cdot (\log n)^{\mathcal{O}(1)} $	<ul> <li>✓</li> </ul>
063	Latent Subspace + Generic Hölder Smooth	$\mathcal{O}\Big( \left(\log\left(rac{1}{\epsilon} ight) ight)^{d_0}/\sigma_t^4 \Big)$	$\left  n^{-o(1/d_0)} \cdot (\log n)^{\mathcal{O}(d_0)} \right $	$ \qquad \qquad$	×
065	Latent Subspace + Stronger Hölder Smooth	$\left(\log\left(\frac{1}{\epsilon}\right)\right)^{\mathcal{O}(1)}/\sigma_t^2$	$ \qquad \qquad n^{-o(1)} \cdot \left(\log n\right)^{\mathcal{O}(1)} $	$ \qquad \qquad n^{-o(1)} \cdot \left(\log n\right)^{\mathcal{O}(1)} $	~



Figure 1: Conditional DiT Network Architecture. The architecture includes a reshape layer R, its reverse <sup>1</sup>, and embedding layers for label y and timestep t. The model concatenates the embeddings with input  $R^{-}$ sequences and processes them through a transformer network  $f_{\mathcal{T}}$ .

where  $\tau$  denotes the conditional or unconditional version. The empirical loss is  $\widehat{\mathcal{L}}(s_W) =$  $\frac{1}{n}\sum_{i=1}^{n}\ell(x_{0,i}, y_i; s_W).$ 

**Conditional Diffusion Transformer Networks.** We use a transformer network as a score estimator  $s_W$ , following notation from (Hu et al., 2024b). The transformer block consists of self-attention and feed-forward layers. The self-attention layer is defined as:

$$f^{(SA)}(Z) = Z + \sum_{i=1}^{h} W_{O}^{i}(W_{V}^{i}Z) \operatorname{Softmax}\left[(W_{K}^{i}Z)^{\top}(W_{Q}^{i}Z)\right],$$
(2.1)

where  $W_V^i, W_K^i, W_Q^i \in \mathbb{R}^{s \times d}$  and  $W_Q^i \in \mathbb{R}^{d \times s}$  are weight matrices. The feed-forward layer is:

$$f^{(\text{FF})}(Z) = Z + W_2 \text{ReLU}(W_1 Z + b_1) + b_2, \qquad (2.2)$$

where  $W^{(1)} \in \mathbb{R}^{r \times d}$ ,  $W^{(2)} \in \mathbb{R}^{d \times r}$ ,  $b^{(1)} \in \mathbb{R}^r$ , and  $b^{(2)} \in \mathbb{R}^d$  are weights and biases.

Definition 2.1 (Transformer Block and Network Function Class). We define a transformer block of h-head, s-hidden dimension, r-feedforward dimension, with positional encoding  $E \in \mathbb{R}^{d \times L}$  as the function:

$$f^{h,s,r}\left(Z\right) \coloneqq f^{(\mathrm{FF})}\left(f^{(\mathrm{SA})}\left(Z+E\right)\right) : \mathbb{R}^{d \times L} \mapsto \mathbb{R}^{d \times L}.$$

The transformer network function class  $\mathcal{T}^{h,s,r}$  consists of all functions that are compositions of one or more such transformer blocks. Formally,

$$\mathcal{T}^{h,s,r} \coloneqq \left\{ \tau : \mathbb{R}^{d \times L} \mapsto \mathbb{R}^{d \times L} \mid \tau = f^{h,s,r} \circ \cdots \circ f^{h,s,r} \right\}$$

**Conditional Diffusion Transformer (DiT).** We consider a transformer network f in the class  $\mathcal{T}^{h,s,r}$ , and we take an input data point (x, y, t) in  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \times [t_0, T]$ . We adopt the "in-context conditioning" approach for conditional DiT networks as described in (Peebles & Xie, 2023) and shown in Figure 1. We reshape a vector input  $x \in \mathbb{R}^{d_x}$  into a sequential matrix input format  $Z \in \mathbb{R}^{d \times L}$ , where  $d_x = d \cdot L$ . **Definition 2.2** (DiT Reshape Layer  $R(\cdot)$ ). Let  $R(\cdot) : \mathbb{R}^{d_x} \to \mathbb{R}^{d \times L}$  reshape a  $d_x$ -dimensional input into a  $d \times L$  matrix. For an image input with  $d_x = i \times i$ , it transforms the input into a sequence representation where feature dimension  $d = p^2$  and sequence length  $L = (i/p)^2$ . The reverse reshape (flatten) layer is defined as  $R^{-1}(\cdot) : \mathbb{R}^{d \times L} \to \mathbb{R}^{d_x}$ .

## **3** STATISTICAL LIMITS OF CONDITIONAL DITS

We first introduce the definition of Hölder space and Hölder ball following (Fu et al., 2024b).

**Definition 3.1** (Hölder Space). Let  $\alpha \in \mathbb{Z}^d_+$  and  $\beta = k_1 + \gamma$  with  $k_1 = \lfloor \beta \rfloor, \gamma \in [0, 1)$ . The Hölder space  $\mathcal{H}^{\beta}(\mathbb{R}^d)$  consists of all  $\alpha$ -differentiable functions  $f : \mathbb{R}^d \to \mathbb{R}$  with finite Hölder norm:

$$f\|_{\mathcal{H}^{\beta}(\mathbb{R}^{d})} \coloneqq \max_{\|\alpha\|_{1} \leq k_{1}} \sup_{x} |\partial^{\alpha} f(x)| + \max_{\|\alpha\|_{1} = k_{1}} \sup_{x \neq x'} \frac{|\partial^{\alpha} f(x) - \partial^{\alpha} f(x')|}{\|x - x'\|_{\infty}^{\gamma}}.$$

The Hölder ball of radius B is define as  $\mathcal{H}^{\beta}(\mathbb{R}^d, B) \coloneqq \{f : \|f\|_{\mathcal{H}^{\beta}(\mathbb{R}^d)} < B\}.$ 

Let  $x_0 \in \mathbb{R}^{d_x}$  denote the initial data, and  $y \in [0, 1]^{d_y}$  the conditional label. With Definition 3.1, we state the generic and stronger Hölder assumption on the conditional distribution of initial data  $x_0$ .

Assumption 3.1 (Hölder Smooth Data). The conditional density function  $p_0(x_0|y)$  is defined on the domain  $\mathbb{R}^{d_x} \times [0,1]^{d_y}$  and belongs to Hölder ball of radius B > 0 for Hölder index  $\beta > 0$ , denoted by  $p_0(x_0|y) \in \mathcal{H}^{\beta}(\mathbb{R}^{d_x} \times [0,1]^{d_y}, B)$ . We consider two cases:

- (Generic) For any  $y \in [0,1]^{d_y}$ , there exist positive constants  $C_1, C_2$  such that  $p_0(x_0|y) \leq C_1 \exp\left(-C_2 \|x_0\|_2^2/2\right)$ .
- (Stronger) Given a constant radius B, positive constants C and  $C_2$ , we assume  $p(x_0|y) = \exp\left(-C_2||x_0||_2^2/2\right) \cdot f(x_0, y)$  where  $f \in \mathcal{H}^{\beta}(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$  and  $f(x_0, y) \geq C$  for all  $(x_0, y) \in \mathbb{R}^{d_x} \times [0, 1]^{d_y}$ .

We state our main result of score approximation using transformers under Assumption 3.1 as follows:

**Theorem 3.1** (Conditional Score Approximation under Assumption 3.1). For any precision parameter  $0 < \epsilon < 1$  and smoothness parameter  $\beta > 0$ , let  $\epsilon \leq \mathcal{O}(N^{-\beta})$  for some  $N \in \mathbb{N}$ . For some positive constants  $C_{\alpha}, C_{\sigma} > 0$ , for any  $y \in [0, 1]^{d_y}$  and  $t \in [N^{-C_{\sigma}}, C_{\alpha} \log N]$ , there exists a  $\mathcal{T}_{score}(x, y, t) \in \mathcal{T}_{B}^{h,s,r}$  such that:

$$\int_{\mathbb{R}^{d_x}} \|\mathcal{T}_{\text{score}}(x, y, t) - \nabla \log p_t(x|y)\|_2^2 p_t(x|y) \, \mathrm{d}x = \mathcal{O}\left(\frac{B^2}{\sigma_t^{\zeta}} \cdot N^{-\omega} \cdot (\log N)^{\phi}\right),$$

where the parameters  $\zeta$ ,  $\omega$ , and  $\phi$  are defined as follows:

> - (Generic)  $\zeta = 4$ ,  $\omega = \frac{\beta}{d_x + d_y}$ , and  $\phi = d_x + \frac{\beta}{2} + 1$ . - (Stronger)  $\zeta = 2$ ,  $\omega = \frac{2\beta}{d_x + d_y}$ , and  $\phi = \beta + 1$ .

Building on our approximation results from Theorem 3.1, next we evaluate the performance of the score estimator  $\hat{s}$  trained with finite samples by optimizing the empirical loss. To quantify this, we introduce the notion of score estimation risk and characterize its upper bound.

**Definition 3.2** (Conditional Score Risk). Given a score estimator  $\hat{s}$ , we define the risk as:

$$\mathcal{R}(\hat{s}) := \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_t, y} \left[ \| \hat{s}(x_t, y, t) - \nabla \log p_t(x_t | y) \|_2^2 \right] \mathrm{d}t.$$

**Theorem 3.2** (Conditional Score Estimation with Transformer). Consider  $y \in [0, 1]^{d_y}$  and  $t \in [t_0, T]$  with  $t_0 = N^{-C_{\sigma}}$  and  $T = C_{\alpha} \log N$ , where  $C_{\sigma}, C_{\alpha}$  are positive constants such that  $t_0 < 1$  holds.

• Assume 
$$d_x = \Omega\left(\sqrt{\frac{\log N}{\log \log N}}\right)$$
 and generic Assumption 3.1. By taking  $N = n^{\frac{d_x+d_y}{d_x+d_y+\beta}}$ , it holds  

$$\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\mathcal{R}(\hat{s})\right] = \mathcal{O}\left(\frac{1}{t_0}n^{-\frac{\min\left(\beta,(1-\nu_1)\left(d_x+d_y\right)-3\beta\right)}{\left(d_x+d_y+\beta\right)}}\left(\log n\right)^{\nu_2+2}\right).$$



Figure 2: Network Architecture of Latent Conditional DiT. The overall architecture consists of linear layer of encoder and decoder, reshaping layer  $\tilde{R}(\cdot)$  and  $\tilde{R}^{-1}(\cdot)$ , embedding layer for label y and timestep t. The embedding concatenates with input sequences and processes by the adapted transformer network.

where  $\nu_1 = \frac{68\beta}{(d_x + d_y)} + 104C_{\sigma}$  and  $\nu_2 = 12d_x + 12\beta + 2$ .

• Under stronger Assumption 3.1. For all  $x \in \mathbb{R}^{d_x}$ , by taking  $N = n^{\frac{d_x + d_y}{d_x + d_y + 2\beta}}$ , it holds  $\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \mathcal{R}(\hat{s}) \right] = \mathcal{O}\left( \log \frac{1}{t_0} n^{-\frac{\min\left(2\beta, (1-\nu_3)(d_x + d_y) - 2\beta\right)}{(d_x + d_y + 2\beta)}} (\log n)^{\max(12,\beta+1)} \right).$ where  $\nu_3 = 4(12\beta d_x + 31\beta d + 6\beta)/d(d_x + d_y) + 12C_{\alpha}(12d_x + 25d + 6)/d + 72C_{\sigma}.$ 

Theorem 3.2 provides a straightforward basis for deriving the distribution estimation theorem presented in Table 1. Furthermore, we show the minimax optimality of the unconditional DiT architecture under stronger Assumption 3.1. Specifically, we obtain the distribution estimation error of unconditional DiTs by removing the condition y and let  $d_y = 0$ . With the condition  $d_x = o\left(\sqrt{\log n/\log \log n}\right)$ , then the distribution estimation error becomes  $\widetilde{\mathcal{O}}(n^{-\frac{\min(\beta,(1-\nu_3)(d_x+d_y)/2-\beta)}{d_x+2\beta}})$ . Unconditional DiT is the minimax optimal distribution estimator under  $(1-\nu_3)(d_x+d_y)/2-\beta > \beta$ .

## 4 LATENT CONDITIONAL DITS

This section builds on Section 3 by exploring latent conditional DiTs. We consider raw data  $x \in \mathbb{R}^{d_x}$ residing in a low-dimensional subspace under Assumption 4.1, represented by latent variables  $h \in \mathbb{R}^{d_0}$  with  $d_0 \leq d_x$ . Adapting the approach from Peebles & Xie (2023), we employ a transformer network to approximate score functions on these latents (see Figure 2). The network includes a reshape layer converting vector inputs h into matrix form  $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$ , with reshaping operations  $\tilde{R}$ and its inverse, under constraints  $d_0 \leq d_x$ ,  $\tilde{d} \leq d$ , and  $\tilde{L} \leq L$ . Linear transformations  $W_U^{\top}$  and  $W_U$ encode raw data x into latents h such that x = Uh, satisfying the conditions of Assumption 4.1.

Assumption 4.1 (Low-Dimensional Linear Latent Space). The data x can be represented through a latent variable  $h \in \mathbb{R}^{d_0}$  such that x = Uh, where  $U \in \mathbb{R}^{d_x \times d_0}$  is a matrix with orthonormal columns. The latent variable h follows a distribution  $P_h$  characterized by the density function  $p_h$ .

The approximation and estimation results closely follows Theorem 3.1, with differences highlighted in low-dimensional data subspace assumption and Hölder smoothness on latent representation. We arrive the results by replacing the input dimension d, L to  $\tilde{d}$  and  $\tilde{L}$ , and the input dimension  $d_x$  with  $d_0$  in Theorem 3.1, and under the the  $\beta_0$ -Hölder smoothness assumption.

204 205

169

170

171 172 173

174

181

182

183

185

186 187

188

## 5 DISCUSSION AND CONCLUSION

We examine the approximation and estimation rates of conditional DiT and its latent setting within
the "in-context" framework introduced by Peebles & Xie (2023), and conduct a comprehensive
analysis under various common data conditions. Notably, we establish the minimax optimality of
unconditional DiTs' estimation by reducing our analysis from conditional to unconditional settings.
Our approach employs a refined score decomposition scheme that enhances transformers' universal
approximation compared to earlier methods derived from the universal approximation results in (Yun et al., 2020) by Hu et al. (2024b).

213

## BOARDER IMPACT

This theoretical work explores the foundational aspects of generative diffusion models and is anticipated to have no adverse societal effects.

## 216 REFERENCES

- Anh-Dung Dinh, Daochang Liu, and Chang Xu. Rethinking conditional diffusion sampling with progressive guidance. *Advances in Neural Information Processing Systems*, 36, 2023.
- Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable
   creation in self-attention mechanisms. In *International Conference on Machine Learning (ICML)*,
   pp. 5793–5831. PMLR, 2022.
- Hengyu Fu, Zehao Dou, Jiawei Guo, Mengdi Wang, and Minshuo Chen. Diffusion transformer captures spatial-temporal dependencies: A theory for gaussian process data. *arXiv preprint arXiv:2407.16134*, 2024a.
- Hengyu Fu, Zhuoran Yang, Mengdi Wang, and Minshuo Chen. Unveil conditional diffusion models
   with classifier-free guidance: A sharp statistical theory. *arXiv preprint arXiv:2403.11968*, 2024b.
- Jerry Yao-Chieh Hu, Wei-Po Wang, Ammar Gilani, Chenyang Li, Zhao Song, and Han Liu. Fundamental limits of prompt tuning transformers: Universality, capacity and efficiency. *arXiv preprint arXiv:2411.16525*, 2024a.
- Jerry Yao-Chieh Hu, Weimin Wu, Zhuoru Li, Sophia Pi, , Zhao Song, and Han Liu. On statistical rates
   and provably efficient criteria of latent diffusion transformers (dits). In *Thirty-eighth Conference on Neural Information Processing Systems (NeurIPS)*, 2024b.
- Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank
   weight matrices universal approximators? In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution
   estimators. In *International Conference on Machine Learning*, pp. 26517–26582. PMLR, 2023.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4195–4205, 2023.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation
   function. *The Annals of Statistics*, 2020, 2020.
- Matus Telgarsky. Neural networks and rational functions. In *International Conference on Machine Learning*, pp. 3387–3393. PMLR, 2017.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computa- tion*, 23(7):1661–1674, 2011.
  - Junde Wu, Wei Ji, Huazhu Fu, Min Xu, Yueming Jin, and Yanwu Xu. Medsegdiff-v2: Diffusion-based medical image segmentation with transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 6030–6038, 2024.
- Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Zou,
   and Stefano Ermon. Tfg: Unified training-free guidance for diffusion models. *arXiv preprint arXiv:2409.15761*, 2024.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are trans formers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations (ICLR)*, 2020.
- 266

256

257

258

259

- 267 268
- 200

## Appendix

A	A Notation	7
F	3 Universal Approximation of Transformers	8
	B.1 Transformers as Universal Approximators	8
	B.2 Parameter Norm Bounds for Transformer Approximation	8
(	<b>C</b> Proof of Theorem 3.1 under Generic Assumption	10
	C.1 Auxiliary Lemmas	10
	C.2 Main Proof of Theorem 3.1 under Generic Assumption	22
Ι	Proof of Theorem 3.1 under Stronger Assumption	26
	D.1 Auxiliary Lemmas	26
	D.2 Main Proof of Theorem 3.1 under Stronger Assumption	36
F	E Proof of Theorem 3.2	38
	E.1 Auxiliary Lemmas for Theorem 3.2	38
	E.2 Proof of Theorem 3.2	41

## <sup>324</sup> A NOTATION

The index set  $\{1, ..., I\}$  is denoted by [I], where  $I \in \mathbb{N}^+$ . We denote (column) vectors by lower case letters, and matrices by upper case letters. Let a[i] denote the *i*-th component of vector *a*. Let  $A_{ij}$ denotes the (i, j)-th entry of matrix A. ||x||,  $||x||_1$  and  $||x||_{\infty}$  denote the Euclidean norm, 1-norm, and infinite norm.  $||W||_2$  and  $||W||_F$  denote the spectral norm and Frobenius norm, and  $||W||_{p,q}$  denotes the (p, q)-norm where *p*-norm is over columns and *q*-norm is over rows. We summarize our notations in the following table for easy reference.

