

ON STATISTICAL RATES OF CONDITIONAL DIFFUSION TRANSFORMERS: APPROXIMATION AND ESTIMATION

Anonymous authors

Paper under double-blind review

ABSTRACT

We investigate the approximation and estimation rates of conditional diffusion transformers (DiTs) with classifier-free guidance. We present a comprehensive analysis for “in-context” conditional DiTs under four common data assumptions. We show that both conditional DiTs and their latent variants lead to the minimax optimality of unconditional DiTs under identified settings. Specifically, we discretize the input domains into infinitesimal grids and then perform a term-by-term Taylor expansion on the conditional diffusion score function under Hölder smooth data assumption. This enables fine-grained use of transformers’ universal approximation through a more detailed piecewise constant approximation and hence obtains tighter bounds. Additionally, we extend our analysis to the latent setting under the linear latent subspace assumption. We not only show that latent conditional DiTs achieve lower bounds than conditional DiTs both in approximation and estimation, but also show the minimax optimality of latent unconditional DiTs. Our findings establish statistical limits for conditional and unconditional DiTs, and offer practical guidance toward developing more efficient and accurate DiT 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 distribution estimation guarantees for both conditional DiTs and their latent variants. We provide a comprehensive analysis under various data conditions. Moreover, we show that both conditional DiTs and their latent variants lead to the minimax optimality of unconditional DiTs under identified settings. This analysis is not only practical but also timely. Transformer-based conditional diffusion models are at the forefront of generative AI due to their success as scalable and flexible backbones for image (Wu et al., 2024a; Bao et al., 2023; Batzolis et al., 2021) and video generation (Liu et al., 2024; Ni et al., 2023; Saharia et al., 2022; Voleti et al., 2022). However, the theoretical understanding of conditional DiTs remains limited. On the one hand, while prior work by Hu et al. (2024) reports approximation and estimation rates of DiTs using the established universality of transformers (Yun et al., 2020), their results are not tight and are limited to unconditional diffusion. On the other hand, existing theoretical works on conditional diffusion models only focus on ReLU networks (Fu et al., 2024a; Yuan et al., 2023), model-free settings (Ye et al., 2024; Guo et al., 2024) or generative sampling process (Dinh et al., 2023), without considering the transformer architectures. This work addresses this gap by providing a timely analysis of the statistical limits of conditional DiTs.

In this work, we present a comprehensive analysis of conditional DiT and its latent setting under four common data assumptions. We also establish the minimax optimality of unconditional DiT and its latent version by deriving the tight distribution estimation error bounds. Our techniques include two key parts: (i) Discretizing the input domains into infinitesimal grids. (ii) On each grid, performing a term-by-term Taylor expansion on the conditional diffusion score function under generic and stronger Hölder smooth data assumptions, motivated by the local diffused polynomial analysis (Fu et al., 2024a; Oko et al., 2023). These techniques leverage the nice regularity of the score function imposed by the Hölder smoothness data assumptions and hence enable fine-grained use of transformers’ universal approximation (Kajitsuka and Sato, 2024; Yun et al., 2020) through a more detailed piecewise constant approximation. Consequently, we obtain tighter bounds.

Contributions. We summarize the theoretical results in Table 1. Our contributions are threefold:

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

Assumption	Score Approximation	Score Estimation	Dist. Estimation (Total Variation Distance)	Minimax Optimality
Generic Hölder Smooth Data Dist. (Sections 3.1 and 3.3)	$\mathcal{O}((\log(\frac{1}{\epsilon}))^{d_x}/\sigma_t^4)$	$n^{-\Theta(1/d_x)} \cdot (\log n)^{\mathcal{O}(d_x)}$	$n^{-\Theta(1/d_x)} \cdot (\log n)^{\mathcal{O}(d_x)}$	✗
Stronger Hölder Smooth Data Dist. (Sections 3.2 and 3.3)	$(\log(\frac{1}{\epsilon}))^{\mathcal{O}(1)}/\sigma_t^2$	$n^{-\Theta(1/d_x^2)} \cdot (\log n)^{\mathcal{O}(1)}$	$n^{-\Theta(1/d_x)} \cdot (\log n)^{\mathcal{O}(1)}$	✓
Latent Subspace + Generic Hölder Smooth Data Dist. (Section 4)	$\mathcal{O}((\log(\frac{1}{\epsilon}))^{d_0}/\sigma_t^4)$	$n^{-\Theta(1/d_0)} \cdot (\log n)^{\mathcal{O}(d_0)}$	$n^{-\Theta(1/d_0)} \cdot (\log n)^{\mathcal{O}(d_0)}$	✗
Latent Subspace + Stronger Hölder Smooth Data Dist. (Section 4)	$(\log(\frac{1}{\epsilon}))^{\mathcal{O}(1)}/\sigma_t^2$	$n^{-\Theta(1/d_0^2)} \cdot (\log n)^{\mathcal{O}(1)}$	$n^{-\Theta(1/d_0)} \cdot (\log n)^{\mathcal{O}(1)}$	✓

- **Score Approximation.** We characterize the approximation limit of matching the conditional DiT score function with a transformer-based score estimator. The approximation results explain the expressiveness of conditional DiT and its latent version, and guide the score network’s structural configuration for practical implementations (Theorems 3.1, 3.2 and 4.1). The results also show that the latent version achieves a better approximation for the score function.
- **Score and Distribution Estimation.** We study the score and distribution estimation of conditional DiTs in practical training scenarios. Specifically, we provide a sample complexity bound for score estimation (Theorems 3.3 and E.3), using norm-based covering number bound of transformer architecture. Additionally, we show that the learned score estimator can recover the initial data distribution in both conditional DiT and its latent setting (Theorems 3.4 and 4.2).
- **Minimax Optimal Estimator.** We extend our analysis to unconditional DiT and investigate whether the generated data distribution achieves the minimax optimality in the total variation distance. Specifically, we show that the upper bounds on the distribution estimation error match the lower bounds under stronger Hölder smooth data distribution (Corollary 3.4.2 and Remark 4.3).