Table 2. Mathematical Inotations and Symbol	Table 2	Mathematical	Notations	and Symbol
---	---------	--------------	-----------	------------

Symbol	Description
[I]	The index set $\{1,, I\}$ , where $I \in \mathbb{N}^+$
a[i]	The <i>i</i> -th component of vector <i>a</i>
$A_{ij}$	The $(i, j)$ -th entry of matrix A
	Euclidean norm of vector x
$  x  _1$	1-norm of vector x
$  x  _2$	2-norm of vector x
$\  x \ _{\infty}$	Spectral norm of matrix W
$\ W\ _{F}^{2}$	Frobenius norm of matrix W
$\ W\ _{p,q}^r$	(p,q)-norm of matrix W, where p-norm is over columns and q-norm is over rows
$  f(x)  _{L^2}$	$L^2$ -norm, where f is a function
$\ f(x)\ _{L^2(P)}$	$L^{2}(P)$ -norm, where f is a function and P is a distribution
$\ f(\cdot)\ _{Lip}$	Lipschitz-norm, where $f$ is a function
$d_p(f,g)$	<i>p</i> -norm of the difference between functions f and g defined as $d_p(f,g) = \left(\int  f(x) - g(x) ^p dx\right)^{1/p}$
$f_{\sharp}P$	Pushforward measure, where $f$ is a function and $P$ is a distribution
$\mathrm{KL}(P,Q)$	Kullback-Leibler (KL) divergence between distributions $P$ and $Q$
$\mathrm{TV}(P,Q)$	Total variation (TV) distance between distributions $P$ and $Q$
$N(\mu, \sigma^2)$	Normal distribution with mean $\mu$ and variance $\sigma^2$
$a \gtrsim b$	There exist constants $C > 0$ such that $a \le Cb$
n	Sample size
x	Data point in original data space, $x \in \mathbb{R}^{d_x}$
y	Conditioning Label, $x \in \mathbb{R}^{d_y}$
$\frac{n}{L}$	Latent variable in low-dimensional subspace, $n \in \mathbb{R}^{n_0}$
n n	$h = U^{-x}$
$U^{p_n}$	The matrix with orthonormal columns to transform h to x, where $U \in \mathbb{R}^{d \times d_0}$
В	Radius of Hölder ball for conditional density function $p(x y)$
$B_0$	Radius of Hölder ball for latent conditional density function $p(\bar{w} y)$
$\beta$	Hölder index for conditional density function $p(x y)$
$\beta_0$	Hölder index for latent conditional density function $p(\overline{h} y)$
D	Granularity in the construction of the transformer universal approximation
N	Resolution of the discretization of the input domain
$\mathcal{R}$	Score risk (expectation of squared $\ell^2$ difference between score estimator and ground truth)
$\mathcal{N}(\epsilon, \mathcal{F}, \ \cdot\ )$	Covering number of collection $\mathcal{F}$
Т	Stopping time in the forward process of diffusion model
$t_0$	Stopping time in the backward process of diffusion model
$\mu$	Discretized step size in backward process
$p_t(\cdot)$	The density function of $\overline{b}$ at time t
$\psi^{Pt}(\cdot)$	(Conditional) Gaussian density function
$\tau$ $\tau h.s.r$	
$f_{h,s,r}$	Transformer block of h head e hidden size e MI P dimension
d d	Industoring office of <i>n</i> -field, s-fielden size, <i>i</i> -field fullension Input dimension of each token in the transformer network of DiT
$\tilde{L}$	Token length in the transformer network of DiT
$\widetilde{d}$	Latent data input dimension of each token in the transformer network of DiT
$\tilde{\tilde{L}}$	Latent data token length in the transformer network of DiT
X	Sequence input of transformer network in DiT, where $X \in \mathbb{R}^{d \times L}$
Η	Sequence latent data input of transformer network in DiT. where $X \in \mathbb{R}^{d \times L}$
E	Position encoding, where $E \in \mathbb{R}^{d \times L}$
$R(\cdot)$	Reshape layer in DiT, $R(\cdot) : \mathbb{R}^{d_x} \to \mathbb{R}^{d \times L}$
$\widetilde{R}(\cdot)$	Reshape layer in Dif. $\widetilde{R}(\cdot)$ : $\mathbb{R}^{d_0} \to \mathbb{R}^{\widetilde{d} \times \widetilde{L}}$
	By the problem layer is $\mathbb{D}^{T}$ , $\mathbb{D}^{-1}(\cdot)$ , $\mathbb{D}^{d} \times L \to \mathbb{D}^{d}$
$R^{-1}(\cdot)$	Reverse reshape layer in D11, $\Lambda$ (·): $\mathbb{R}^{-1} \to \mathbb{R}^{-2}$
$\frac{R^{-1}(\cdot)}{\widetilde{R}^{-1}(\cdot)}$	Reverse reshape layer in DT, $\widetilde{R}^{-1}(\cdot) : \mathbb{R}^{\widetilde{d} \times \widetilde{L}} \to \mathbb{R}^{d_0}$

#### 378 В UNIVERSAL APPROXIMATION OF TRANSFORMERS 379

380 In this section, we discuss the universal approximation theory of transformers.

In Appendix B.1, we present the universal approximation results of transformers for score approxi-382 mation. We emphasize that most of the material in Appendix B.1 is not original and is drawn from (Hu et al., 2024a; Kajitsuka & Sato, 2024; Yun et al., 2020). 384

In Appendix B.2, we compute the parameter norm bounds of the transformers used for score approximation. These bounds are crucial for calculating the covering number of the transformers and are essential for score and distribution estimation.

**B.1** TRANSFORMERS AS UNIVERSAL APPROXIMATORS

385

386

387 388

389 390

391

392

393

394

395

403

404

Theorem B.1 (Transformers with 1-Layer Self-Attention are Universal Approximators, Modified from Proposition 1 of (Kajitsuka & Sato, 2024)). Let  $0 \le p < \infty$  and  $f^{(\text{FF})}$ ,  $f^{(\text{SA})}$  be feed-forward neural network layers and a single-head self-attention layer with softmax function respectively. Then, for any permutation equivariant, continuous function f with compact support and  $\epsilon > 0$ , there exists  $f' \in \mathcal{T}_R^{\hat{h},s,r}$  such that  $d_p(f,f') < \epsilon$  holds

Lastly, we provide the next corollary stating that the required transformer configuration (h, s, r) for 397 universal approximation. 398

Corollary B.1.1 (Universal Approximation of Transformers). From Theorem B.1, for any permu-399 tation equivariant, continuous function f with compact support and  $\epsilon > 0$ , a transformer network  $f' \in \mathcal{T}_R^{1,1,4}$  with MLP dimension (width) r = 4 and  $= \mathcal{O}((1/\epsilon)^{dL})$  FFN layers is sufficient to 400 401 approximate f such that  $d_p(f, f') < \epsilon$ . 402

**Remark B.1.** We remark that  $\mathcal{T}_{R}^{1,1,4}$  belongs to the considered transformer network function class Definition 2.1. 405

406 We establish in Corollary B.1.1 the minimal transformer configuration required to achieve universal 407 approximation for compactly supported functions. We remark that this configuration is minimally sufficient but not necessary. More complex configurations can also achieve transformer universality, 408 as reported in (Hu et al., 2024b; Kajitsuka & Sato, 2024; Yun et al., 2020). 409

410 Throughout this paper, unless otherwise specified, we use the transformer class  $\mathcal{T}_{R}^{1,1,4}$  to construct 411 score function approximations. 412

413 **B.2** PARAMETER NORM BOUNDS FOR TRANSFORMER APPROXIMATION 414

415 In the analysis of the approximation ability of transformers in (Kajitsuka & Sato, 2024), universal approximation is ensured by using a sufficiently large granularity D, a sufficiently small  $\delta$  in  $f_1^{(FF)}$ , 416 and an appropriate scaling factor R in  $f_2^{(FF)}$ . Here, we provide a detailed discussion on parameter bounds for matrices in  $\mathcal{T}_R^{h,r,s}$ , focusing on the choice of granularity and scaling factor. 417 418 419

420 Lemma B.1 (Order of Granularity and Scaling Factor). Consider the universal approximation 421 theorem for transformers in Theorem B.1. The order for the granularity and the scaling factor 422 follows  $D = \mathcal{O}(\epsilon^{-1/d})$  and  $R = \mathcal{O}(D)$ , and the parameter  $\delta$  for the first feed-forward layer follows  $\delta = o(D^{-1}).$ 423

424 Building upon Lemma B.1, we extend the results to derive explicit parameter bounds for matrices 425 regarding the transformer-based universal approximation framework. That is, we ensure a more 426 precise quantification of parameter constraints across the architecture. 427

428 **Lemma B.2** (Transformer Matrices Bounds). Consider an input sequence  $Z \in [0,1]^{d \times L}$ . Let 429  $f(Z): [0,1]^{d \times L} \to \mathbb{R}^{d \times L}$  be any permutation equivariant and continuous sequence-to-sequence function on compact support  $[0,1]^{d \times L}$ . For the transformer network  $f' \in \mathcal{T}_R^{r,h,s}$  to approximate 430 f within  $\epsilon$  precision, i.e.,  $d_p(f, f') < \epsilon$ , the following parameter bounds must hold for  $d \ge 1$  and 431

432	
433	$L \ge 2$ :
434	(2dL+1), 1
435	$  W_Q  _2 =   W_K  _2 = \mathcal{O}(d \cdot e^{-(\underline{d})})(\log L)^{\frac{1}{2}});$
436	$  W_{O}  _{\bullet} =   W_{V}  _{\bullet} = \mathcal{O}(d^{\frac{3}{2}} \cdot e^{-(\frac{2dL+1}{d})}(\log L)^{\frac{1}{2}})$
437	$\  (VQ) \ _{2,\infty} = \  (VX) \ _{2,\infty} = C ((0, 0, 0)) $
438	$\left\ W_{O}\right\ _{2} = \mathcal{O}\left(\sqrt{d}\epsilon^{\frac{1}{d}}\right); \left\ W_{O}\right\ _{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{d}}\right);$
430	
440	$\ W_V\ _2 = \mathcal{O}(\sqrt{d}); \ W_V\ _{2,\infty} = \mathcal{O}(d);$
441	$\ W_1\ _{\epsilon} = \mathcal{O}\left(d\epsilon^{-\frac{1}{d}}\right) \ W_1\ _{\epsilon} = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}}\right)$
442	$\  \mathcal{V}_1 \ _2 = \mathcal{O}\left( \left( \begin{array}{c} u \\ y \end{array} \right), \  \mathcal{V}_1 \ _{2,\infty} = \mathcal{O}\left( \left( \begin{array}{c} u \\ y \end{array} \right), \\ (u \\ y \end{array} \right),$
443	$\ W_2\ _2 = \mathcal{O}\left(d\epsilon^{-\frac{1}{d}}\right); \ W_2\ _2 = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}}\right);$
444	(1,3)
445	$\left\ E^+ ight\ _{2,\infty}=\mathcal{O}\left(d^{rac{1}{2}}L^{rac{2}{2}} ight).$
446	· · · · · · · · · · · · · · · · · · ·
447	For the case $L = 1$ , the parameter bounds remain valid with the substitution of $\log L$ with 1.
448	
449	
450	
451	
452	
453	
454	
455	
456	
457	
458	
459	
460	
461	
462	
403	
404	
465	
467	
468	
469	
470	
471	
472	
473	
474	
475	
476	
477	
478	
479	
480	
481	
482	
483	
484	
400	

## 486 C PROOF OF THEOREM 3.1 UNDER GENERIC ASSUMPTION

Our proof builds on the local smoothness properties of functions within Hölder spaces and the universal approximation of transformers. While the universal approximation theory of transformers ensures arbitrarily small errors, it does not account for the smoothness of functions in the result. To incorporate the smoothness assumptions of interest, we propose the following three steps to integrate function smoothness into approximation theory of transformer architectures.

• Step 1. Consider the integral form of  $p_t(x_t|y)$  in (C.1). We clip the input domain  $\mathbb{R}^{d_x}$  into closed and bounded region  $B_{x,N}$ . This facilitates the error analysis for the Taylor expansion approximation in the next step. The clipping error arises from the integral over the region outside  $B_{x,N}$ . We specify the clipping error in Lemma C.1.

$$p_t(x_t|y) = \int_{\mathbb{R}^{d_x}} \frac{\mathrm{d}x_0}{\sigma_t^{d_x}(2\pi)^{d_x/2}} \cdot \underbrace{p_0(x_0|y)}_{\approx k_1 \text{-order Taylor polynomial}} \cdot \underbrace{\exp\left(-\frac{\|\alpha_t x_0 - x_t\|^2}{2\sigma_t^2}\right)}_{\approx k_2 \text{-order Taylor polynomial}}.$$
 (C.1)

- Step 2. We employ  $k_1$ -order and  $k_2$ -order Taylor expansion for  $p(x_0|y)$  and  $\exp(\cdot)$  in (C.1). We construct the diffused local polynomial in Lemma C.2 based on the Taylor expansion. We approximate  $p_t$  and  $\nabla p_t$  with the diffused local polynomial  $f_1(x, y, t) \in \mathbb{R}$  and  $f_2(x, y, t) \in \mathbb{R}^{d_x}$  in Lemma C.3 and Lemma C.4.
- Step 3. We approximate  $f_1(x, y, t)$ ,  $f_2(x, y, t)$  with transformers in Lemmas C.5 and C.6. To construct the final score approximator with the transformer, we approximate necessary algebraic operators in Lemmas C.7 to C.11. We provide the output bound of our transformer model in Lemma C.12. We combine all components into Lemma C.13, and complete the proof of Theorem 3.1.

Noe that the proof under latent subspace assumption in Table 1 closely follows the proof in this section, with the input dimension d, L to  $\tilde{d}$  and  $\tilde{L}$ , and the input dimension  $d_x$  with  $d_0$  in Theorem 3.1, and consider under the  $\beta_0$ -Hölder smoothness assumption on latent data.

**Organization.** Appendix C.1 includes details regarding the three steps with auxiliary lemmas for supporting our proof. Appendix C.2 includes the main proof of Theorem 3.1.

C.1 AUXILIARY LEMMAS

 Step 1: Clip  $\mathbb{R}^{d_x} \times [0,1]^{d_y}$  for  $p_t(x|y)$ . We introduce a helper lemma on the clipping integral.

**Lemma C.1** (Approximating Clipped Multi-Index Gaussian Integral, Lemma A.8 of (Fu et al., 2024b)). Under generic Assumption 3.1. Consider any integer vector  $\kappa \in \mathbb{Z}_+^{d_x}$  with  $\|\kappa\|_1 \leq n$ . There exists a constant  $C(n, d_x) \geq 1$ , such that for any  $x \in \mathbb{R}^{d_x}$  and  $0 < \epsilon \leq 1/e$ , it holds

$$\int_{\mathbb{R}^{d_x} \setminus B_x} \left| \left( \frac{\alpha_t x_0 - x}{\sigma_t} \right)^{\kappa} \right| \cdot p(x_0 | y) \cdot \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left( -\frac{\left\| \alpha_t x_0 - x \right\|^2}{2\sigma_t^2} \right) \mathrm{d}x_0 \le \epsilon, \tag{C.2}$$

where  $\left(\frac{\alpha_t x_0 - x}{\sigma_t}\right)^{\kappa} \coloneqq \left(\left(\frac{\alpha_t x_0[1]_1 - x[1]}{\sigma_t}\right)^{\kappa[1]}, \left(\frac{\alpha_t x_0[2] - x[2]}{\sigma_t}\right)^{\kappa[2]}, \dots, \left(\frac{\alpha_t x_0[d_x] - x[d_x]}{\sigma_t}\right)^{\kappa[d_x]}\right)$  is a *multi-indexed* vector and

$$B_x \coloneqq \left[\frac{x - \sigma_t C(n, d_x) \sqrt{\log(1/\epsilon)}}{\alpha_t}, \frac{x + \sigma_t C(n, d_x) \sqrt{\log(1/\epsilon)}}{\alpha_t}\right]$$
$$\bigcap \left[-C(n, d_x) \sqrt{\log(1/\epsilon)}, C(n, d_x) \sqrt{\log(1/\epsilon)}\right]^{d_x}.$$

**Remark C.1.**  $B_x$  is a bounded domain. Lemma C.1 provides the difference between integrals of the form (C.2) on  $\mathbb{R}^{d_x}$  and on  $B_x$ . The difference becomes arbitrarily small with precision  $\epsilon = 1/N$ .

**Step 2:** Approximate  $p_t(x|y)$  and  $\nabla p_t(x|y)$  with Taylor Expansion. We begin with the definition.

**Definition C.1** (Normalization of  $B_{x,N}$ ). Consider the clipping in Lemma C.1 and the initial conditional distribution  $p(x_0|y)$  with closed and bounded support  $B_{x,N} \times [0,1]^{d_y}$ . We define  $R_B := (2C(0,d)\sqrt{\beta \log N})$  and  $x'_0 := x_0/R_B + 1/2$ . Moreover, we define  $M(x'_0, y) := p(R_B(x'_0 - 1/2)|y)$ .

**Remark C.2.** The purpose of Definition C.1 is to simplify the process of discretizing  $B_{x,N} \times [0,1]^{d_y}$ into  $N^{d_x+d_y}$  hypercubes. In particular,  $M(x'_0, y)$  has compact support on  $[0,1]^{d_x+d_y}$ , where  $R_B$ denotes the length of each coordinate of  $B_{x,N}$ , and  $x'_0 \in [0,1]^{d_x}$  represents  $x_0$  normalized on  $B_{x,N}$ . **Remark C.3.** The only difference between  $M(x'_0, y)$  and  $p(x_0|y)$  lies in their respective domains, leading to the difference in the size of the Hölder ball radius. Recall that under generic Assumption 3.1, we have  $p(x_0|y) \in \mathcal{H}^{\beta}(\mathbb{R}^{d_x} \times [0,1]^{d_y}, B)$ . Here we have  $M(x'_0, y) \in \mathcal{H}([0,1]^{d_x+d_y}, BR^{k_1}_B)$ . This follows from the fact that  $p(\cdot|y)$  is  $k_1$ -time differentiable so that the radius scale by a factor of  $R^{k_1}_B$ .

**Lemma C.2** (Diffused Local Polynomial, Modified from (Fu et al., 2024a)). Under generic Assumption 3.1. We write  $p_t(x|y)$  into the product of  $p(x_0|y)$  and  $\exp(\cdot)$ :

$$p_t(x|y) = \int_{\mathbb{R}^{d_x}} p(x_0|y) p_t(x|x_0) dx_0 = \int_{\mathbb{R}^{d_x}} \frac{1}{\sigma_t^{d_x} (2\pi)^{d_x/2}} p(x_0|y) \exp\left(-\frac{\|\alpha_t x_0 - x\|^2}{2\sigma_t^2}\right) dx_0.$$

Then we approximate  $p(x_0|y)$  and  $\exp\left(-\frac{\|\alpha_t x_0 - x\|^2}{2\sigma_t^2}\right)$  with  $k_1$ -order Taylor polynomial and  $k_2$ -order Taylor polynomial within  $B_{x,N}$  respectively. Altogether, we approximate  $p_t(x|y)$  with the following *diffused local polynomial* with the bounded domain  $B_{x,N}$  around x:

$$f_{1}(x,y,t) = \sum_{v \in [N]^{d}, w \in [N]^{d_{y}}} \sum_{\|n_{x}\|_{1} + \|n_{y}\|_{1} \le k_{1}} \frac{R_{B}^{\|n_{x}\|}}{n_{x}!n_{y}!} \frac{\partial^{n_{x}+n_{y}}p}{\partial x^{n_{x}}\partial y^{n_{y}}} \bigg|_{x = R_{B}(\frac{v}{N} - \frac{1}{2}), y = \frac{w}{N}} \Phi_{n_{x},n_{y},v,w}(x,y,t),$$
(C.3)

where

•  $\phi(\cdot)$  is the trapezoid function.

• 
$$g(x, n_x, v, k_2) \coloneqq \frac{1}{\sigma_t \sqrt{2\pi}} \int \left(\frac{x_0}{R} + \frac{1}{2} - \frac{v}{N}\right)^{n_x} \frac{1}{k_2!} \left(-\frac{|x - \sigma_t x_0^2|}{2\sigma_t^2}\right)^{k_2} \mathrm{d}x_0.$$
  
•  $\Phi_{n_x, n_y, v, w}(x, y, t) \coloneqq \left(y - \frac{w}{N}\right)^{n_y} \prod_{j=1}^{d_y} \phi\left(3N(y[j] - \frac{w}{N})\right) \prod_{i=1}^{d_x} \sum_{k_2 < p} g(x[i], n_x[i], v[i], k_2).$ 

**Remark C.4.** The form of the diffused local polynomial arises from the Taylor expansion approximation applied on each grid point within  $[0, 1]^{d_x+d_y}$ , with  $v \in [N]^{d_x}$  and  $w \in [N]^{d_y}$  denoting the specific grid point undergoing approximation.

**Remark C.5.** The Hölder space assumption in generic Assumption 3.1 establishes an upper bound on the error arising from the remainder term in the Taylor expansion. This ensures the approximation accuracy is well-controlled.

We specifies the error from the approximation of  $p_t$  and  $\nabla p_t$  with  $f_1$  and  $f_2$  in Lemmas C.3 and C.4.

**Lemma C.3** (Approximation of  $p_t(x|y)$  by Polynomials, Lemma A.4 of (Fu et al., 2024b)). Under generic Assumption 3.1. For any  $x \in \mathbb{R}^{d_x}$ ,  $y \in [0, 1]^{d_y}$ , t > 0, and a sufficiently larger N > 0, there exists a diffused local polynomial  $f_1(x, y, t)$  with at most  $N^{d_x+d_y}(d_x + d_y)^{k_1}$  monomials such that

 $|f_1(x, y, t) - p_t(x|y)| \lesssim BN^{-\beta} \log^{\frac{d_x + k_1}{2}} N.$ 

**Lemma C.4** (Approximation of  $\nabla \log p_t(x|y)$  by Polynomials, Lemma A.6 of (Fu et al., 2024b)). Under generic Assumption 3.1. For any  $x \in \mathbb{R}^{d_x}$ ,  $y \in [0,1]^{d_y}$ , t > 0, and a sufficiently larger N > 0, there exists  $f_2 \coloneqq (f_2[1], \ldots, f_2[d_x])^\top \in \mathbb{R}^{d_x}$  with local diffused polynomial  $f_2[i]$  such that

$$|f_2(x, y, t)[i] - \sigma_t \nabla p_t(x|y)[i]| \lesssim BN^{-\beta} \log^{\frac{d_x+k_1+1}{2}} N,$$

594 where each  $f_2[i]$  contains at most  $N^{d_x+d_y}(d_x+d_y)^{k_1}$  monomials. 595 596 We have finished the approximation of  $p_t$  and  $\nabla p_t$  with diffused local polynomial  $f_1$  and  $f_2$ . 597 Step 3. Approximate Diffused Local Polynomials and Algebraic Operators with Transformers. 598 First, we utilize universal approximation capabilities of transformers to deal with  $f_1, f_2$  established in previous step. Second, we employ similar scheme to approximate several algebraic operators 600 necessary in final score approximation. Lastly, we present the incorporation of these components in 601 Lemma C.13 with a unified transformer architecture and corresponding parameter configuration. 602 603 • Step 3.1: Approximate the Diffused Local Polynomials  $f_1$  and  $f_2$ . 604 We invoke the universal approximation theorem of transformer (Theorem B.1). We utilize 605 network consisting of one transformer block and one feed-forward layer. Lemma C.5 (Approximate Scalar Polynomials with Transformers). Under generic Assump-607 tion 3.1. Consider the diffused local polynomial  $f_1$  in Lemma C.3. For any  $\epsilon > 0$ , there exists 608 a transformer  $\mathcal{T}_{f_1} \in \mathcal{T}_R^{h,s,r}$ , such that for any  $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}, y \in [0,1]^{d_y}$ 609 and  $t \in [N^{-C_{\sigma}}, C_{\alpha} \log N]$  it holds 610 611  $|f_1(x, y, t) - \mathcal{T}_{f_1}(x, y, t)[d_x]| < \epsilon.$ 612 613 The parameter bounds in the Transformer network class satisfy 614 615  $\|W_Q\|_2, \|W_K\|_2 = \mathcal{O}\left(d\epsilon^{-\frac{2dL+4d+1}{d}} (\log L)^{\frac{1}{2}}\right);$ 616  $\|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(d^{\frac{3}{2}} \epsilon^{-\frac{2dL+4d+1}{d}} (\log L)^{\frac{1}{2}}\right);$ 617 618  $||W_V||_2 = \mathcal{O}(\sqrt{d}); ||W_V||_2 = \mathcal{O}(d);$ 619 620  $\left\|W_{O}\right\|_{2} = \mathcal{O}\left(\sqrt{d}\epsilon^{\frac{1}{d}}\right); \left\|W_{O}\right\|_{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{d}}\right);$ 621  $\left\|W_{1}\right\|_{2} = \mathcal{O}\left(d\epsilon^{-\frac{1}{d}} \cdot \log N\right); \left\|W_{1}\right\|_{2,\infty} = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}} \cdot \log N\right);$ 622 623  $\left\|W_{2}\right\|_{2} = \mathcal{O}\left(d\epsilon^{-\frac{1}{d}}\right); \left\|W_{2}\right\|_{2,\infty} = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}}\right); \left\|E^{\top}\right\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right).$ 624 625 626 Similarly, we have the corresponding  $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$  for the approximation of  $f_2(x, y, t)$ . 627 Lemma C.6 (Approximate Vector-Valued Polynomials with Transformers). Under generic 628 Assumption 3.1 and consider  $f_2(x, y, t) \in \mathbb{R}^{d_x}$  with every entry  $f_2[1], \ldots, f_2[d_x]$  is a local diffused polynomial defined in Lemma C.2. For any  $\epsilon > 0$ , there exists a transformer 630  $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$  such that 631 632  $\|f_2(x, y, t) - \mathcal{T}_{f_2}\|_{\infty} \le \epsilon,$ 633 634 for any  $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}, y \in [0, 1]^{d_y}$  and  $t \in [N^{-C_\sigma}, C_\alpha \log N]$ . The 635 parameter bounds in the transformer network class follows Lemma C.5. 636 So far, we have obtained approximation results for  $f_1$  and  $f_2$ . To complete the full approxi-637 mation of the score decomposition  $\nabla \log p = \frac{\nabla p}{p}$ , we still need to approximate several key 638 algebraic operators, including the product (Lemma C.8), inverse (Lemma C.9)...etc. 639 We establish their approximations as follows. 640 641 Step 3.2: Approximate Algebraic Operators with Transformers. 642 We give transformer approximation theory for the clipping operator, the inverse operator, 643 the product operator, and functions that evolve with time t: - Clipping operation (Lemma C.7) 645 Product operation (Lemma C.8) 646 - Inverse operation (Lemma C.9) 647

- Mean  $\alpha_t = \exp(-t/2)$  (Lemma C.10)

- Standard deviation  $\sigma_t = \sqrt{1 - e^{-t}}$  (Lemma C.11)

The approximations for these operators are common with the network structure consisting of ReLU activation function and fully connected feed-forward layers, such as the product approximation by Schmidt-Hieber (2020) and the inverse approximation by Telgarsky (2017).

The following lemma provides a network that executes the clipping operation.

**Lemma C.7** (Clipping Operation, Lemma F.4 of (Oko et al., 2023)). For any  $a, b \in \mathbb{R}^d$  with  $a[i] \leq b[i]$  for all  $i \in [d]$ , there exist a neural network  $\phi_{\text{clip}}(x; a, b) \in \Phi(L, W, S, B)$  such that for all  $i \in [d]$ , it holds

$$\phi_{\mathsf{clip}}(x;a,b)[i] = \min(b[i], \max(x[i], a[i])),$$

with

$$L = 2, \quad W = (d, 2d, d)^{\top}, \quad S = 7d, \quad B = \max_{1 \le i \le d} \max(|a[i]|, b[i]).$$
 (C.4)

Moreover, suppose a[i] = c and b[i] = C for all  $i \in [d]$  with c and C being some constant,  $\phi_{clip}(x; a, b)$  is denoted as  $\phi_{clip}(x; c, C)$ .

Next, we deal with the approximation of products with Transformer.

**Lemma C.8** (Approximation of the Product Operator with Transformer.). Let  $m \ge 2$  and  $C \ge 1$ . For any  $0 < \epsilon_{\text{mult}} < 1$ , there exists  $\mathcal{T}_{\text{mult}}(\cdot) \in \mathcal{T}_R^{h,s,r}$  such that for all  $x \in [-C, C]^m$ ,  $x' \in \mathbb{R}^m$  with  $||x - x'||_{\infty} \le \epsilon_{\text{error}}$ , it holds

$$\left|\mathcal{T}_{\text{mult}}(x') - \prod_{i=1}^{m} x_i\right| \le \epsilon_{\text{mult}} + mC^{m-1}\epsilon_{\text{error}}.$$

The parameter bounds in the transformer network class  $\mathcal{T}_{R}^{h,s,r}$  satisfy

$$\begin{split} \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\text{mult}}^{-(2m+1)}(\log m)^{\frac{1}{2}}\right); \\ \|W_O\|_2, \|W_O\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\text{mult}}^m\right); \quad \|W_V\|_2, \|W_V\|_{2,\infty} &= \mathcal{O}(1); \\ \|W_1\|_2, \|W_1\|_{2,\infty} &= \mathcal{O}\left(C\epsilon_{\text{mult}}^{-m}\right); \quad \|W_2\|_2, \|W_2\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\text{mult}}^{-m}\right). \end{split}$$

Next, we introduce the next lemma to approximate the inverse operator.

**Lemma C.9** (Approximation of the Reciprocal Function with Transformer.). For any  $0 < \epsilon_{\text{rec}} < 1$  there exists a  $\mathcal{T}_{\text{rec}}(\cdot) \in \mathcal{T}_R^{h,s,r}$  such that for all  $x \in [\epsilon_{\text{rec}}, \epsilon_{\text{rec}}^{-1}]$  and  $x' \in \mathbb{R}$ . It holds that

$$\left|\mathcal{T}_{\text{rec}}(x') - \frac{1}{x}\right| \le \epsilon_{\text{rec}} + \frac{|x - x'|}{\epsilon_{\text{rec}}^2}$$

The parameter bounds in the Transformer network class satisfy

$$\begin{aligned} \|W_Q\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_2, \|W_K\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\rm rec}^{-3}\right); \\ \|W_O\|_2, \|W_O\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\rm rec}\right); \|W_V\|_2, \|W_V\|_{2,\infty} &= \mathcal{O}(1); \\ \|W_1\|_2, \|W_1\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\rm rec}^{-2}\right); \|W_2\|_2, \|W_2\|_{2,\infty} &= \mathcal{O}\left(\epsilon_{\rm rec}^{-1}\right). \end{aligned}$$

Next, we state approximation results using Transformer for  $\alpha_t$  and  $\sigma_t$ . Note that we have  $\alpha_t = \exp(-t/2)$  and  $\sigma_t = \sqrt{1 - \alpha_t^2}$ .

702 **Lemma C.10** (Approximation of  $\alpha_t = \exp(-t/2)$  with Transformer.). For any  $\epsilon_{\alpha} \in (0, 1)$ , there exists Transformer  $\mathcal{T}_{\alpha}(t) \in \mathcal{T}_{R}^{h,s,r}$  such that for all  $t \geq 0$ , we have 704 705  $|\mathcal{T}_{\alpha}(t) - \alpha_t| < \epsilon_{\alpha}.$ 706 The parameter bounds in the Transformer network class satisfy 708  $||W_Q||_2, ||W_Q||_{2,\infty}, ||W_K||_2, ||W_K||_{2,\infty} = \mathcal{O}(\epsilon_{\alpha}^{-3});$ 709 710  $||W_O||_2, ||W_O||_2 = \mathcal{O}(\epsilon_{\alpha}^{-1}); ||W_V||_2, ||W_V||_2 = \mathcal{O}(1);$ 711  $||W_1||_2, ||W_1||_{2,\infty} = \mathcal{O}\left((\log \epsilon_{\alpha}^{-1})\epsilon_{\alpha}^{-1}\right); ||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(\epsilon_{\alpha}^{-1}\right).$ 712 713 **Lemma C.11** (Approximation of  $\sigma_t = \sqrt{1 - e^{-t}}$  with transformer). For any  $\sigma_{\sigma} \in (0, 1)$ , 714 there exists a transformer  $\mathcal{T}_{\sigma}(t) \in \mathcal{T}_{R}^{h,s,r}$  such that for any  $t \in [t_0,T]$  with  $t_0 < 1$  we have 715 716  $|\mathcal{T}_{\sigma}(t) - \sigma_t| < \epsilon_{\sigma}.$ 717 718 The parameter bounds in the transformer network class satisfy 719 720  $||W_Q||_2, ||W_Q||_2, \dots, ||W_K||_2, ||W_K||_2, \dots = \mathcal{O}(\epsilon_{\sigma}^{-3});$ 721  $||W_O||_2, ||W_O||_2 = \mathcal{O}(\epsilon_{\sigma}); ||W_V||_2, ||W_V||_2 = \mathcal{O}(1);$ 722  $\|W_1\|_2 = \mathcal{O}\left(T\epsilon_{\sigma}^{-1}\right); \quad \|W_1\|_2 = \mathcal{O}\left(T\epsilon_{\sigma}^{-1}\right);$ 723 724  $\|W_2\|_2 = \mathcal{O}\left(\epsilon_{\sigma}^{-1}\right); \quad \|W_2\|_2 = \mathcal{O}\left(\epsilon_{\sigma}^{-1}\right).$ 725 726 We have finished the approximation of every key component for the proof of Theorem 3.1. 727 We now proceed to the detailed assembly and integration of these components to finalize the 728 proof. 729 Step 3.3: Unified Transformer-Based Score Function Approximation. 730 First, we establish a theoretical upper bound for transformer model output by analyzing the 731 upper bound of the score function in  $\ell_{\infty}$  distance under generic Assumption 3.1 as follows. 732 – Bound on  $p_t(x|y)$ : 733 Recall that the conditional distribution at time t has the form: 734 735  $p_t(x|y) = \frac{1}{\sigma^d (2\pi)^{\frac{d}{2}}} \int p(x_0|y) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) \mathrm{d}x_0.$ 738 Applying the light tail property in generic Assumption 3.1, the upper bound follows: 739 740  $p_t(x|y) \le \frac{C_1}{\sigma_t^d(2\pi)^{\frac{d}{2}}} \int \exp\left(-\frac{C_2 \|x_0\|^2}{2}\right) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) \mathrm{d}x_0.$ (C.5) 741 742 743 On the other hand, the lower bound follows: 744 745  $p_t(x|y) \ge \frac{1}{\sigma_t^d(2\pi)^{\frac{d}{2}}} \int_{\|x_0\| \le 1} p(x_0|y) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) \mathrm{d}x_0.$ 746 (C.6) 747 748 - Bound on  $\nabla p_t(x|y)$ : The first element of the gradient has the form: 749 750  $|(\nabla p_t)[1]| = \frac{1}{\sigma_t^2 (2\pi)^{\frac{d}{2}}} \cdot \left| \int \left( \frac{x[1] - \alpha_t x_0[1]}{\sigma_t^2} \right) p(x_0|y) \exp\left( -\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2} \right) \mathrm{d}x_0 \right|.$ 751 752 753 754 The  $\ell_{\infty}$  bound on  $\nabla p_t$  follows by applying light tail property to each coordinate as in (C.5).

 Combining (C.5), (C.6) and (C.7), we provide the  $\ell_{\infty}$  bounds on the score.

**Lemma C.12** (Bounds on Score, Lemma A.10 of (Fu et al., 2024b)). Assume generic Assumption 3.1. There exists a constant *K* such that

$$\|\nabla \log p_t(x|y)\|_{\infty} \le \frac{K}{\sigma_t^2}(\|x\|+1).$$

Further details regarding the derivation are in Appendix A.7 of (Fu et al., 2024b). Next lemma incorporates previous approximation results into an unified transformer architecture.

**Lemma C.13** (Approximation Score Function with Transformer on Supported Domain). Under generic Assumption 3.1. Consider  $t \in [N^{-C_{\sigma}}, C_{\alpha} \log N]$ , for constant  $C_{\sigma}, C_{\alpha}$ , and  $(x, y) \in -[C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x} \times [0, 1]^{d_y}$ , where  $N \in \mathbb{N}$  and  $C_x$  depends on  $d, \beta, B, C_1, C_2$ . There exist a transformer network  $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_R^{h,s,r}$  such that

$$p_t(x|y) \|\nabla \log p_t(x|y) - \mathcal{T}_{\text{score}}(x,y,t)\|_{\infty} \lesssim \frac{B}{\sigma_t^2} N^{-\beta} (\log N)^{\frac{d_x+k_1+1}{2}}.$$

The parameter bounds in the Transformer network class satisfy

$$\begin{split} \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} &= \mathcal{O}\left(N^{(7\beta+6C_{\sigma})}\right);\\ \|W_O\|_2, \|W_O\|_{2,\infty} &= \mathcal{O}\left(N^{-(3\beta+6C_{\sigma})}(\log N)^{3(d_x+\beta)}\right);\\ \|W_V\|_2 &= \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} &= \mathcal{O}(d); \|E^{\top}\|_{2,\infty} &= \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);\\ \|W_1\|_2, \|W_1\|_{2,\infty} &= \mathcal{O}\left(N^{(2\beta+4C_{\sigma})}\right); C_{\mathcal{T}} &= \mathcal{O}\left(\sqrt{\log N}/\sigma_t^2\right);\\ \|W_2\|_2, \|W_2\|_{2,\infty} &= \mathcal{O}\left(N^{(3\beta+2C_{\sigma})}\right). \end{split}$$

Proof of Lemma C.13. Our poof follows the structure of Fu et al. (2024b, Proposition A.3). Recall that from Lemma C.12, we have  $\|\nabla \log p_t(x|y)\|_{\infty} \leq K(C_x\sqrt{d_x \log N} + 1)/\sigma_t^2$ , along with the diffused local polynomial  $f_1$  and  $f_2$ , we define first-step score approximator  $f_3(x, y, t)$  as

$$f_3(x, y, t) = \min\left(\frac{f_2}{\sigma_t f_{1, \text{clip}}}, \frac{K}{\sigma_t^2} (C_x \sqrt{d_x \log N} + 1)\right)$$

where we set  $f_{1,\text{clip}} = \{f_1, \epsilon_{\text{low}}\}$  to prevent score from blowing up and we set  $\epsilon_{\text{low}}$  later. We proceed with the following three steps:

– Step A. Approximate Score Function with  $f_3$ .

Without loss of generality, we first derive error bound on the difference between the first component in  $f_3$  and the score.

$$\begin{aligned} |(\nabla \log p_t)[1] - f_3[1]| &\leq \left| (\nabla \log p_t)[1] - \frac{f_2[1]}{\sigma_t f_{1,\text{clip}}} \right| \\ &\leq \left| \frac{(\nabla p_t)[1]}{p_t} - \frac{(\nabla p_t)[1]]}{f_{1,\text{clip}}} \right| + \left| \frac{(\nabla p_t)[1]}{f_{1,\text{clip}}} - \frac{f_2[1]}{\sigma_t f_{1,\text{clip}}} \right|. \end{aligned}$$

From Lemma C.12, the bound on the score implies  $(\nabla p_t)[1] \leq K(\sqrt{d_x \log N} + 1)p_t/\sigma_t^2$ .

Therefore,

$$\begin{aligned} &|(\nabla \log p_t)[1] - f_3[1]| \\ &\leq \frac{K}{\sigma_t^2} (\sqrt{d \log N} + 1) p_t \left| \frac{1}{p_t} - \frac{1}{f_{1,\text{clip}}} \right| + \frac{1}{f_{1,\text{clip}}} \left| \frac{(\nabla \sigma_t p_t)[1] - f_2[1]}{\sigma_t} \right| \end{aligned}$$

 $\lesssim \frac{1}{f_{1,\text{clip}}} \left( \frac{1}{\sigma_t^2} \sqrt{\log N} |p_t - f_{1,\text{clip}}| + \frac{(\nabla \sigma_t p_t)[1] - f_2[1]}{\sigma_t} \right).$ (By dropping Constant Terms)

From Lemma C.5, we have

$$|f_1 - p_t| \le BN^{-\beta} \log^{\frac{d_x + k_1}{2}} N.$$

We set  $\epsilon_{\text{low}} = C_3 N^{-\beta} \log^{(d_x+k_1)/2} N \le p_t$  such that  $f_1 \ge p_t/2$  by the choice of constant  $C_3$ .

$$\begin{split} &|(\nabla \log p_t)[1] - f_3[1]| \\ \lesssim \frac{1}{p_t} \left( \frac{1}{\sigma_t^2} \sqrt{\log N} |p_t - f_{1,\text{clip}}| + \frac{(\nabla \sigma_t p_t)[1] - f_2[1]}{\sigma_t} \right) \quad \text{(By the choice of } \epsilon_{\text{low}} \text{)} \\ \lesssim \frac{B}{\sigma_t^2 p_t} N^{-\beta} (\log N)^{\frac{d_x + k_1 + 1}{2}}. \quad \text{(By Lemma C.3 and Lemma C.4)} \end{split}$$

By the symmetry of each coordinate, the infinity bound for the score holds as well:

$$\left\|\nabla \log p_t - f_3\right\|_{\infty} \lesssim \frac{B}{\sigma_t^2 p_t} N^{-\beta} (\log N)^{\frac{d_x + k_1 + 1}{2}}.$$
 (C.8)

## – Step B: Approximate $f_3$ with Transformer $\mathcal{T}_{score}$ .

In this step, we utilize transformers to approximate  $f_3$  to an accuracy of order  $N^{-\beta}$ such that it aligns with the error order in (C.8).

Since  $f_3$  is the minimum between two components, we approximate each of them as follows.

\* Step B.1: Approximate  $\frac{1}{\sigma_t} \cdot \frac{f_2}{f_{1, \text{clip}}}$ . First, we utilize  $\mathcal{T}_{f_1}, \mathcal{T}_{f_2}$  and  $\mathcal{T}_{\sigma, 1}$  in Lemma C.5, Lemma C.6, and Lemma C.11 for  $f_1$ ,  $f_2$ , and  $\sigma_t$  respectively. This gives error  $\epsilon_{f_1}$ ,  $\epsilon_{f_2}$  and  $\epsilon_{\sigma,1}$ , and we address the clipping of  $f_1$  in later paragraph. Next, We utilize  $\mathcal{T}_{rec,1}$  and  $\mathcal{T}_{rec,2}$  in Lemma C.9 for the approximation of the inverse of  $f_1$  and  $\sigma_t$ .

This gives error

$$\left|\mathcal{T}_{\text{rec},1} - \frac{1}{f_1}\right| \le \epsilon_{\text{rec},1} + \frac{|\mathcal{T}_{f_1} - f_1|}{\epsilon_{\text{rec},1}^2} \le \epsilon_{\text{rec},1} + \frac{\epsilon_{f_1}}{\epsilon_{\text{rec},1}^2}$$

and

$$\left|\mathcal{T}_{\text{rec},2} - \frac{1}{\sigma_t}\right| \leq \epsilon_{\text{rec},2} + \frac{|\mathcal{T}_{\sigma,1} - \sigma_t|}{\epsilon_{\text{rec},2}^2} \leq \epsilon_{\text{rec},2} + \frac{\epsilon_{\sigma,1}}{\epsilon_{\text{rec},2}^2}$$

Note that all the approximation error propagates to the next approximation. Next, we utilize  $\mathcal{T}_{\text{mult},1}$  in Lemma C.8 for the approximation of the product of  $f_1^{-1}$ ,  $f_2$  and  $\sigma_t^{-1}$ . This gives error of

$$\left| \mathcal{T}_{\text{mult},1} - \frac{f_2}{\sigma_t f_1} \right| \le \epsilon_{\text{mult},1} + 3K_2^2 \max\left( \epsilon_{\text{rec},1} + \frac{\epsilon_{f_1}}{\epsilon_{\text{rec},1}^2}, \epsilon_{f_2}, \epsilon_{\text{rec},2} + \frac{\epsilon_{\sigma,1}}{\epsilon_{\text{rec},2}^2} \right) = \epsilon_{\text{mult},1} + 3K_2^2 \epsilon_1,$$

and  $K_2$  is a positive constant. From Lemma C.8 we require that  $[-K_2, K_2]$  covers the domain for all of  $f_1^{-1}$ ,  $f_2$  and  $f_{\sigma}^{-1}$ .

To be more specific, we reiterate three facts that determines the choice of  $K_2$ . • Recall that in the **Step A.**, we set  $f_{1,clip} = \{f_1, \epsilon_{low}\}$ . • Lemma C.12 states  $K(C_x\sqrt{d_x \log N} + 1)/\sigma_t^2$  is the  $\ell_\infty$  bound on the score. • The maximum value of  $\sigma_t^{-1}$  happens at  $t = t_0$ . As a result, we set  $K_2$  as  $K_2 = \max\left(\frac{1}{\epsilon_{low}}, \frac{K}{\sigma_{t_0}}(C_x\sqrt{d_x \log N} + 1), \frac{1}{\sigma_{t_0}}\right)$ . By the earlier choice of  $\epsilon_{low}$ , we have  $\epsilon_{low}^{-1} = \mathcal{O}(N^\beta \log N^{-(d_x+k_1)/2})$ , and next we expand  $\sigma_{t_0}$ .  $\sigma_{t_0} = \sqrt{1 - \exp(N^{-C_\sigma})} = 1 - (1 - \mathcal{O}(N^{-C_\sigma}))$ .

Therefore we have  $\sigma_{t_0}^{-1} = \mathcal{O}(N^{C_{\sigma}})$ . Putting all together, we have

$$K_2 = \mathcal{O}\left(N^{\beta + C_{\sigma}} \log^{-\frac{d_x + \beta}{2}} N\right), \tag{C.9}$$

where we use  $k_1 \leq \beta$ .

\* Step B.2 : Approximate  $K(C_x\sqrt{d_x \log N} + 1)/\sigma_t^2$ .

We invoke  $\mathcal{T}_{\sigma,2}$  in Lemma C.11 for the approximation of  $\sigma_t$ , and this gives error  $\epsilon_{\sigma,2}$ .

Next, we utilize  $T_{rec,3}$  in Lemma C.8 for the approximation of the inverse of  $\sigma_t$ . This gives error

$$\left|\mathcal{T}_{\text{rec},3} - \frac{1}{\sigma_t}\right| \le \epsilon_{\text{rec},3} + \frac{\left|\mathcal{T}_{\sigma,3} - \sigma_t\right|}{\epsilon_{\text{rec},3}^2} \le \epsilon_{\text{rec},3} + \frac{\epsilon_{\sigma,2}}{\epsilon_{\text{rec},3}^2}$$

Next, we utilize  $T_{\text{mult},2}$  for the approximation of the square of  $\sigma_t^{-1}$ . This gives error of

$$\left|\mathcal{T}_{\text{mult},2} - \left(\frac{1}{\sigma_t}\right)^2\right| \le \epsilon_{\text{mult},2} + 2K_1 \left(\epsilon_{\text{rec},3} + \frac{\epsilon_{\sigma,2}}{\epsilon_{\text{rec},3}^2}\right).$$

and  $K_1$  is constant to be chosen such that  $\sigma_t \in [-K_1, K_1]$ . With the same argument for  $K_2$ , it suffices to take  $\mathcal{O}(\sigma_t^{-1})$ :

$$K_1 = \mathcal{O}\left(N^{C_\sigma}\right). \tag{C.10}$$

\* Step B.3: Error Bound on Every Approximation Combined. Combining Step B.1 and Step B.2, the total error is bounded by

$$\epsilon_{\text{score}} \leq \max\left(\epsilon_{\text{mult},2} + 2K_1\left(\epsilon_{\text{rec},3} + \frac{\epsilon_{\sigma,2}}{\epsilon_{\text{rec},3}^2}\right), \epsilon_{\text{mult},1} + 3K_2^2\epsilon_1\right).$$

The goal is to guarantee the final error  $\epsilon_{\text{score}} \leq N^{-\beta}$  such that it matches the order of the approximation error in **Step A.** We list all the error choice to achieve the goal.<sup>1</sup>

· For the Error of the First Two Inverse Operators:

$$\epsilon_{\mathrm{rec},1}, \epsilon_{\mathrm{rec},2} = \mathcal{O}\left(N^{-(3\beta+2C_{\sigma})}(\log N)^{(d_x+\beta)}\right).$$

<sup>&</sup>lt;sup>1</sup>Further details regarding the choice of each one of  $\epsilon$  are in Appendix F.4 of (Fu et al., 2024b).

 For the Error of the Last Inverse Operator:  $\epsilon_{\mathrm{rec},3} = \mathcal{O}\left(N^{-(\beta+2C_{\sigma})}\right).$ • For the Error of  $f_1$ :  $\epsilon_{f_1} = \mathcal{O}\left(N^{-(9\beta + 6C_{\sigma})} (\log N)^{3(d_x + \beta)}\right).$ • For the Error of  $f_2$ :  $\epsilon_{f_2} = \mathcal{O}\left(N^{-(3\beta+2C_{\sigma})}(\log N)^{(d_x+\beta)}\right).$ · For the Error of the First Variance:  $\epsilon_{\sigma,1} = \mathcal{O}\left(N^{-(9\beta + 6C_{\sigma})} (\log N)^{3(d_x + \beta)}\right).$ · For the Error of the Second Variance:  $\epsilon_{\sigma,2} = \mathcal{O}\left(N^{-(7\beta+5C_{\sigma})}(\log N)^{2(d_x+\beta)}\right).$ · For the Error of the Two Product Operators:  $\epsilon_{\text{mult},1}, \epsilon_{\text{mult},2} = \mathcal{O}(N^{-\beta}).$ The above error choice renders  $\epsilon_{\text{score}} \leq N^{-\beta}$ . Therefore we conclude that there exist a transformer  $\mathcal{T}_{score} \in \mathcal{T}_{R}^{h,s,r}$  such that  $\left\|\mathcal{T}_{\text{score}}(x, y, t) - f_3(x, y, t)\right\|_{\infty} \le N^{-\beta}.$ (C.11)Combining (C.8) and (C.11) we obtain  $\|\nabla \log p_t - \mathcal{T}_{\text{score}}(x, y, t)\|_{\infty} \lesssim \frac{1}{p_t} \frac{B}{\sigma_t^2} N^{-\beta} (\log N)^{\frac{d_x + k_1 + 1}{2}}.$ We have completed the first part of the proof. We next give the norm bounds for the transformer parameters. Specifically, we select the parameter bounds that are consistent across all operations. including Lemma C.5, Lemma C.6, Lemma C.8, Lemma C.9 and Lemma C.11. Step C: Transformer Parameter Bound. Our result highlights the influence of N under varying  $d_x$ . Therefore, for the trans-former parameter bounds, we keep terms with  $d_x, d, L$  appearing in the exponent of N and  $\log N$ . Note that the following parameter selection is based on high-dimensional case where  $\log N$  term dominates N term. \* Parameter Bound on  $W_Q$  and  $W_K$ . Given error  $\epsilon$ , the bound on each operation follows: • For  $\epsilon_{f_1}$ : By Lemma C.5, we have  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta + 6C_{\sigma}) \cdot \frac{2dL + 4d + 1}{d}} \cdot (\log N)^{-3(d_x + \beta) \cdot \frac{2dL + 4d + 1}{d}}\right).$ • For  $\epsilon_{f_2}$ : By Lemma C.6, we have  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_{\sigma})\cdot\frac{2dL+4d+1}{d}} \cdot (\log N)^{-(d_x+\beta)\cdot\frac{2dL+4d+1}{d}}\right).$ 

972 • For  $\epsilon_{mult,1}$ : By Lemma C.8 with m = 3, we have 973 974  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_K||_{2,\infty} = \mathcal{O}(N^{7\beta}).$ 975 976 • For  $\epsilon_{mult,2}$ : By Lemma C.8 with m = 2, we have 977  $||W_Q||_2, ||W_K||_2, ||W_Q||_2, \dots, ||W_K||_2 = \mathcal{O}(N^{5\beta}).$ 978 979 • For  $\epsilon_{rec,1}$ ,  $\epsilon_{rec,2}$ : By Lemma C.9, we have 980 981  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta + 6C_{\sigma})} (\log N)^{-3(d_x+\beta)}\right).$ 982 983 • For  $\epsilon_{rec,3}$ : By Lemma C.9, we have 984 985  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_K||_{2,\infty} = \mathcal{O}\left(N^{(3\beta+6C_{\sigma})}\right).$ 986 987 988 • For  $\epsilon_{\sigma_1}$ : By Lemma C.11, we have 989  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_Q\|_{2,\infty} = \mathcal{O}\left(N^{(27\beta + 18C_{\sigma})}(\log N)^{-9(d_x+\beta)}\right).$ 991 992 • For  $\epsilon_{\sigma_2}$ : By Lemma C.11, we have 993  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_Q\|_{2,\infty} = \mathcal{O}\left(N^{(21\beta+15C_{\sigma})}(\log N)^{-6(d_x+\beta)}\right).$ 994 995 996 We select the largest parameter bound from  $\epsilon_{mult,1}$  and  $\epsilon_{rec,3}$  that remains valid 997 across all other approximations. That is, we take  $N^{(7\beta+6C_{\sigma})}$  as the upper-bound. 998 999 \* Parameter Bound on  $W_O$  and  $W_V$ . Given error  $\epsilon$ , the bound on each operation follows: 1002 • For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1004  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}\left(N^{-\frac{(9\beta+6C_{\sigma})}{d}}(\log N)^{\frac{3(d_x+\beta)}{d}}\right).$ • For  $\epsilon_{f_2}$ : By Lemma C.6, we have 1008  $\left\|W_O\right\|_2, \left\|W_O\right\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(3\beta+2C_{\sigma})}{d}}(\log N)^{\frac{(d_x+\beta)}{d}}\right).$ 1010 1011 • For  $\epsilon_{\text{mult},1}$ : By Lemma C.8 with m = 3, we have 1012 1013  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}(N^{-3\beta}).$ 1014 1015 • For  $\epsilon_{\text{mult},2}$ : By Lemma C.8 with m = 2, we have 1016 1017  $||W_O||_2, ||W_O||_2 = \mathcal{O}(N^{-2\beta}).$ • For  $\epsilon_{rec,1}$ ,  $\epsilon_{rec,2}$ : By Lemma C.9, we have 1021  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+C_{\sigma})}(\log N)^{d_x+\beta}\right).$ 1023 • For  $\epsilon_{rec,3}$ : By Lemma C.9, we have 1024 1025  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}\left(N^{-(\beta+2C_{\sigma})}\right).$ 

1026 • For  $\epsilon_{\sigma_1}$ : By Lemma C.11, we have 1027 1028  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}\left(N^{-(9\beta+6C_{\sigma})}(\log N)^{3(d_x+\beta)}\right).$ 1029 1030 • For  $\epsilon_{\sigma_2}$ : By Lemma C.11, we have 1031 1032  $||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}\left(N^{-(7\beta+5C_{\sigma})}(\log N)^{2(d_x+\beta)}\right).$ 1033 1034 Note that only  $\epsilon_{f_1}$  and  $\epsilon_{f_2}$  involve the reshape operation. From Lemma B.2, we take 1035  $\mathcal{O}(\sqrt{d})$  and  $\mathcal{O}(d) \|W_V\|_2$  and  $\|W_V\|_{2,\infty}$ . Moreover, We select the largest parameter bound from  $\epsilon_{rec,1}$  and  $\epsilon_{\sigma_1}$  that remains valid across all other approximations. That is, we take  $N^{-(3\beta+6C_{\sigma})}(\log N)^{3(d_x+\beta)}$  as the upper-bound. 1039 \* Parameter Bound on W<sub>1</sub>. Given error  $\epsilon$ , the bound on each operation follows: 1041 • For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1043  $\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_{\sigma})}{d}}(\log N)^{-\frac{3(d_x+\beta)}{d}} \cdot (\log N)\right).$ 1045 1046 • For  $\epsilon_{f_2}$ : By Lemma C.6, we have 1047 1048  $\|W_1\|_2, \|W_1\|_{2\infty} = \mathcal{O}\left(N^{\frac{(3\beta+2C_{\sigma})}{d}}(\log N)^{-\frac{(d_x+\beta)}{d}} \cdot (\log N)\right).$ 1049 1050 • For  $\epsilon_{\text{mult},1}$ : By Lemma C.8 with m = 3 and  $C = K_2$  in (C.9), we have 1051 1052  $\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(K_2 \cdot N^{3\beta}\right) = \mathcal{O}\left(N^{(4\beta + C_{\sigma})} (\log N)^{-\frac{1}{2}(d_x + \beta)}\right).$ • For  $\epsilon_{\text{mult},2}$ : By Lemma C.8 with m = 2 and  $C = K_1$  in (C.10), we have 1056  $||W_1||_2, ||W_1||_{2,\infty} = \mathcal{O}(K_1 \cdot N^{2\beta}) = \mathcal{O}(N^{(2\beta + C_{\sigma})}).$ 1058 • For  $\epsilon_{rec,1}$ ,  $\epsilon_{rec,2}$ : By Lemma C.9, we have  $||W_1||_2, ||W_1||_{2,\infty} = \mathcal{O}\left(N^{(6\beta+4C_{\sigma})}(\log N)^{-2(d_x+\beta)}\right).$ 1062 • For  $\epsilon_{rec,3}$ : By Lemma C.9, we have 1064 1065  $||W_1||_2, ||W_1||_{2,\infty} = \mathcal{O}\left(N^{(2\beta+4C_{\sigma})}\right)$ 1067 • For  $\epsilon_{\sigma_1}$ : By Lemma C.11, we have 1068 1069  $||W_1||_2, ||W_1||_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_{\sigma})}(\log N)^{-3(d_x+\beta)} \cdot \log N\right).$ 1070 1071 • For  $\epsilon_{\sigma_2}$ : By Lemma C.11, we have  $\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+5C_{\sigma})}(\log N)^{-2(d_x+\beta)} \cdot \log N\right).$ 1075 We select the largest parameter bound from  $\epsilon_{rec,3}$  that remains valid across all other approximations. That is, we take  $N^{(2\beta+4C_{\sigma})}$  as the upper-bound. 1077 \* Parameter Bound for W<sub>2</sub>. 1079 Given error  $\epsilon$ , the bound on each operation follows:

1080 • For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1082  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_{\sigma})}{d}}(\log N)^{-3\frac{(d_x+\beta)}{d}}\right).$ • For  $\epsilon_{f_2}$ : By Lemma C.6, we have For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1085  $\left\|W_2\right\|_2, \left\|W_2\right\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+2C_{\sigma})}{d}} (\log N)^{-\frac{(d_x+\beta)}{d}}\right).$ 1087 1088 • For  $\epsilon_{\text{mult},1}$ : By Lemma C.8 with m = 3, we have 1089  $||W_2||_2, ||W_2||_2 = \mathcal{O}(N^{3\beta}).$ • For  $\epsilon_{\text{mult},2}$ : By Lemma C.8 with m = 2, we have 1093  $||W_2||_2, ||W_2||_2 = \mathcal{O}(N^{2\beta}).$ 1095 • For  $\epsilon_{rec,1}$ ,  $\epsilon_{rec,2}$ : By Lemma C.9, we have  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_{\sigma})}(\log N)^{-(d_x+\beta)}\right).$ 1099 1100 • For  $\epsilon_{rec,3}$ : By Lemma C.9, we have 1101 1102  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{(\beta+2C_{\sigma})}\right).$ 1103 1104 • For  $\epsilon_{\sigma_1}$ : By Lemma C.11, we have 1105  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_{\sigma})}(\log N)^{-3(d_x+\beta)}\right).$ 1106 1107 1108 • For  $\epsilon_{\sigma_2}$ : By Lemma C.11, we have 1109 1110  $\|W_2\|_2, \|W_2\|_2 = \mathcal{O}\left(N^{(7\beta+5C_{\sigma})}(\log N)^{-2(d_x+\beta)}\right).$ 1111 1112 We select the largest parameter bound from  $\epsilon_{mult,1}$  and  $\epsilon_{rec,3}$  that remains valid 1113 across all other approximations. That is, we take  $N^{(3\beta+2C_{\sigma})}$  as the upper-bound. 1114 **Parameter Bound for** *E*. \* 1115 Since only  $\epsilon_{f_1}$  and  $\epsilon_{f_2}$  involve the reshape operation. From Lemma B.2, we take 1116  $\mathcal{O}(d^{\frac{1}{2}}L^{\frac{3}{2}}) \text{ for } \left\| E^{\top} \right\|_{2,\infty}.$ 1117 1118 By integrating results above, we derive the following parameter bounds for the transformer network, ensuring valid approximation across all nine approximations. 1119 1120 1121  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+6C_{\sigma})}\right);$ 1122 1123  $\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+6C_{\sigma})}(\log N)^{3(d_x+\beta)}\right);$ 1124 1125  $\|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|E^{\top}\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);$ 1126  $\left\|W_{1}\right\|_{2}, \left\|W_{1}\right\|_{2,\infty} = \mathcal{O}\left(N^{\left(2\beta+4C_{\sigma}\right)}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_{t}^{2}\right);$ 1127 1128  $||W_2||_2, ||W_2||_2 = \mathcal{O}\left(N^{(3\beta+2C_{\sigma})}\right).$ 1129 1130 1131 The last network output bound  $C_{\mathcal{T}} = \mathcal{O}(\sqrt{d_x \log N}/\sigma_t^2)$  follows the entry-wise mini-1132 mum bounds  $K(C_x\sqrt{d\log N}+1)/\sigma_t^2$  in  $\ell_\infty$  distance by Lemma C.12. 1133

This completes the proof.

#### 1134 C.2 MAIN PROOF OF THEOREM 3.1 UNDER GENERIC ASSUMPTION 1135

1136 In Lemma C.13, we establish the score approximation with transformer that incorporates every 1137 essential components and encodes the Hölder smoothness in the final result. However, it is only valid within the input domain  $[C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x} \times [0, 1]^{d_y}$ , and we also excludes region  $p_t < \epsilon_{\text{low}}$ 1138 where the problem of score explosion remains unaddressed. 1139

1140 To combat this, we introduce two additional lemmas. The first lemma gives us the error caused by 1141 the truncation of  $\mathbb{R}^{d_x}$  within a radius  $R_1$  in  $\ell_2$  distance. 1142

Lemma C.14 (Truncate x for Score Function, Lemma A.1 of (Fu et al., 2024b)). Under generic 1143 Assumption 3.1. For any  $R_1 > 1$ , y, t > 0 we have 1144

1145 1146

1147

1148

1150

1149

$$\begin{split} &\int_{\|x\|_{\infty} \ge R_{1}} p_{t}(x|y) dx \le R_{1} \exp\left(-C_{2}^{\prime}R_{1}^{2}\right), \\ &\int_{\|x\|_{\infty} \ge R_{1}} \|\nabla \log p_{t}(x|y)\|_{2}^{2} p_{t}(x|y) dx \le \frac{R_{1}^{3}}{\sigma_{t}^{4}} \exp\left(-C_{2}^{\prime}R_{1}^{2}\right), \end{split}$$

where  $C'_2 = C_2/(2 \max(C_2, 1))$ . 1151

1152 **Remark C.6.** Because we only impose assumption on the light tail property of the conditional 1153 distribution in generic Assumption 3.1, the unboundedness of x necessitates a truncation for integrals 1154 regarding x, or else the result would diverge. 1155

Furthermore, we address the explosion of score function with the second lemma. 1156

1157 Lemma C.15 (Lemma A.2 of (Fu et al., 2024b)). Under generic Assumption 3.1. For any 1158  $R_2, y, \epsilon_{\text{low}} > 0$  we have 1159

1160 1161

1162

1163 1164

1167

1169

1170 1171

1172 1173

1174

1175

1176 1177

1178 1179

1180

1181

 $\int_{\|x\|_{\infty} \leq R_2} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \cdot p_t(x|y) dx \leq R_2^{d_x} \epsilon_{\text{low}},$  $\int_{\|x\|_{-} \leq R_{2}} \mathbb{1}\{|p_{t}(x|y)| < \epsilon_{\text{low}}\} \cdot \|\nabla \log p_{t}(x|y)\|_{2}^{2} p_{t}(x|y) dx \leq \frac{1}{\sigma_{t}^{4}} R_{2}^{d_{x}+2} \epsilon_{\text{low}}.$ 

1165 **Remark C.7.** Recall that the score function has the form  $\nabla \log p_t(x|y) = \nabla p_t(x|y)/p_t(x|y)$ . It is 1166 essential to set a threshold for  $p_t(x|y)$  prevents the explosion of the score function.

We begin the proof of Theorem 3.1. 1168

*Proof Sketch of Theorem 3.1.* In the following proof, we give error bound for the three terms:

• (A.1): The approximation for  $||x||_{\infty} > R_1$ .

This step controls the error from truncation of  $\mathbb{R}^{d_x}$  with radius  $R_1$  in  $\ell_2$  distance. We approximate the error with Lemma C.14

• (A.2): The approximation for  $\mathbf{1}\{p_t(x|y) < \epsilon_{\text{low}}\}\$  and  $\{\|x\|_{\infty} \le R_1\}$ .

This step controls the error from setting a threshold to prevent score explosion within the bounded domain  $||x||_{\infty} \leq R_1$ . We approximate the error with Lemma C.15.

• (A.3) The approximation for  $1\{p_t(x|y) \ge \epsilon_{low}\}$  and  $\{\|x\|_{\infty} \le R_1\}$ .

With previous two steps ensuring the bounded domain and preventing the divergence of score function, we approximate with Lemma C.13.

*Proof of Theorem 3.1.* We apply  $N = N^{1/(d_x+d_y)}$  in Lemma C.13. Throughout the proof, we use 1186 N as a notational simplification, with the understanding that N represents  $N^{1/(d_x+d_y)}$  in full form. 1187 At the end of the proof we replace N by  $N^{1/(d_x+d_y)}$ .

 $\int_{\mathbb{R}^{d_x}} \|s(x,y,t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx$ 

 $=\underbrace{\int_{\|x\|_{\infty}>\sqrt{\frac{2\beta}{C_{2}^{\prime}}\log N}}\|s(x,y,t)-\nabla\log p_{t}(x|y)\|_{2}^{2}\cdot p_{t}(x|y)\mathrm{d}x}_{(A_{1})},$ 

To begin with, we set  $R_1 = R_2 = \sqrt{2\beta \log N/C_2'}$  in Lemma C.14 and Lemma C.15, and we expand the target into three parts  $(A_1)$ ,  $(A_2)$ , and  $(A_3)$ :

$$+ \underbrace{\int_{\|x\|_{\infty} \leq \sqrt{\frac{2\beta}{C_{2}'} \log N}} \mathbb{1}\{|p_{t}(x|y)| < \epsilon_{\text{low}}\} \|s(x,y,t) - \nabla \log p_{t}(x|y)\|_{2}^{2} \cdot p_{t}(x|y) dx}_{(A_{2})}}_{(A_{2})} + \underbrace{\int_{\|x\|_{\infty} \leq \sqrt{\frac{2\beta}{C_{2}'} \log N}} \mathbb{1}\{|p_{t}(x|y)| \geq \epsilon_{\text{low}}\} \|s(x,y,t) - \nabla \log p_{t}(x|y)\|_{2}^{2} \cdot p_{t}(x|y) dx}_{(A_{3})}}_{(A_{3})}$$

1207 We derive the bound for  $(A_1), (A_2), (A_3)$  and combine these results.

• Bounding (A<sub>1</sub>). We apply Lemma C.14. Note that we have  $||s(x, y, t)||_{\infty} \leq \sqrt{\log N}/\sigma_t^2$  from the construction of the score estimator in Lemma C.13.

• Bounding (A<sub>2</sub>). We apply Lemma C.15. Note that we set  $\epsilon_{\text{low}} = C_3 N^{-\beta} (\log N)^{(d_x+k_1)/2}$  in Lemma C.13.

$$\int_{\|x\|_{\infty} \le \sqrt{\frac{2\beta}{C_{2}'} \log N}} \mathbb{1}\{|p_{t}(x|y)| < \epsilon_{\text{low}}\} \|s(x,y,t) - \nabla \log p_{t}(x|y)\|_{2}^{2} \cdot p_{t}(x|y) dx$$

(By expanding the  $\ell_2$  norm)

Combining  $(A_1)$ ,  $(A_2)$  and  $(A_3)$ , we have

$$\int_{\mathbb{R}^d} \left\| s(x,y,t) - \nabla \log p_t(x|y) \right\|_2^2 p_t(x|y) \mathrm{d}x$$

1296  
1297 
$$\lesssim \frac{N^{-2\beta} (\log N)^{\frac{3}{2}}}{\sigma^4} + \frac{\epsilon_{\text{low}} (\log N)^{\frac{d_x+2}{2}}}{\sigma^4} + \frac{B^2 d_x}{\sigma^4 c} N^{-2\beta} (\log N)^{\frac{3d_x}{2} + k_1 + 1}$$

$$\underbrace{\frac{\sigma_{t}^{4}}{\sigma_{t}^{4}}}_{(A_{1})} + \underbrace{\frac{\sigma_{t}^{4}}{\sigma_{t}^{4}}}_{(A_{2})} + \underbrace{\frac{\sigma_{t}^{4}}{\sigma_{t}^{4}}}_{(A_{3})} + \underbrace{\frac{\sigma_{t}^{4}}{\sigma_{t}^{4}}}_{(A_{3})} N^{-2\beta} (\log N)^{\frac{1}{2} + k_{1} + 1}}_{(A_{3})}$$

By replacing  $\epsilon_{\text{low}}$  with  $C_3 N^{-\beta} (\log N)^{d_x + k_1/2}$  and using the relation  $k_1 \leq \beta$ ,<sup>2</sup> we obtain

$$\int_{\mathbb{R}^d} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx = \mathcal{O}\left(\frac{B^2}{\sigma_t^4} N^{-\beta} (\log N)^{d_x + \frac{\beta}{2} + 1}\right).$$

1306 Replacing N with  $N^{1/(d_x+d_y)}$  completes the first part of the proof.

1307 The transformer parameter norm bounds follow Lemma C.13, with the replacement of N with 1308  $N^{1/(d_x+d_y)}$  as well. Note that this results in  $t \in [N^{-C_\alpha/(d_x+d_y)}, C_\sigma/((d_x+d_y))\log N]$ . For better 1309 interpretation of the cutoff and early stopping time parameter, we reset  $C_\alpha$  as  $(d_x + d_y)C_\alpha$  and  $C_\sigma$ 1310 as  $(d_x + d_y)C_\sigma$  such that  $t \in [N^{-C_\alpha}, C_\sigma \log N]$ .

This completes the proof.

<sup>&</sup>lt;sup>2</sup>Recall the definition of the Hölder smoothness from Definition 3.1.

#### 1350 D **PROOF OF THEOREM 3.1 UNDER STRONGER ASSUMPTION** 1351

1352 We state the proof of Theorem 3.1 under stronger Hölder assumption as follows.

- Step 0. We decompose the density function and the score function under stronger Assumption 3.1. In Lemma D.1, we provide details regarding the decomposed form of the score function. We specify the upper and lower bound on h and  $\nabla h$  in Lemma D.2.
  - Step 1. Similar to the domain discretization in the proof of previous main result, we discretize the input domain of the decomposed density function in Lemma D.3.
  - Step 2. We construct polynomial approximation based on Taylor expansion of h and  $\nabla h$  in Lemmas D.4 and D.5. The approximation result captures the local Hölder smoothness, with improved precision relative to the analogous step in Lemma C.3 and Lemma C.4.
  - Step 3. We approximate h and  $\nabla h$  with transformer in Lemmas D.6 and D.7. In order to construct the score approximator with transformer, we approximate several additional algebraic operators with transformer in Lemma D.8, Lemma D.9 and Lemma D.10. We incorporate these results into a unified transformer architecture in Lemma D.11.

**Organization.** Appendix D.1 includes the four steps and auxiliary lemmas supporting our proof. Appendix D.2 includes the formal version and main proof of Theorem 3.1.

1370 D.1 AUXILIARY LEMMAS

1372 Step 0: Decompose the Score with Stronger Hölder Smoothness Assumption. We utilize the condition assumed in stronger Assumption 3.1 to achieve the decomposition. 1373

**Lemma D.1** (Lemma B.1 of Fu et al. (2024b)). Under stronger Assumption 3.1. The conditional distribution at time t has the following expression: 1376

$$p_t(x|y) = \frac{1}{(\alpha_t^2 + C_2 \sigma_t^2)^{d_x/2}} \exp\left(-\frac{C_2 ||x||_2^2}{2(\alpha_t^2 + C_2 \sigma_t^2)}\right) h(x, y, t).$$

1380 Moreover, the score function has the following expression:

$$\nabla \log p_t(x|y) = \frac{-C_2 x}{\alpha_t^2 + C_2 \sigma_t^2} + \frac{\nabla h(x, y, t)}{h(x, y, t)},$$

1384 1385 1386

1387

1388

1389

1396

1353

1354

1355

1356 1357

1358

1359

1363

1365

1367

1369

1371

1374

1375

1377

1381 1382

where 
$$h(x, y, t) = \int \frac{f(x_0, y)}{\hat{\sigma}_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|x_0 - \hat{\alpha}_t x\|^2}{2\hat{\sigma}_t^2}\right) \mathrm{d}x_0, \ \hat{\sigma}_t = \frac{\sigma_t}{(\alpha_t^2 + C_2 \sigma_t^2)^{1/2}}, \text{ and } \hat{\alpha}_t = \frac{\alpha_t}{\alpha_t^2 + C_2 \sigma_t^2}.$$

Next, we provide lemma that provides bound on h(x, y, t) and  $\nabla h(x, y, t)$  in Lemma D.1

Lemma D.2 (Lemma B.8 of (Fu et al., 2024b)). Under stronger Assumption 3.1, we have the following bounds for h(x, y, t) and  $\frac{\widehat{\sigma}_t}{\widehat{\alpha}_t} \nabla h(x, y, t)$ 

$$C_1 \le h(x, y, t) \le B, \quad \left\| \frac{\widehat{\sigma}_t}{\widehat{\alpha}_t} \nabla h(x, y, t) \right\|_{\infty} \le \sqrt{\frac{2}{\pi}} B_t$$

where  $C_1$  and B are the hyperparameters of  $\mathcal{H}^{\beta}(\mathbb{R}^{d_x} \times [0,1]^{d_y}, B)$  in stronger Assumption 3.1. 1395

**Remark D.1** (Bound on h and  $\nabla h$ ). We reiterate that Lemma D.2 drives the key distinction between the analyses in Theorem 3.1 and Theorem 3.1 under stronger assumption. Specifically, in 1398 Appendix C.2, the decomposed term containing the threshold  $\epsilon_{low}$  results in lower approximation rate, 1399 while bounds on h and  $\nabla h$  eliminate the need of the threshold with h's lower bound  $C_1$ , rendering 1400 faster approximation rate. 1401

1402 Step 1: Discretize  $\mathbb{R}^{d_x} \times [0,1]^{d_y}$  for h(x,y,t). This step parallels Lemma C.1; however, the 1403 discretization differs due to the structure of h.

1407 Lemma D.3 (Clipping Integral, Lemma B.10 of Fu et al. (2024b)). Under stronger Assumption 3.1. Consider any integer vector  $\kappa \in \mathbb{Z}_+^{d_x}$  with  $\|\kappa\|_1 \leq n$ . There exists a constant  $C(n, d_x)$ , such that for any  $x \in \mathbb{R}^{d_x}$  and  $0 < \epsilon \leq 0.99$ , it holds

$$\int_{\mathbb{R}^{d_x} \setminus B_x} \left| \left( \frac{\widehat{\alpha}_t x_0 - x}{\widehat{\sigma}_t} \right)^{\kappa} \right| \cdot p(x_0 | y) \cdot \frac{1}{\widehat{\sigma}_t^d (2\pi)^{d/2}} \exp\left( -\frac{\left\| \widehat{\alpha}_t x_0 - x \right\|^2}{2\widehat{\sigma}_t^2} \right) \mathrm{d}x_0 \le \epsilon, \qquad (D.1)$$
where  $\left( \frac{\widehat{\alpha}_t x_0 - x}{\widehat{\sigma}_t} \right)^{\kappa} \coloneqq \left( \left( \frac{\widehat{\alpha}_t x_0 [1]_1 - x[1]}{\widehat{\sigma}_t} \right)^{\kappa[1]}, \left( \frac{\widehat{\alpha}_t x_0 [2] - x[2]}{\widehat{\sigma}_t} \right)^{\kappa[2]}, \dots, \left( \frac{\widehat{\alpha}_t x_0 [d_x] - x[d_x]}{\widehat{\sigma}_t} \right)^{\kappa[d_x]} \right) \text{ and}$ 

$$B_x \coloneqq \left[ \widehat{\alpha}_t x - C(n, d) \widehat{\sigma}_t \sqrt{\log \epsilon^{-1}}, \widehat{\alpha}_t x + C(n, d) \widehat{\sigma}_t \sqrt{\log \epsilon^{-1}} \right]^{d_x}.$$

v

Step 2: Approximate h and  $\nabla h$  with Polynomials. Similar to the construction of the diffused local polynomials in Lemma C.5 and Lemma C.6, the following two lemmas render the first step approximation for h(x, y, t) and  $\nabla h(x, y, t)$  that captures the local smoothness.

**Lemma D.4** (Approximation with Diffused Local Polynomials, Lemma B.4 of (Fu et al., 2024b)). 1422 Under stronger Assumption 3.1. For sufficiently larger N > 0 and constant  $C_2$ , there exists a diffused 1423 local polynomial  $f_1(x, y, t)$  with at most  $N^{d+d_y}(d+d_y)^{k_1}$  monomials such that

$$|f_1(x, y, t) - h(x, y, t)| \lesssim BN^{-\beta} \log^{\frac{\kappa_1}{2}} N,$$

1427 for any  $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}, y \in [0, 1]^{d_y}$  and t > 0.

**Lemma D.5** (Counterpart of Lemma D.4, Lemma B.6 of (Fu et al., 2024b)). Under stronger Assumption 3.1. For sufficiently larger N > 0 and constant  $C_2$ , there exists a diffused local polynomial  $f_2(x, y, t) \in \mathcal{T}_R^{h,s,r}$  with at most  $N^{d_x+d_y}(d_x + d_y)^{k_1}$  monomials  $f_2[i](x, y, t)$  such that

$$\left|f_2[i](x,y,t) - \left(\frac{\widehat{\sigma}_t}{\widehat{\alpha}_t} \nabla h(x,y,t)\right)[i]\right| \lesssim B N^{-\beta} \log^{\frac{k_1+1}{2}} N,$$

for any  $x \in \mathbb{R}^{d_x}$ ,  $y \in [0, 1]^{d_y}$  and t > 0.

1438 Step 3: Approximate Diffused Local Polynomials and Algebraic Operators with Transformers. 1439 First, we apply the universal approximation theory of transformers to  $f_1$  and  $f_2$ . Second, we adopt a 1440 comparable approach to approximate the algebraic operators essential for the final score computation. 1441 Last, we introduce Lemma D.11 that outlines how these components fit into a single transformer 1442 architecture with a specified parameter configuration.

## • Step 3.1: Approximate the Diffused Local Polynomials $f_1$ and $f_2$ .

We invoke the universal approximation theorem of transformer Theorem B.1. We utilize network consisting of one transformer block and one feed-forward layer.

**Lemma D.6** (Approximate Scalar Polynomials with Transformers). Under stronger Assumption 3.1. Consider the diffused local polynomial  $f_1$  in Lemma D.4. For any  $\epsilon > 0$ , there exists a transformer  $\mathcal{T}_{f_1} \in \mathcal{T}_R^{h,s,r}$ , such that for any  $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}, y \in [0, 1]^{d_y}$  and  $t \in [N^{-C_\sigma}, C_\alpha \log N]$ , it holds

$$|f_1(x, y, t) - \mathcal{T}_{f_1}(x, y, t)[d_x]| \le \epsilon,$$

The parameter bounds in the transformer network class follows Lemma C.5.

**Lemma D.7** (Approximate Vector-Valued Polynomials with Transformers). Under stronger Assumption 3.1 and consider  $f_2(x, y, t) \in \mathbb{R}^{d_x}$  in Lemma D.5. For any  $\epsilon > 0$ , there exists a

1458	transformer $\mathcal{T}_{f_2} \in \mathcal{T}_{P}^{h,s,r}$ such that
1460	
1461	$\left\ f_2(x, y, t) - \mathcal{T}_{f_2}\right\ _{\infty} \le \epsilon,$
1462	for any $x \in [-C_{\sigma}\sqrt{\log N} C_{\sigma}\sqrt{\log N}]^{d_x}$ $y \in [0, 1]^{d_y}$ and $t \in [N^{-C_{\sigma}} C_{\sigma} \log N]$ The
1463	parameter bounds in the transformer network class follows Lemma C.5.
1464	Stop 3.2: Approximate Algebraic Operators with Transformers
1465	Next we introduce lemmas regarding the function of time. These are also key components
1467	to the proof of Theorem D.1.
1468	
1469	<b>Lemma D.8</b> (Approximation of $\alpha^2$ with Transformer). For $t \in [t_0, T]$ with $t_0 < 1$ , there
1470	exists Transformer $\mathcal{T}_{\alpha^2}(t) \in \mathcal{T}_R^{n,s,v}$ such that
1471	$ \mathcal{T}_{\epsilon^2} - \alpha^2  < \epsilon_{\widehat{\alpha}}.$
1472	$ \alpha - \alpha  = \alpha$
1473	The parameter bounds in the Transformer network class follow Lemma C.11.
1475	Also, we approximate $\hat{\alpha}$ and $\hat{\sigma}_t$ as well.
1476	Lemma D.0 (Approximation of $\hat{\alpha}$ with Transformer). Consider $\hat{\alpha}_{t} = \alpha_{t}$ for $t \in C$
1477	Lemma D.9 (Approximation of $\alpha$ with transformer). Consider $\alpha_t = \frac{1}{\alpha_t^2 + \tilde{C}_2 \sigma_t^2}$ , for $t \in C_2$
1478	$[t_0, T]$ with $t_0 < 1$ , there exists Transformer $\mathcal{T}_{\widehat{\alpha}}(t) \in \mathcal{T}_R^{n,o,r}$ such that
1479	$ \mathcal{T}_{\widehat{lpha}}-\widehat{lpha}  \leq \epsilon_{\widehat{lpha}}.$
1480	
1481	The parameter bounds in the transformer network class follow Lemma C.11.
1483	Lemma D 10 (Approximation of $\hat{\sigma}$ with Transformer) Consider $\hat{\sigma} = \sigma_t$ for
1484	<b>Lemma D.10</b> (Approximation of $\delta$ with Hansformer). Consider $\delta_t = \frac{1}{(\alpha_t^2 + C_2 \sigma_t^2)^{1/2}}$ , for
1485	$t \in [t_0, T]$ with $t_0 < 1$ , there exists Transformer $\mathcal{T}_{\widehat{\sigma}}(t) \in \mathcal{T}_R^{n,s,r}$ such that
1486	$ \mathcal{T}_{\widehat{\sigma}} - \widehat{\sigma}  \leq \epsilon_{\widehat{\sigma}}.$
1487	
1488	The parameter bounds in the transformer network class follow Lemma C.11.
1490	We have finished establishing the approximation with transformer for every key component for the proof of Theorem 3.1.
1491 1492	Step 3.3: Unified Transformer-Based Score Function Approximation.
1493	We introduce the counterpart of Lemma C.13. It is the core of the proof for Theorem 3.1.
1494	Lamma D 11 (Sector Approximation with Transformer). Under stronger Assumption 2.1
1495	<b>Lemma D.11</b> (Score Approximation with Transformer). Under stronger Assumption 3.1 For sufficiently large integer N there exists a mapping from transformer $\mathcal{T} = \mathcal{T}^{h,s,r}$
1496	such that
1497	
1498	$\left\  \mathcal{T}_{\text{score}} - \nabla \log h(x, y, t) + \frac{C_2 x}{2 + C_2 - 2} \right\  \le \frac{B}{2} N^{-\beta} (\log N)^{\frac{k_1 + 1}{2}},$
1500	$\  \qquad \alpha_t^2 + C_2 \sigma_t^2 \ _{\infty}  \sigma_t$
1501	for any $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}, y \in [0, 1]^{d_y}$ and $t \in [N^{-C_\sigma}, C_\alpha \log N]$ .
1502	The parameter bounds in the transformer network class satisfy
1503	$\begin{pmatrix} & (a a + a g + 2dL + 4d + 1) \end{pmatrix}$
1504	$\ W_Q\ _2, \ W_K\ _2, \ W_Q\ _{2,\infty}, \ W_K\ _{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_{\sigma})}\right);$
1505	$  W_V  _2 = \mathcal{O}(\sqrt{d});   W_V  _{2-\alpha} = \mathcal{O}(d);   W_O  _2,   W_O  _{2-\alpha} = \mathcal{O}(N^{-\beta});$
1507	$(1, 1) = \left( \frac{1}{2} \right) \left( \frac{1}{2} \right) \left( \frac{1}{2} \right) \left( \frac{3}{2} \right) \left( \frac{3}{2} \right) \left( \frac{3}{2} \right) \left( \frac{1}{2} \right) \left( $
1508	$\ W_1\ _{2}, \ W_1\ _{2,\infty} = \mathcal{O}\left(N^{\frac{1}{2}p + 3C_{\sigma} + \frac{1}{2}} \cdot \log N\right); \ E^+\ _{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{1}{2}}\right);$
1509	$\ W_2\ _{2,\tau}\ W_2\ _{2,\tau} = \mathcal{O}\left(N^{4\beta+9C_{\sigma}+\frac{3C_{\alpha}}{2}}\right) : C_{\tau} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t\right).$
1510	$\int \frac{1}{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{\sqrt{$
1511	<b>Proof</b> Our proof follows the proof structure of (Eq. et al. 2024b) Proposition R 3)
	$1700$ . Our proof follows the proof structure of (1 $\mu$ et al., 20240, 110position D.3).

Recall the decomposed score function presented in **Step 0**, we establish the first-step approximator  $f_3$  with the form:

$$f_3(x, y, t) \coloneqq \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \cdot \frac{f_2(x, y, t)}{f_1(x, y, t)} - \frac{C_2 x}{\alpha_t^2 + C_2 \sigma_t^2}$$

We derive the error bound on the approximation of the first term containing Taylor polynomials in  $f_3$ . We incorporate second term containing the linear function in x into the the transformer architecture.

We proceed as follows:

- 1. Step A: Approximate  $\nabla \log p_t(x|y)$  with  $f_3$ .
- 2. **Step B:** Approximate  $f_3$  with  $\mathcal{T}_{\text{score}} \in \mathcal{T}_R^{h,s,r}$ .

3. Step C: Derive the final Parameter Configuration

## – Step A. Approximate Scroe Function with $f_3$ .

We first construct  $f_1(x, y, t)$  and  $f_2(x, y, t)$  from Lemma D.4 and Lemma D.5 to approximate h(x, y, t) and  $\nabla h(x, y, t)$  respectively.

From Lemma D.2, we have  $C_1 \le h \le B$  and  $\left\| \frac{\hat{\sigma}_t \nabla h}{\hat{\alpha}_t} \right\|_{\infty} \le \sqrt{\frac{2}{\pi}} B$ .

Next, by Lemma D.4 and Lemma D.5, we select a sufficiently large N such that  $\frac{C_1}{2} \leq f_1 \leq 2B$  and  $f_2 \leq B$ .

Without loss of generality, we begin by bounding the first coordinate of  $\nabla h$ , denoted as  $\nabla h[1]$ :

$$\begin{split} \left| \frac{\nabla h[1]}{h} - \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \frac{f_2[1]}{f_1} \right| &\leq \left| \frac{\nabla h[1]}{h} - \frac{\nabla h[1]]}{f_1} \right| + \left| \frac{\nabla h[1]}{f_1} - \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \frac{f_2[1]]}{f_1} \right|, \\ &\leq \left| \frac{\nabla h[1]]}{h \cdot f_1} \right| \left| f_1 - h \right| + \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \left| \frac{1}{f_1} \right| \left| f_2 - \frac{\widehat{\sigma}_t}{\widehat{\alpha}_t} \nabla h[1] \right| \right), \\ &\lesssim \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \left( \left| f_1 - h \right| + \left| f_2 - \frac{\widehat{\sigma}_t}{\widehat{\alpha}_t} \nabla h[1] \right| \right), \\ &\qquad (\text{By bounds on } h, \nabla h, f_1, f_2) \\ &\leq \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \left( BN^{-\beta} (\log N^{\frac{k_1}{2}} + BN^{-\beta} (\log N^{\frac{k_1+1}{2}}) \right), \\ &\qquad (\text{By Lemma D.4 and Lemma D.5)} \\ &\lesssim \frac{1}{\sigma_t} \left( BN^{-\beta} (\log N^{\frac{k_1+1}{2}}) \right). \end{split}$$

Note that in the last line, we utilize

$$\frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} = \frac{\alpha_t}{\sigma_t} \frac{1}{\sqrt{\alpha_t^2 + C_2 \sigma_t^2}} = \frac{1}{\sigma_t} \frac{1}{\sqrt{1 + C_2 \left(\sigma_t / \alpha_t\right)^2}} = \frac{1}{\sigma_t} \frac{1}{\sqrt{1 + C_2 \frac{\sigma_t^2}{1 - \sigma_t^2}}} = \mathcal{O}(\sigma_t^{-1}).$$

By the symmetry of each coordinate in  $\nabla h$ , we obtain the  $\ell_{\infty}$  bounds:

$$\left\|\frac{\nabla h(x,y,t)}{h(x,y,t)} - \frac{\widehat{\alpha}_t}{\widehat{\sigma}_t} \frac{f_2(x,y,t)}{f_1(x,y,t)}\right\|_{\infty} \lesssim \frac{B}{\sigma_t} N^{-\beta} (\log N)^{\frac{k_1+1}{2}}.$$
 (D.2)

– Step B. Approximate  $f_3$  with Transformer  $\mathcal{T}_{score}$ .

Next, we prove that there exist Transformer networks  $\mathcal{T}_{score} \in \mathcal{T}_{R}^{h,s,r}$  that approximates  $f_{3}(x, y, t)$  with error of order  $N^{-\beta}$ .

In the following, we construct a transformer approximating the two terms in  $f_3$ , and incorporate the result into a unified network architecture.

\* Step B.1: Approximation for  $\frac{\hat{\alpha}_t f_2}{\hat{\sigma}_t f_1}$ .

We utilize  $\mathcal{T}_{f_1}$ ,  $\mathcal{T}_{f_2}$ ,  $\mathcal{T}_{\hat{\alpha}}$  and  $\mathcal{T}_{\hat{\sigma}}$  in Lemma C.5, Lemma C.6, Lemma D.9 and Lemma D.10 to approximate each one of the component. This gives error  $\epsilon_{f_1}$ ,  $\epsilon_{f_2}$ ,  $\epsilon_{\hat{\alpha}}$  and  $\epsilon_{\hat{\sigma}}$  respectively.

Next we utilize  $\mathcal{T}_{\text{rec},2}$  and  $\mathcal{T}_{\text{rec},3}$  in Lemma C.9 for the approximation of the inverse of  $f_1$  and  $\hat{\sigma}_t$ . This gives error

$$\left|\mathcal{T}_{\text{rec},2} - \frac{1}{f_1}\right| \le \epsilon_{\text{rec},2} + \frac{|\mathcal{T}_{f_1} - f_1|}{\epsilon_{\text{rec},2}^2} \le \epsilon_{\text{rec},2} + \frac{\epsilon_{f_1}}{\epsilon_{\text{rec},2}^2}$$

and

$$\left|\mathcal{T}_{\text{rec},3} - \frac{1}{\widehat{\sigma}_t}\right| \leq \epsilon_{\text{rec},3} + \frac{|\mathcal{T}_{\widehat{\sigma}} - \widehat{\sigma}_t|}{\epsilon_{\text{rec},2}^2} \leq \epsilon_{\text{rec},3} + \frac{\epsilon_{\widehat{\sigma}}}{\epsilon_{\text{rec},3}^2}$$

Next we utilize  $\mathcal{T}_{\text{mult},1}$  in Lemma C.8 for the approximation of the product of  $f_1^{-1}$ ,  $f_2$ ,  $\hat{\alpha}_t$  and  $\hat{\sigma}_t^{-1}$ . This gives error

$$\begin{split} \left| \mathcal{T}_{\text{mult},1} - \frac{\widehat{\alpha}_t f_2}{\widehat{\sigma}_t f_1} \right| \\ \leq \epsilon_{\text{mult},1} + 4K_4^3 \max\left( \epsilon_{\text{rec},2} + \frac{\epsilon_{f_1}}{\epsilon_{\text{rec},2}^2}, \epsilon_{f_2}, \epsilon_{\widehat{\alpha}}, \epsilon_{\text{rec},3} + \frac{\epsilon_{\widehat{\sigma}}}{\epsilon_{\text{rec},3}^2} \right) \\ \underbrace{ = \epsilon_2} \\ \coloneqq \\ \vdots = \epsilon_2 \end{split} \coloneqq \\ \end{split}$$

and  $K_3$  is a positive constant.

From Lemma C.8, we require  $[-K_4, K_4]$  to cover the domain of  $f_1^{-1}$ ,  $f_2$ ,  $\hat{\alpha}$ , and  $\hat{\sigma}_t$ . Recall that we give the upper and lower bounds for  $f_1^{-1}$  and  $f_2$  in the beginning of **Step 1.** Thus, we set  $K_4 = \max(\hat{\sigma}_t^{-1}, \hat{\alpha}_t)$ .

To derive the asymptotic behavior of  $K_4$ , we set the positive constant  $C_2 = 2$  without loss of generality and note that the maximum occurs at  $t = t_0$ . We then expand  $\hat{\sigma}_{t_0}$  and  $\hat{\alpha}_{t_0}^{-1}$ :

$$\widehat{\sigma}_{t_0} = \left(\frac{1 - \exp(-t_0)}{2 - \exp(-t_0)}\right)^{\frac{1}{2}} = \left(1 - \frac{1}{2 - \exp(-t_0)}\right)^{\frac{1}{2}} = \mathcal{O}\left(N^{-C_{\sigma}}\right).$$

and

$$\widehat{\alpha}_{t_0}^{-1} = \left(\frac{2 - \exp(-t_0)}{\exp\left(-\frac{t_0}{2}\right)}\right) = 2\exp\left(\frac{t_0}{2}\right) - \exp\left(-\frac{t_0}{2}\right) = \mathcal{O}\left(N^{-C_{\sigma}}\right)$$

So we take  $K_4 = \mathcal{O}(N^{C_{\sigma}})$ .

\* Step B.2: Approximation for  $-C_2 x/(\alpha_t^2 + C_2 \sigma_t^2)$ . We use  $\alpha_t^2 + \sigma_t^2 = 1$  to rewrite  $(\alpha_t^2 + C_2 \sigma_t^2)^{-1}$  as  $(C_2 + (1 - C_2)\alpha_t^2)^{-1}$ . We first utilize  $\mathcal{T}_{\alpha^2}$  in Lemma D.8 for the approximation of  $\alpha_t^2$ . This gives error  $\epsilon_{\alpha^2}$ . Next, we utilize  $\mathcal{T}_{\alpha \alpha 1}$  in Lemma C.8 for the approximation of the inverse of  $\alpha_t^2$ .

Next, we utilize  $\mathcal{T}_{\text{rec},1}$  in Lemma C.8 for the approximation of the inverse of  $\alpha_t^2$ . This gives error

$$\mathcal{T}_{\text{rec},1} - \frac{1}{\alpha_t^2} \left| \le \epsilon_{\text{rec},1} + \frac{\left| \mathcal{T}_{\alpha_t^2} - \alpha_t^2 \right|}{\epsilon_{\text{rec},1}^2} \le \epsilon_{\text{rec},1} + \frac{\epsilon_{\alpha^2}}{\epsilon_{\text{rec},1}^2} \right|$$

Next, we utilize  $\mathcal{T}_{\text{mult},2}$  for the approximation of the product of  $(C_2 + (1 - C_2)\alpha_t^2)^{-1}$ and x. This gives error

$$\left|\mathcal{T}_{\text{mult},2} - \left(\frac{x}{C_2 + (1 - C_2)\alpha_t^2}\right)\right| \le \epsilon_{\text{mult},2} + 2K_3\left(\epsilon_{\text{rec},1} + \frac{\epsilon_{\alpha^2}}{\epsilon_{\text{rec},1}^2}\right),$$

and from Lemma C.8,  $K_3$  is positive constant such that  $x \in [-K_3, K_3]$ and  $\alpha_t^{-1} \in [-K_3, K_3]$ . Since  $x \in [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]$  and  $\alpha_T^{-1} = (\exp(-C_\alpha \log N/2))^{-1} = N^{C_\alpha/2}$ , we take  $K_3 = N^{C_\alpha/2}$ .

## \* Step B.3: Error Bound on Every Approximation Combined.

Combining **Step B.1** and **Step B.2**, we obtain the total network with error bounded by

$$\epsilon_{\text{score}} \leq \epsilon_{\text{mult},2} + 2K_3 \left( \epsilon_{\text{rec},1} + \frac{\epsilon_{\alpha^2}}{\epsilon_{\text{rec},1}^2} \right) + \epsilon_{\text{mult},1} + 4K_4^3 \epsilon_2.$$

Next, we specify on the choice of  $\epsilon$  in each approximation to attain a final approximation error of order  $N^{-\beta}$ .

· For the Error of the First Inverse Operator:

$$\epsilon_{\mathrm{rec},1} = \mathcal{O}\left(N^{-(\beta + \frac{1}{2}C_{\alpha})}\right)$$

· For the Error of the Second and Third Inverse Operator:

$$\epsilon_{\mathrm{rec},2}, \epsilon_{\mathrm{rec},3} = \mathcal{O}\left(N^{-(\beta+3C_{\sigma})}\right).$$

• For the Error of  $f_1$ :

$$\epsilon_{f_1} = \mathcal{O}\left(N^{-(3\beta + 9C_{\sigma})}\right)$$

• For the Error of  $f_2$ :

$$\epsilon_{f_2} = \mathcal{O}\left(N^{-(\beta+3C_{\sigma})}\right)$$

· For the Error of  $\hat{\sigma}$ :

$$\epsilon_{\widehat{\sigma}} = \mathcal{O}\left(N^{-(3\beta + 9C_{\sigma})}\right)$$

• For the Error of  $\widehat{\alpha}$ :

$$\epsilon_{\widehat{\alpha}} = \mathcal{O}\left(N^{-(\beta+3C_{\sigma})}\right).$$

• For the Error of  $\alpha^2$ :

$$\epsilon_{\alpha^2} = \mathcal{O}\left(N^{-(3\beta + \frac{3}{2}C_{\alpha})}\right).$$

• For the Error of the Two Product Operators:

$$\epsilon_{\text{mult},1}, \epsilon_{\text{mult},2} = \mathcal{O}(N^{-\beta}).$$

With above error choice, we have

$$\mathcal{T}_{\text{score}}(x, y, t) - f_3(x, y, t) | \le N^{-\beta}.$$
(D.3)

1674 Combining (D.2), (D.3) and dropping lower order term, we obtain 1675 1676  $\|\mathcal{T}_{\text{score}} - \nabla \log p_t(x|y)\|_{\infty} \lesssim \frac{B}{\sigma_t} N^{-\beta} (\log N)^{\frac{k_1+1}{2}}.$ 1677 1678 We have completed the first part of the proof. Next, we select the parameter bounds 1679 based on all the above approximations. **Step C: Transformer Parameter Bound.** 1681 Our result highlights the influence of N under varying  $d_x$ . Therefore, for the transformer parameter bounds, we keep terms with  $d_x$ , d, L appearing in the exponent of N and  $\log N$ . 1683 – Parameter Bound on  $W_Q$  and  $W_K$ . Given error  $\epsilon$ , the bound on each operation follows: \* For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1687  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_{\sigma})\frac{2dL+4d+1}{d}}\right).$ 1689 \* For  $\epsilon_{f_2}$ : By Lemma C.6, we have  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(\beta+3C_{\sigma})\frac{2dL+4d+1}{d}}\right).$ 1693 \* For  $\epsilon_{mult,1}$ : By Lemma C.8 with m = 4, we have 1695 1696  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_K||_{2,\infty} = \mathcal{O}(N^{9\beta}).$ 1698 \* For  $\epsilon_{\text{mult},2}$ : By Lemma C.8 with m = 2, we have 1700  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_K||_{2,\infty} = \mathcal{O}(N^{5\beta}).$ \* For  $\epsilon_{rec,1}$ : By Lemma C.9, we have 1702  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{3\beta + \frac{3C_{\alpha}}{2}}\right)$ 1704 1706 \* For  $\epsilon_{rec,2}$  and  $\epsilon_{rec,3}$ : By Lemma C.9, we have 1707  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_K||_{2,\infty} = \mathcal{O}(N^{3\beta+9C_{\sigma}}).$ 1708 1709 \* For  $\epsilon_{\hat{\alpha}}$ : By Lemma C.11, we have 1710 1711  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_Q||_{2,\infty} = \mathcal{O}(N^{3\beta+9C_{\sigma}}).$ 1712 1713 \* For  $\epsilon_{\alpha^2}$ : By Lemma C.11, we have 1714 1715  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2,\infty}, ||W_Q||_{2,\infty} = \mathcal{O}\left(N^{9\beta + \frac{9C_\alpha}{2}}\right).$ 1716 1717 \* For  $\epsilon_{\widehat{\sigma}}$ : By Lemma C.11, we have 1718 1719  $||W_Q||_2, ||W_K||_2, ||W_Q||_{2\infty}, ||W_Q||_{2\infty} = \mathcal{O}\left(N^{9\beta+27C_{\sigma}}\right).$ We select the largest parameter bound from  $\epsilon_{f_1}$  that remains valid across all other 1722 approximations. 1723 - Parameter Bound on  $W_O$  and  $W_V$ . 1724 Given error  $\epsilon$ , the bound on each operation follows: 1725 \* For  $\epsilon_{f_1}$ : By Lemma C.5, we have 1726 1727

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(3\beta+9C_{\sigma})}{d}}\right).$$

1728 1729	* For $\epsilon_{f_2}$ : By Lemma C.6, we have
1730	$(\beta + 3C_{\sigma}))$
1731	$\left\ W_{O}\right\ _{2}, \left\ W_{O}\right\ _{2,\infty} = \mathcal{O}\left(N^{-\frac{d}{d}}\right).$
1732	* For $c_{1}$ , $\cdot$ By Lemma C 8 with $m = 4$ we have
1733	* For $\epsilon_{mult,1}$ . By Lemma C.8 with $m = 4$ , we have
1734	$\ W_O\ _2, \ W_O\ _{2,\infty} = \mathcal{O}\left(N^{-4\beta}\right).$
1735	
1730	* For $\epsilon_{mult,2}$ : By Lemma C.8 with $m = 2$ , we have
1738	$\ \mathbf{U}\  \ \mathbf{U}\  = \mathcal{O}(N^{-2\beta})$
1739	$\ WO\ _2, \ WO\ _{2,\infty} = O(W^{-1}).$
1740	* For $\epsilon_{rec,1}$ : By Lemma C.9, we have
1741	
1742	$\left\ W_O\right\ _2, \left\ W_O\right\ _{2,\infty} = \mathcal{O}\left(N^{-(\beta + \frac{C_\alpha}{2})}\right).$
1743	
1744	* For $\epsilon_{rec,2}$ and $\epsilon_{rec,3}$ : By Lemma C.9, we have
1745	$(\beta_1, \beta_2, \beta_3)$
1740	$\ W_O\ _2, \ W_O\ _{2,\infty} = \mathcal{O}\left(N^{-(\beta+\beta C_\sigma)}\right).$
1748	t For c. t Py Lomma C 11, we have
1749	* For $\epsilon_{\hat{\alpha}}$ : By Lemma C.11, we have
1750	$  W_{\alpha}  ,   W_{\alpha}  , = \mathcal{O}\left(N^{-(\beta+3C_{\sigma})}\right)$
1751	$\  (O_{12}, \  (O_{12}, \infty) - O_{12}) \ _{2,\infty} = O_{12} (1, 1, 1)$
1752	* For $\epsilon_{\alpha^2}$ : By Lemma C.11, we have
1753	
1755	$\left\ W_O\right\ _2, \left\ W_O\right\ _{2,\infty} = \mathcal{O}\left(N^{-(3\beta + rac{3C_{\alpha}}{2})} ight).$
1756	
1757	* For $\epsilon_{\hat{\sigma}}$ : By Lemma C.11, we have
1758	$  \mathbf{u}_{\mathbf{r}}   =   \mathbf{u}_{\mathbf{r}}   = \langle \mathbf{o} \left( \mathbf{v}_{\mathbf{r}} - (3\beta + 9C_{-}) \right)$
1759	$\ W_O\ _2, \ W_O\ _{2,\infty} = O\left(N^{-(0,p+0,0,p)}\right).$
1760	Since we do not impose any relation on $C$ , $C$ , and $\beta$ , we simply take looser bound
1761	$\ W_{\Omega}\ _{2}, \ W_{\Omega}\ _{2} = N^{-\beta}$ . Moreover, since only $\epsilon_{f_{1}}$ and $\epsilon_{f_{2}}$ involve the reshape
1762	operation From Lemma B 2, we take $\mathcal{O}(\sqrt{d})$ and $\mathcal{O}(d)   W_V  _2$ and $  W_V  _2$
1764	D = (A - B) = A C - B C
1765	- Parameter Bound for $W_1$ .
1766	Given error $\epsilon$ , the bound on each operation follows:
1767	* For $\epsilon_{f_1}$ : By Lemma C.5, we have
1768	$(1, \dots, 1) = (1, \dots, 1) = (2 (3\beta + 9C_{\sigma}))$
1769	$\ W_1\ _2, \ W_1\ _{2,\infty} = \mathcal{O}\left(N^{-d} \cdot \log N\right).$
1770	<b>For</b> c. <b>P</b> <sub>U</sub> Lemma C. 6, we have
1771	* For $e_{f_2}$ . By Lemma C.0, we have
1772	$\ W_1\  = \ W_1\  = -\mathcal{O}\left(N^{\frac{(\beta+3C_{\sigma})}{d}} \cdot \log N\right)$
1774	$\  (r_1 \ _2) \  (r_1 \ _{2,\infty}) = O((r_1 r_1 r_2))$
1775	* For $\epsilon_{\text{mult,1}}$ : By Lemma C.8 with $m = 4$ and $C = K_4$ in (C.9), we have
1776	
1777	$\left\ W_{1}\right\ _{2}, \left\ W_{1}\right\ _{2,\infty} = \mathcal{O}\left(K_{4} \cdot N^{4\beta}\right) = \mathcal{O}\left(N^{\left(4\beta + C_{\sigma}\right)}\right).$
1778	
1779	* For $\epsilon_{\text{mult},2}$ : By Lemma C.8 with $m = 2$ and $C = K_3$ in (C.10), we have
1/80	$\  \mathbf{T} \mathbf{T} \  = \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \right) \right) \right) \right) \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \  \mathbf{T} \  + \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  \mathbf{T} \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \right) \  \mathbf{T} \  \mathbf{T} \  = \left( \frac{1}{2} \left( \frac{1}{2} \right) \  \mathbf{T} \  $
101	$  W_1  _2,   W_1  _{2,\infty} = \mathcal{O}\left(K_3 \cdot N^{2\rho}\right) = \mathcal{O}\left(N^{(2\rho + \frac{1}{2})}\right).$

1782	* For $c \rightarrow By I$ amma $C = 0$ we have
1783	* For erec,1. By Lemma C.9, we have
1784	$\ W_1\ _{2^*} \ W_1\ _{2^* - \pi} = \mathcal{O}\left(N^{2\beta + C_{\alpha}}\right).$
1785	$11  112711  112,\infty  ( )$
1786	* For $\epsilon_{rec,2}$ and $\epsilon_{rec,3}$ : By Lemma C.9, we have
1787	
1788	$\ W_1\ _2, \ W_1\ _{2,\infty} = \mathcal{O}\left(N^{(2\beta+6C_{\sigma})}\right).$
1700	
1791	* For $\epsilon_{\hat{\alpha}}$ : By Lemma C.11, we have
1792	$\ \mathbf{U}_{\mathcal{I}}\  = \ \mathbf{U}_{\mathcal{I}}\  \qquad $
1793	$\ W_1\ _2, \ W_1\ _{2,\infty} = O\left(N^{-1} + \log N\right).$
1794	* For c a. By Lemma C 11 we have
1795	* For $c_{\alpha}^2$ . By Echinic C.11, we have
1796	$\ W_1\ _{\infty}, \ W_1\ _{\infty} = \mathcal{O}\left(N^{(3\beta + \frac{3C_{\alpha}}{2})} \cdot \log N\right).$
1797	$  \cdot  _{1} _{2},   \cdot  _{1} _{2,\infty} = (-1)^{-1}$
1798	* For $\epsilon_{\hat{\sigma}}$ : By Lemma C.11, we have
1/99	
1801	$\left\ W_{1}\right\ _{2}, \left\ W_{1}\right\ _{2,\infty} = \mathcal{O}\left(N^{(3eta+9C_{\sigma})} \cdot \log N ight).$
1802	
1803	We select the largest parameter bound from $\epsilon_{f_1}$ that remains valid across all other
1804	approximations.
1805	
1806	- Parameter Bound for $W_2$ .
1807	Given error $\epsilon$ , the bound on each operation follows:
1808	
1809	* For $\epsilon_{f_1}$ : By Lemma C.5, we have
1811	$\begin{pmatrix} & (28\pm0C_{-}) \end{pmatrix}$
1812	$\left\ W_2\right\ _2, \left\ W_2\right\ _{2,\infty} = \mathcal{O}\left(N^{\frac{(\beta)+\beta \in \mathcal{O}(\beta)}{d}}\right).$
1813	
1814	* For $\epsilon_{f_2}$ : By Lemma C.6, we have
1815	$\left    \mathbf{U}_{\mathbf{U}}   \right  = \left    \mathbf{U}_{\mathbf{U}}   \right  = \left  \mathbf{O} \left( \mathbf{V}^{\left( \beta + 3C_{\sigma} \right)} \right) \right $
1816	$\ W_2\ _2, \ W_2\ _{2,\infty} = \mathcal{O}\left(N  d \right).$
1817	* For $c_{1}$ , By Lemma C 8 with $m = 4$ , we have
1818	* For $\epsilon_{\text{mult},1}$ . By Lemma C.8 with $m = 4$ , we have
1819	$\ W_2\ _2, \ W_2\ _2 = \mathcal{O}(N^{4\beta}).$
1821	
1822	* For $\epsilon_{mult,2}$ : By Lemma C.8 with $m = 2$ , we have
1823	$\ W_{L}\  = \ W_{L}\  = O(N^{2\beta})$
1824	$\ W_2\ _2, \ W_2\ _{2,\infty} = O(1^{V+1}).$
1825	* For $\epsilon_{rec.1}$ : By Lemma C.9, we have
1826	
1827	$\left\ W_2\right\ _2, \left\ W_2\right\ _{2,\infty} = \mathcal{O}\left(N^{\left(eta+rac{C_{lpha}}{2} ight)} ight).$
1828	
1830	* For $\epsilon_{rec,2}$ and $\epsilon_{rec,3}$ : By Lemma C.9, we have
1831	$(\alpha + 2C)$
1832	$\ W_2\ _2, \ W_2\ _{2,\infty} = \mathcal{O}\left(N^{(\beta+3C_{\sigma})}\right).$
1833	Ear Da Lamma C 11 and have
1834	* For $\epsilon_{\hat{\alpha}}$ : By Lemma C.11, we have
1835	$  W_{\sigma}   =   W_{\sigma}   = - \mathcal{O}\left(N^{(\beta+3C_{\sigma})}\right)$
	$\ vv_2\ _{2}, \ vv_2\ _{2,\infty} = O\left(1v^{-1}\right)$

\* For  $\epsilon_{\alpha^2}$ : By Lemma C.11, we have  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{(3\beta + \frac{3C_{\alpha}}{2})}\right).$ \* For  $\epsilon_{\hat{\sigma}}$ : By Lemma C.11, we have  $||W_2||_2, ||W_2||_{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_{\sigma})}\right).$ We select the largest parameter bound from  $\epsilon_{f_1}$  that remains valid across all other approximations. – Parameter Bound for E. Since only  $\epsilon_{f_1}$  and  $\epsilon_{f_2}$  involve the reshape operation. From Lemma B.2, we take  $\mathcal{O}(d^{1/2}L^{3/2}).$ By integrating results above, we derive the following parameter bounds for the transformer network, ensuring valid approximation across all ten approximations.  $\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_{\sigma})\frac{2dL+4d+1}{d}}\right);$  $||W_V||_2 = \mathcal{O}(\sqrt{d}); ||W_V||_{2,\infty} = \mathcal{O}(d); ||W_O||_2, ||W_O||_{2,\infty} = \mathcal{O}(N^{-\beta});$  $\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{4\beta + 9C_{\sigma} + \frac{3C_{\alpha}}{2}} \cdot \log N\right); \|E^{\top}\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);$  $\left\|W_{2}\right\|_{2}, \left\|W_{2}\right\|_{2,\infty} = \mathcal{O}\left(N^{4\beta + 9C_{\sigma} + \frac{3C_{\alpha}}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_{t}\right).$ This completes the proof. 

#### D.2 MAIN PROOF OF THEOREM 3.1 UNDER STRONGER ASSUMPTION

We state the formal version of Theorem 3.1 under stronger assumption.

Next, similar to the proof of Theorem 3.1, we need the truncation of x due to the unboundedness as well. 

Lemma D.12 (Truncate x, Lemma B.2 of (Fu et al., 2024b).). Under stronger Assumption 3.1. For any  $R_3 > 1$ , we have:

$$\int_{\|x\|_{\infty} \ge R_3} p_t(x|y) \mathrm{d}x \lesssim R_3 \exp\left(-C_2' R_2^2\right).$$

$$\int_{\|x\|_{\infty} \ge R_3} \|\nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \lesssim R_3 \exp\left(-C_2' R_3^2\right) \lesssim \frac{1}{\sigma_t^2} R_3^3 \exp\left(-C_2' R_3^2\right),$$

where  $C'_2 = C_2/(2 \max(1, C_2))$ .

Again, unlike result under generic Assumption 3.1, the explicit form of  $p_t(x|y)$  and the upper and the lower bound of the joint distribution Lemma D.2 automatically allow us to skip the threshold  $\epsilon_{low}$  as in Lemma C.15.

Theorem D.1 (Approximation Score Function with Transformer under Stronger Hölder Assumption (Formal Version of Theorem 3.1)). Under stronger Assumption 3.1 and  $d_x = \Omega(\frac{\log N}{\log \log N})$ . For any precision parameter  $0 < \epsilon < 1$  and smoothness parameter  $\beta > 0$ , let  $\epsilon \leq \mathcal{O}(N^{-\beta})$  for some  $N \in \mathbb{N}$ . For some positive constants  $C_{\alpha}, C_{\sigma} > 0$ , for any  $y \in [0, 1]^{d_y}$  and  $t \in [N^{-C_{\sigma}}, C_{\alpha} \log N]$ , there exists a  $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_{R}^{h,s,r}$  such that the conditional score approximation satisfies 

$$\int_{\mathbb{R}^{d_x}} \|\mathcal{T}_{\text{score}}(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx = \mathcal{O}\left(\frac{B^2}{\sigma_t^2} \cdot N^{-\frac{2\beta}{d_x + d_y}} \cdot (\log N)^{\beta + 1}\right).$$

Notably, for  $\epsilon = \mathcal{O}(N^{-\beta})$ , the approximation error has the upper bound  $\widetilde{\mathcal{O}}(\epsilon^{2/(d_x+d_y)}/\sigma_{\epsilon}^2)$ . The parameter bounds in the transformer network class satisfy 

$$\begin{split} \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} &= \mathcal{O}\left(N^{\frac{3\beta(2d_x+4d+1)}{d(d_x+d_y)} + \frac{9C_{\alpha}(2d_x+4d+1)}{d}}\right);\\ \|W_V\|_2 &= \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{\beta}{d_x+d_y}}\right);\\ \|W_1\|_2, \|W_1\|_{2,\infty} &= \mathcal{O}\left(N^{\frac{4\beta}{d_x+d_y} + 9C_{\sigma} + \frac{3C_{\alpha}}{2}} \cdot \log N\right); \|E^{\top}\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);\\ \|W_2\|_2, \|W_2\|_{2,\infty} &= \mathcal{O}\left(N^{\frac{4\beta}{d_x+d_y} + 9C_{\sigma} + \frac{3C_{\alpha}}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t\right). \end{split}$$

Proof of Theorem 3.1 under Stronger Assumption. For simplicity, we change the variable N to  $N^{\overline{d_x+d_y}}$  in the following subsection. We put the original form back at the end of the proof.

We take 
$$C_x = \sqrt{\frac{2\beta}{C_2'}}$$
 in Lemma D.11 and  $R_3 = C_x \sqrt{\log N}$  in Lemma D.12.

With the transformer parameter bounds in Lemma D.11, we have  $\|\mathcal{T}_{score}\|_2 \leq \sqrt{\log N}/\sigma_t$  for any  $x \in \mathbb{R}^{d_x}, y \in \mathbb{R}^{d_y}$  and t > 0. We start with the truncation on x 

1955 The transformer parameter norm bounds follow Lemma D.11, with the replacement of N with 1956  $N^{1/d_x+d_y}$ . This gives in  $t \in [N^{-C_\alpha/(d_x+d_y)}, C_\sigma(\log N)^{1/(d_x+d_y)}]$ . For a better interpretation of the 1957 cutoff and early stopping time parameter, we reset  $C_\alpha = (d_x + d_y)C_\alpha$  and  $C_\sigma = (d_x + d_y)C_\sigma$  such 1958 that  $t \in [N^{-C_\alpha}, C_\sigma \log N]$ .

This completes the proof.

## <sup>1998</sup> E PROOF OF THEOREM 3.2

## 2000 Overview of Our Proof Strategy of Theorem 3.2.

- Step 0. Preliminaries. We introduce the mixed risk that accounts for risk with the distribution of the mask signal in Definition E.1. We restate the loss function and the score matching technique in Definition E.2.
- Step 1. Truncate the Domain of the Risk. We truncate the domain of the loss function in order
   to obtain finite covering number of transformer network class. Precise definition of the
   truncated loss function class is in Definition E.4. We bound the error from the truncation
   from the assumed light tail condition in Lemma E.1.
- Step 2. Derive the Covering Number of Transformer Network. We introduce the covering number of a given function class in Definition E.5. We provide lemma detailing the calculation of the covering number for transformer architecture in Lemma E.2. We derive the covering numbers under the respective parameter configurations for our two previous main results in Lemma E.3.
- Step 3. Bound the True Risk on Truncated Domain. With the previous steps, we present the upper-bound of the mixed risk in Lemma E.4.

Organization. Appendix E.1 includes auxiliary lemmas for supporting our proof of Theorem 3.2.
 Appendix E.2 includes the main proof of Theorem 3.2.

E.1 AUXILIARY LEMMAS FOR THEOREM 3.2

**Step 0: Preliminary Framework.** We evaluate the quality of the estimator  $s_W$  through the risk:

$$\mathcal{R}(s_W) \coloneqq \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_t, y} \| s_W(x_t, y, t) - \nabla \log p_t(x_t | y) \|_2^2 \mathrm{d}t.$$
(E.1)

**Definition E.1** (Mixed Risk). The risk (E.1) considers guidance y throughout whole the diffusion process. We refer to it as the conditional score risk. In contrast, we have the mixed risk  $\mathcal{R}_m$  that accounts for the distribution of the mask signal  $\tau = \{\emptyset, id\}$  with  $P(\tau = \emptyset) = P(\tau = id) = 0.5$ :

$$\mathcal{R}_m(s_W) \coloneqq \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{(x_t, y, \tau)} \left[ \|s_W(x_t, \tau y, t) - \nabla \log p_t(x_t | \tau y)\|_2^2 \right] \mathrm{d}t, \tag{E.2}$$

**Remark E.1.** Given the score estimator  $\hat{s}$  trained from the empirical loss, the conditional score risk is upper-bounded by twice of the mixed risk. That is, we have  $\mathcal{R}(\hat{s}) \leq 2\mathcal{R}_m(\hat{s})$ . This follows from direct calculation:

$$\mathcal{R}_m(\widehat{s}) = \frac{1}{2} \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_t} \left[ \|\widehat{s}(x_t, \emptyset, t) - \nabla \log p_t(x_t)\|_2^2 \right] \mathrm{d}t + \frac{1}{2} \mathcal{R}(\widehat{s}).$$

**Definition E.2** (Loss Function and Score Matching). Let  $x = x_t | x_0$  denote the random variable following Gaussian distribution  $N(\alpha_t x_0, \sigma_t^2 I_{d_x})$ , we define loss function and score matching loss:

$$\ell(x, y; s_W) \coloneqq \int_{T_0}^T \frac{1}{T - T_0} \mathbb{E}_{\tau, x} \left[ \| s_W(x_t, \tau y, t) - \nabla \log p_t (x_t | x_0) \|_2^2 \right] \mathrm{d}t,$$

$$\mathcal{L}(s_W) \coloneqq \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_0, y} \left[ \mathbb{E}_{\tau, x} \left[ \|s_W(x_t, \tau y, t) - \nabla \log p_t(x_t | x_0) \|_2^2 \right] \right] \mathrm{d}t.$$

**Remark E.2.** Given i.i.d samples  $\{x_{0,i}, y_i\}_{i=1}^n$ , we write  $\ell(x_i, y_i; s_W)$  with the understanding that  $x_i = x_t | x_{0,i}$ . When context is clear, we use  $\ell(x_i, y_i; s_W)$  and  $\ell(x_{0,i}, y_i; s_W)$ ;  $\{x_{0,i}, y_i\}_{i=1}^n$  and  $\{x_i, y_i\}_{i=1}^n$  interchangeably.

**Remark E.3.** By (Vincent, 2011),  $\mathcal{L}(s_W)$  and  $\mathcal{R}_m(s_W)$  differ by a constant that is inconsequential to the minimization. Therefore, minimizing the mixed risk is equivalent to minimizing the score matching loss

2055 2056 2057

2079

2084

2085

2086

2089 2090

2091

2092

2100 2101 2102 **Definition E.3** (Empirical Risk). Consider a score estimator  $s_W \in \mathcal{T}_R^{h,s,r}$ . Recall the definition of empirical loss:  $\widehat{\mathcal{L}}(s_W) = \sum_{i=1}^n \frac{1}{n} \ell(x_i, y_i; s_W)$ . Let  $s^\circ := \nabla \log p_t(x|y)$ , we define empirical risk:

$$\widehat{\mathcal{R}}_m(s_W) \coloneqq \widehat{\mathcal{L}}(s_W) - \widehat{\mathcal{L}}(s^\circ) = \sum_{i=1}^n \frac{1}{n} \ell(x_i, y_i; s_W) - \sum_{i=1}^n \frac{1}{n} \ell(x_i, y_i; s^\circ).$$

**Remark E.4.** The key distinction between  $\mathcal{R}_m$  and  $\mathcal{L}$  lies in their formulations. Specifically,  $\mathcal{R}_m$ measures the expected difference between  $s_W$  and the ground truth  $\nabla \log p_t(x|y)$  with respect to  $(x_t, y, \tau)$ . In contrast, the score matching loss  $\mathcal{L}$  provides an explicit calculation based on the sample  $\{x_{0,i}, y_i\}_{i=1}^n$ . With the tower property of conditional expectation,  $\mathcal{L}$  measures the expected difference between  $s_W$  and  $\nabla \log p_t(x|x_0)$  first with respect to  $(x_t|x_0, \tau)$ , and then with respect to  $x_0$ .

**Remark E.5.** Observe (I):  $s^{\circ} = \nabla \log p_t(x|y)$  is the ground truth of score function with  $\mathcal{R}_m(s^{\circ}) = 0$ , and (II): By (Vincent, 2011),  $\mathcal{R}_m$  and  $\mathcal{L}$  differ by a constant. Based on (I) and (II), we define the empirical risk  $\widehat{\mathcal{R}}_m$  using the score matching loss as an intermediary:  $\mathcal{R}_m(s_W) = \mathcal{R}_m(s_W) - \mathcal{R}_m(s^{\circ}) = \mathcal{L}(s_W) - \mathcal{L}(s^{\circ})$ . This leads to the definition of the empirical risk  $\widehat{\mathcal{R}}_m$  as a practical approximation of the true risk difference  $\mathcal{R}_m(s_W) - \mathcal{R}_m(s^{\circ})$ .

**Remark E.6.** For any score estimator  $s_W \in \mathcal{T}_R^{h,s,r}$  obtained from the training with i.i.d. samples  $\{x_i, y_i\}_{i=1}^n$ , it holds  $\mathbb{E}_{\{x_i, y_i\}_{i=1}^n}[\widehat{\mathcal{R}}_m(s_W)] = \mathcal{R}_m(s_W)$ . This follows from direct calculation with Definition E.3 and the i.i.d. assumption.

Step 1: Domain Truncation of the Risk. We define the loss function with truncated domain. This is essential for obtaining finite covering number for transformer network class.

**Definition E.4** (Truncated Loss). We define the truncated domain of the score function by  $\mathcal{D} := [-R_{\mathcal{T}}, R_{\mathcal{T}}]^{d_x} \times [0, 1]^{d_y} \cup \emptyset$ . Given loss function  $\ell(x, y; s_W)$ , we define the truncated loss:

$${}^{\mathrm{runc}}(x, y; s_W) := \ell(x, y; s_W) \mathbb{1}\{ \|x\|_{\infty} \le R_{\mathcal{T}} \}.$$
(E.3)

Similarly, we define  $\mathcal{L}^{\text{trunc}}(s_W) \coloneqq \mathcal{L}(s_W) \mathbb{1}\{\|x\|_{\infty} \leq R_{\mathcal{T}}\}$ ,  $\mathcal{R}_m^{\text{trunc}}(s_W) \coloneqq \mathcal{R}_m(s_W) \mathbb{1}\{\|x\|_{\infty} \leq R_{\mathcal{T}}\}$  and  $\widehat{\mathcal{R}}_m^{\text{trunc}}(s_W) \coloneqq \widehat{\mathcal{R}}_m(s_W) \mathbb{1}\{\|x\|_{\infty} \leq R_{\mathcal{T}}\}$ . We define the function class of the truncated loss by

 $\mathcal{S}(R_{\mathcal{T}}) \coloneqq \{\ell(\cdot, \cdot; s_W) : \mathcal{D} \to \mathbb{R} \mid s_W \in \mathcal{T}_R^{h, s, r}\}.$ (E.4)

Next, we introduce the following lemma dealing with the error bound for the truncation of the loss.

**Lemma E.1** (Truncation Error, Lemma D.1 of (Fu et al., 2024b)). Consider the truncated loss  $\ell^{\text{trunc}}(x, y; s_W)$  and  $t \in [n^{-\mathcal{O}(1)}, \mathcal{O}(\log n)]$ . Under generic Assumption 3.1, we have  $|\ell(x, y; s_W)| \leq 1/t_0$ . Consider the parameter configuration in Theorem 3.1, it holds:

$$\mathbb{E}_{x,y}\left[\left|\ell(x,y,t)-\ell^{\mathrm{trunc}}(x,y,s)\right|\right] \lesssim \exp\left(-C_2 R_{\mathcal{T}}^2\right) R_{\mathcal{T}}\left(\frac{1}{t_0}\right).$$

Moreover, under stronger Assumption 3.1, we have  $|\ell(x, y; s_W)| \leq \log(1/t_0)$ . Consider the parameter configuration in Theorem D.1, it holds:

$$\mathbb{E}_{x,y}\left[\left|\ell(x,y,t) - \ell^{\text{trunc}}(x,y,s)\right|\right] \lesssim \exp\left(-C_2 R_{\mathcal{T}}^2\right) R_{\mathcal{T}} \log\left(\frac{1}{t_0}\right)$$

2103 Step 2: Covering Number of Transformer Network Class. We begin with the definition.

**Definition E.5** (Covering Number). Given a function class  $\mathcal{F}$  and a data distribution P. Sample n data points  $\{X_i\}_{i=1}^n$  from P, then the covering number  $\mathcal{N}(\epsilon, \mathcal{F}, \{X_i\}_{i=1}^n, \|\cdot\|)$  is the smallest size of

a collection (a cover)  $C \in F$  such that for any  $f \in F$ , there exist  $\hat{f} \in C$  satisfying

$$\max_{i} \left\| f(X_{i}) - \widehat{f}(X_{i}) \right\| \le \epsilon.$$

2111 Further, we define the covering number with respect to the data distribution as

$$\mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|) = \sup_{\{X_i\}_{i=1}^n \sim P} \mathcal{N}(\epsilon, \mathcal{F}, \{X_i\}_{i=1}^n, \|\cdot\|).$$

Next, we introduce the following lemma that provides results for the calculation of the covering number for transformer networks.

Lemma E.2 (Modified from Theorem A.17 of Edelman et al. (2022)).

Let 
$$\mathcal{T}_{R}^{h,s,r}(C_{\mathcal{T}}, C_{Q}^{2,\infty}, C_{Q}, C_{K}^{2,\infty}, C_{K}, C_{V}^{2,\infty}, C_{V}, C_{O}^{2,\infty}, C_{O}, C_{E}, C_{f_{1}}^{2,\infty}, C_{f_{1}}, C_{f_{2}}^{2,\infty}, C_{f_{2}}, L_{\mathcal{T}})$$

represent the class of functions of one transformer block satisfying the norm bound for matrix and Lipsichitz property for feed-forward layers. Then for all data point  $||X||_{2,\infty} \leq R_T$  we have

$$\log \mathcal{N}(\epsilon_{c}, \mathcal{T}_{R}^{n, s, i}, \|\cdot\|_{2}) \leq \frac{\log(nL\tau)}{\epsilon_{c}^{2}} \cdot \left(\alpha^{\frac{2}{3}} \left(d^{\frac{2}{3}} \left(C_{F}^{2, \infty}\right)^{\frac{4}{3}} + d^{\frac{2}{3}} \left(2(C_{F})^{2}C_{OV}C_{KQ}^{2, \infty}\right)^{\frac{2}{3}} + 2\left((C_{F})^{2}C_{OV}^{2, \infty}\right)^{\frac{2}{3}}\right)\right)^{3},$$

where  $\alpha \coloneqq (C_F)^2 C_{OV} (1 + 4C_{KQ}) (R_T + C_E).$ 

With Lemma E.2, we derive the covering number under transformer weights configuration in Theorem 3.1 and Theorem D.1.

**2133** Lemma E.3 (Covering Number for  $S(R_T)$ ). Given  $\epsilon_c > 0$  and consider  $||x||_{\infty} \le R_T$ . With 2134 sample  $\{x_i, y_i\}_{i=1}^n$ , the  $\epsilon_c$ -covering number for  $S(R_T)$  with respect to  $||\cdot||_{L_{\infty}}$  under the network 2135 configuration in Theorem 3.1 satisfies

$$\operatorname{og} \mathcal{N}\left(\epsilon_{c}, \mathcal{S}(R_{\mathcal{T}}), \left\|\cdot\right\|_{\infty}\right) \lesssim \frac{\log n}{\epsilon_{c}^{2}} N^{\nu_{1}} (\log N)^{\nu_{2}} (R_{\mathcal{T}})^{2}$$

where  $\nu_1 = \frac{172\beta}{(d_x + d_y)} + 104C_{\sigma}$  and  $\nu_2 = 12d_x + 12\beta + 2$ . Moreover, under network configuration in Theorem D.1, we have

$$\log \mathcal{N}\left(\epsilon_{c}, S(R_{\mathcal{T}}), \left\|\cdot\right\|_{\infty}\right) \lesssim \frac{\log n}{\epsilon_{c}^{2}} N^{\nu_{3}} (\log N)^{10} (R_{\mathcal{T}})^{2},$$

where  $\nu_3 = 48d\beta(L+2)(d_x+2d+1)/(d_x+d_y) + 144dC_{\sigma}(L+2) - 8\beta$ .

## Step 3: Bound the True Risk on Truncated Domain. We begin with the definition.

**Definition E.6.** Let  $s^{\circ} \coloneqq \nabla \log p_t(x|y)$  denote the ground truth of score function for simplicity. Given i.i.d samples  $\{x_i, y_i\}_{i=1}^n$  and a score estimator  $s_W \in \mathcal{T}_R^{h,s,r}$ , we define the difference function:

$$\Delta_n(s_W, s^\circ) \coloneqq \left| \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \widehat{\mathcal{R}}_m^{\mathrm{trunc}}(s_W) - \mathcal{R}_m^{\mathrm{trunc}}(s_W) \right] \right|.$$

**Remark E.7.** Note that the difference function  $\Delta_n(s_W, s^\circ)$  measures the expected difference between the truncated empirical risk and the truncated mixed risk with respect to the training sample. Since the true risk is unattainable, we construct  $\Delta_n(s_W, s^\circ)$  serving as an intermediate that allows us to derive the upper-bound on the mixed risk. Surprisingly, we are able to handle the upper-bound of the difference function, presented in Lemma E.4. **Definition E.7.** Given the truncated loss function class  $S(R_T)$ , we define its  $\epsilon_c$ -covering with the minimum cardinality in the  $L^{\infty}$  metric as  $\mathcal{L}_{\mathcal{N}} := \{\ell_1, \ell_2, \dots, \ell_{\mathcal{N}}\}$ . Moreover, we define  $\ell_J \in \mathcal{L}_{\mathcal{N}}$  with random variable *J*. By definition, there exist  $\ell_J \in \mathcal{L}_{\mathcal{N}}$  such that  $\|\ell_J - \ell(x_i, y_i; s_W)\|_{\infty} \le \epsilon_c$ .

Note that Lemma E.3 provides the upper-bound on the  $\epsilon_c$ -covering number of  $S(R_T)$  for score estimator trained from transformer network class. Next, we bound the difference function.

**Lemma E.4** (Bound on Difference Function). Consider i.i.d training samples  $\{x_{0,i}, y_i\}_{i=1}^n$  and score estimator  $\hat{s}$ . Under generic Assumption 3.1 and parameter configuration in Theorem 3.1, it holds:

$$\Delta_n(\widehat{s}, s^\circ) \lesssim \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \widehat{\mathcal{R}}_m(\widehat{s}) \right] + \frac{1}{t_0} \left( R_{\mathcal{T}} \exp\left(-C_2 R_{\mathcal{T}}^2\right) + \frac{1}{n} \log \mathcal{N} \right) + 7\epsilon_c,$$

where  $\mathcal{N}(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2)$  is the covering number of transformer network class. Moreover, Under stronger Assumption 3.1 and parameter configuration in Theorem D.1, it holds:

$$\Delta_n(\widehat{s}, s^\circ) \lesssim \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \widehat{\mathcal{R}}_m(\widehat{s}) \right] + \log \frac{1}{t_0} \left( R_{\mathcal{T}} \exp\left(-C_2 R_{\mathcal{T}}^2\right) + \frac{1}{n} \log \mathcal{N} \right) + 7\epsilon_c.$$

### E.2 PROOF OF THEOREM 3.2

*Proof of Theorem 3.2.* For simplicity, we use  $\kappa = 1/t_0$  for the case in Theorem 3.1 and  $\kappa = \log(1/t_0)$  for the case in Theorem D.1. The proof proceeds through the following three steps.

### • Step A: Decompose the mixed risk.

We denote the ground truth by  $s^{\circ}(x, y, t) = \nabla \log p_t(x|y)$ . Moreover, if  $y = \emptyset$  we set  $s^{\circ}(x, y, t) = \nabla \log p_t(x)$ .

Recall Definition E.3 and Lemma E.4. By introducing a different set of i.i.d. samples  $\{x'_i, y'_i\}_{i=1}^n$  from the initial data distribution  $P_0(x, y)$  independent of the training samples, we rewrite the mixed risk:

$$\mathcal{R}_{m}(\widehat{s}) = \mathbb{E}_{\{x'_{i}, y'_{i}\}_{i=1}^{n}} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( \ell(x'_{i}, y'_{i}, \widehat{s}) - \ell(x'_{i}, y'_{i}, s^{\circ}) \right) \right] = \mathbb{E}_{\{x'_{i}, y'_{i}\}_{i=1}^{n}} \left[ \widehat{\mathcal{R}}'_{m}(\widehat{s}) \right],$$

where we use  $\widehat{\mathcal{R}}'_m(\widehat{s})$  to denote the empirical risk of the score estimator  $\widehat{s}$  trained from the i.i.d samples  $\{x'_i, y'_i\}_{i=1}^n$ .

This allows us to do the decomposition of  $\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\widehat{s})]$  as follows.

$$\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}[\mathcal{R}_m(\widehat{s})] = \underbrace{\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\mathbb{E}_{\{x_i',y_i'\}_{i=1}^n}\left[\widehat{\mathcal{R}}_m'(\widehat{s}) - \widehat{\mathcal{R}}_m'^{\mathrm{trunc}}(\widehat{s})\right]\right]}_{(\mathbf{I})}_{(\mathbf{I})} + \underbrace{\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\mathbb{E}_{\{x_i',y_i'\}_{i=1}^n}\left[\widehat{\mathcal{R}}_m'^{\mathrm{trunc}}(\widehat{s}) - \widehat{\mathcal{R}}_m^{\mathrm{trunc}}(\widehat{s})\right]\right]}_{(\mathbf{I})}_{(\mathbf{I})} + \underbrace{\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\widehat{\mathcal{R}}_m^{\mathrm{trunc}}(\widehat{s}) - \widehat{\mathcal{R}}_m(\widehat{s})\right]}_{(\mathbf{I})} + \underbrace{\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\widehat{\mathcal{R}}_m^{\mathrm{trunc}}(\widehat{s}) - \widehat{\mathcal{R}}_m(\widehat{s})\right]}_{(\mathbf{I})} + \underbrace{\mathbb{E}_{\{x_i,y_i\}_{i=1}^n}\left[\widehat{\mathcal{R}}_m(\widehat{s}) - \widehat{\mathcal{R}}_m(\widehat{s})\right]}_{(\mathbf{I})}}_{(\mathbf{I})}$$

## • Step B: Derive the Upper Bound.

## – Step B.1: Bound Each Term.

- \* By Lemma E.1, we have both (I), (III)  $\lesssim \kappa \exp(-C_2 R_T^2) R_T$ .
- \* By Lemma E.4, we have (II)  $\leq (IV) + \kappa \left( R_T \exp\left(-C_2 R_T^2\right) + \frac{1}{n} \log N \right) + 7\epsilon_c$ ,

\* By the following, we have (IV)  $\leq \min_{s_W \in \mathcal{T}_B^{h,s,r}} \mathcal{R}_m(s)$ .

$$(\mathbf{IV}) = \mathbb{E}_{\{z_i\}_{i=1}^n} \left[ \widehat{\mathcal{R}}(\widehat{s}) \right] \le \mathbb{E}_{\{z_i\}_{i=1}^n} \left[ \widehat{\mathcal{R}}_m(s) \right] = \mathcal{R}_m(s)$$

The inequality holds because  $\hat{s}$  is the minimizer of the empirical risk.

- Step B.2: Combine (I), (II), (III), (IV).

Combining these results we obtain

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\hat{s})] \leq 2 \min_{s_W \in \mathcal{T}_R^{h, s, r}} \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{x_t, y, \tau} \left[ \|s(x_t, \tau y, t) - \nabla \log p_t(x_t | \tau y)\|_2^2 \right] \mathrm{d}t \\ + \mathcal{O}\left(\frac{\kappa}{n} \log \mathcal{N}\right) + \mathcal{O}(\exp\left(-C_2 R_\mathcal{T}^2\right) \kappa) + \mathcal{O}\left(\epsilon_c\right).$$
(E.5)

By taking  $R_{\mathcal{T}} = \sqrt{\frac{(C_{\sigma} + 2\beta)\log N}{C_2(d_x + d_y)}}$  we have

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\widehat{s})] \leq 2 \min_{s \in \mathcal{T}_R^{h, s, r}} \int_{t_0}^T \frac{1}{T - t_0} \mathbb{E}_{\tau, x_t, y} \left[ \|s(x, \tau y, t) - \nabla \log p_t(x|y)\|_2^2 \right] dt$$
$$\mathcal{O}\left(\frac{\kappa}{n} \log \mathcal{N}\right) + \mathcal{O}\left(N^{-\frac{2\beta}{d_x + d_y}}\right) + \mathcal{O}\left(\epsilon_c\right). \tag{E.6}$$

where we use  $\kappa \lesssim \frac{1}{t_0} = N^{C_{\sigma}}$  by Lemma E.1 to obtain the third term on the RHS.

## Step C: Altogether.

To apply the previous approximation theorems (Theorem 3.1 and Theorem D.1) to the first term on the RHS of (E.5), we rewrite the expectation as

$$\mathbb{E}_{x_t,y,\tau} \left[ \| s(x_t,\tau y,t) - \nabla \log p_t(x_t | \tau y) \|_2^2 \right]$$
(E.7)  
=  $\frac{1}{2} \int_{\mathbb{R}^{d_x}} \| s(x,\emptyset,t) - \nabla \log p_t(x|y) \|_2^2 p_t(x) dx + \frac{1}{2} \mathbb{E}_y \left[ \int_{\mathbb{R}^{d_x}} \| s(x,y,t) - \nabla \log p_t(x|y) \|_2^2 p_t(x|y) dx \right]$ 

Since the marginal distribution  $p_t(x)$  also satisfies the subgaussian property, the previous result of the conditional score estimation applies to its unconditional counterpart by removing the label throughout the whole process.

- Step C.1: Result under generic Assumption 3.1.

By Theorem 3.1, we rewrite (E.6) as

$$\mathbb{E}_{\{z_i\}_{i=1}^n}[\mathcal{R}_m(\widehat{s})] \lesssim \underbrace{\mathcal{O}\left(N^{-\frac{\beta}{d_x+d_y}}(\log N)^{d_x+\frac{\beta}{2}+1}\right)}_{(\mathbf{i})} + \underbrace{\mathcal{O}\left(N^{-\frac{2\beta}{d_x+d_y}}\right)}_{(\mathbf{i}\mathbf{i})} + \underbrace{\mathcal{O}\left(\frac{\kappa}{n}\log\mathcal{N}\right)}_{(\mathbf{i}\mathbf{i}\mathbf{i})} + \underbrace{\mathcal{O}\left(\epsilon_c\right)}_{(\mathbf{i}\mathbf{v})}.$$

Moreover, from Lemma E.1 we have  $\kappa = O(1/t_0)$  and from Lemma E.3 we have

$$\log \mathcal{N}\left(\epsilon_{c}, \mathcal{S}(R_{\mathcal{T}}), \left\|\cdot\right\|_{\infty}\right) \lesssim \frac{\log n}{\epsilon_{c}^{2}} N^{\frac{68\beta}{d_{x}+d_{y}}+104C_{\sigma}} (\log N)^{12d_{x}+12\beta+2} (R_{\mathcal{T}})^{2}$$
$$\coloneqq \frac{\log n}{\epsilon_{c}^{2}} N^{\nu_{1}} (\log N)^{\nu_{2}} (R_{\mathcal{T}})^{2},$$

where  $\nu_1 = 68\beta/(d_x + d_y) + 104C_\sigma$  and  $\nu_2 = 12d_x + 12\beta + 2$ . By taking  $N = n^{\frac{d_x + d_y}{(d_x + d_y + \beta)}}$  and  $\epsilon_c = N^{-\frac{2\beta}{(d_x + d_y)}}$ , we have error: \* (i) =  $\mathcal{O}\left((\log n)^{d_x + \frac{\beta}{2} + 1}n^{-\frac{\beta}{(d_x + d_y + \beta)}}\right)$ . \* (ii) =  $\mathcal{O}\left(n^{-\frac{2\beta}{(d_x + d_y + \beta)}}\right)$  \* (iii) = 
$$\mathcal{O}\left(\kappa n^{-1} \cdot \underbrace{n^{\frac{4\beta}{d_x + d_y + \beta}}}_{\epsilon_c^{-2}} \cdot (\log n) \cdot \underbrace{n^{\frac{\nu_1(d_x + d_y)}{d_x + d_y + \beta}}}_{N^{\nu_1}} \cdot \underbrace{(\log n)^{\nu_2}}_{(\log N)^{\nu_2}} \cdot \underbrace{(\log n)}_{R_{\mathcal{T}}^2}\right)$$
 with  $\kappa = \frac{1/t_0}{1}$ .

Rearranging the expression, we have (iii) =  $\mathcal{O}\left(\frac{1}{t_0}n^{-\frac{(1-\nu_1)(d_x+d_y)-3\beta}{d_x+d_y+\beta}}(\log n)^{\nu_2+2}\right)$ \* (iv) =  $\mathcal{O}\left(n^{-\frac{2\beta}{d_x+d_y+\beta}}\right)$ 

\* (iv) = 
$$\mathcal{O}\left($$

We take the mixture of (i) and (iii) as the final error bound:

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \mathcal{R}(\hat{s}) \right] = \mathcal{O}\left( \frac{1}{t_0} n^{-\frac{\min\left(\beta, (1-\nu_1)(d_x+d_y)-3\beta\right)}{(d_x+d_y+\beta)}} (\log n)^{\nu_2+2} \right).$$

- Step C.2: Result under stronger Assumption 3.1.

With Theorem D.1, we further write (E.6) as

Moreover, by Lemma E.1 we have  $\kappa = \mathcal{O}(\log \frac{1}{t_0})$ , and by Lemma E.3 we have:

$$\log \mathcal{N}\left(\epsilon_{c}, \mathcal{S}(R_{\mathcal{T}}), \left\|\cdot\right\|_{\infty}\right) \lesssim \frac{\log n}{\epsilon_{c}^{2}} N^{\nu_{3}} (\log N)^{10} (R_{\mathcal{T}})^{2}.$$

where 
$$\nu_{3} = \frac{4(12\beta d_{x}+31\beta d+6\beta)}{d(d_{x}+d_{y})} + \frac{12(12C_{\alpha}d_{x}+25C_{\alpha}\cdot d+6C_{\alpha})}{d} + 72C_{\sigma}.$$
  
By taking  $N = n^{\frac{(d_{x}+d_{y})}{(d_{x}+d_{y}+2\beta)}}$  and  $\epsilon_{c} = N^{-\frac{2\beta}{(d_{x}+d_{y})}}$ , we have error:  
\* (i)  $= \mathcal{O}\left((\log n)^{\beta+1}n^{-\frac{2\beta}{(d_{x}+d_{y}+2\beta)}}\right).$   
\* (ii)  $= \mathcal{O}\left(n^{-\frac{2\beta}{(d_{x}+d_{y}+2\beta)}}\right).$   
\* (iii)  $= \mathcal{O}\left(\kappa n^{-1}\underline{n\frac{4\beta}{d_{x}+d_{y}+2\beta}}\cdot(\log n)\cdot\underline{n\frac{\nu_{3}(d_{x}+d_{y})}{N^{\nu_{3}}}}(\log n)^{10}\underline{(\log n)}\right)$  with  $\kappa = \log(1/t_{0}).$   
Rearranging the expression we have (iii)  $= \mathcal{O}\left(\log \frac{1}{t_{0}}n^{-\frac{(1-\nu_{3})(d_{x}+d_{y})-2\beta}{d_{x}+d_{y}+2\beta}}(\log n)^{12}\right).$   
\* (iv)  $= \mathcal{O}\left(n^{-\frac{2\beta}{d_{x}+d_{y}+2\beta}}\right).$   
We take the mixture of (i) and (iii) as the final error bound:

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[ \mathcal{R}(\hat{s}) \right] = \mathcal{O}\left( \log \frac{1}{t_0} n^{-\frac{\min\left(2\beta, (1-\nu_3)(d_x+d_y)-2\beta\right)}{(d_x+d_y+2\beta)}} (\log n)^{\max(12,\beta+1)} \right).$$

This completes the proof.