Organization. Section 2 presents preliminaries and the problem setup. Section 3 presents the results of conditional DiTs. Section 4 presents the results of latent conditional DiTs. Appendix C.1 presents related works’ discussions. The appendix contains an extended and improved version of (Hu et al., 2024) on conditional DiTs (Appendix F), additional results, and detailed proofs.

Notations. The index set $\{1, \dots, I\}$ 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.

2 BACKGROUND AND PRELIMINARIES

In this section, we provide a high-level overview of the conditional diffusion model with classifier-free guidance in Section 2.1 and conditional Diffusion Transformer (DiT) networks in Section 2.2.

2.1 CONDITIONAL DIFFUSION MODEL WITH CLASSIFIER-FREE GUIDANCE

Forward and Backward Conditional Diffusion Process. In the *forward* process, conditional diffusion models gradually add noise to the original data $x_0 \in \mathbb{R}^{d_x}$. Give a condition $y \in \mathbb{R}^{d_y}$, and $x_0 \sim P_0(\cdot|y)$. Let x_t denote the noisy data at the timestamp t , with marginal distribution and density as $P_t(\cdot|y)$ and $p_t(\cdot|y)$. The conditional distribution $P_t(x_t|y)$ follows $N(\alpha_t x_0, \sigma_t^2 I_{d_x})$, where $\alpha_t = e^{-t/2}$, $\sigma_t^2 = 1 - e^{-t}$, and $w(t) > 0$ is a nondecreasing weighting function. In practice, the forward process terminates at a large enough T such that P_T is close to $N(0, I_{d_x})$. In the *backward* process, we obtain x_t^- by reversing the forward process. The generation of x_t^- depends on the score function $\nabla \log p_t(\cdot|y)$. See Appendix G.1 for the details. In below, when the context is clear, we suppress the notation dependence of x_t on the time step t .

Classifier-Free Guidance. Classifier-free guidance (Ho and Salimans, 2022) is the standard workhorse for training condition diffusion models. It approximates both conditional and unconditional score functions using neural networks s_W with parameters W . It uses the following loss function:

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

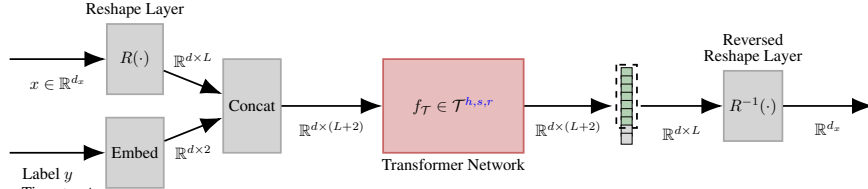


Figure 1: **Conditional DiT Network Architecture.** The architecture consists of a reshape layer $R(\cdot)$, a reversed reshape layer $R^{-1}(\cdot)$, and the embedding layers for label y and timestep t . The embeddings of y and t are concatenated with input sequences and then processed by a transformer network $f_{\mathcal{T}} \in \mathcal{T}^{h,s,r}$.

where $\nabla_{x_t} \log \phi_t(x_t|x_0) = -(x_t - \alpha_t x_0)/\sigma_t^2$, t_0 is a small cutoff to stabilize training¹. $\tau = \emptyset$ denotes the unconditional version, $\tau = \text{id}$ denotes the conditional version, and $P(\tau = \emptyset) = P(\tau = \text{id}) = 0.5$. To train s_W , we select n i.i.d. samples $\{x_{0,i}, y_i\}_{i=1}^n$, where $x_{0,i} \sim P_0(\cdot|y_i)$. We use

$$\widehat{\mathcal{L}}(s_W) := \frac{1}{n} \sum_{i=1}^n \ell(x_{0,i}, y_i; s_W), \quad (2.1)$$

as the empirical loss. In addition, we denote population loss as $\mathcal{L}(s_W)$. See Appendix G.2 for details.

2.2 CONDITIONAL DIFFUSION TRANSFORMER NETWORKS

We use a transformer network as a score estimator s_W . Our notation follows (Hu et al., 2024).

Transformer Block. Let $f^{(\text{SA})} : \mathbb{R}^{d \times L} \rightarrow \mathbb{R}^{d \times L}$ denote the self-attention layer. Let h and s denote the number of heads and hidden dimension in the self-attention layer, and then we have

$$f^{(\text{SA})}(Z) := Z + \sum_{i=1}^h W_O^i (W_V^i Z) \text{Softmax} [(W_K^i Z)^\top (W_Q^i Z)], \quad (2.2)$$

where $W_V^i, W_K^i, W_Q^i \in \mathbb{R}^{s \times d}$, and $W_O^i \in \mathbb{R}^{d \times s}$ are the weight matrices. Next, we define the feed-forward layer with MLP dimension r :

$$f^{(\text{FF})}(Z) := Z + W_2 \text{ReLU}(W_1 Z + b_1) + b_2, \quad (2.3)$$

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

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

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

Now, we define the transformer networks as compositions of transformer blocks.

Definition 2.2 (Transformer Network Function Class). Let $\mathcal{T}^{h,s,r}$ denote the transformer network function class where each function $\tau \in \mathcal{T}^{h,s,r}$ is a composition of transformer blocks $f^{h,s,r}$, i.e.,

$$\mathcal{T}^{h,s,r} := \{\tau : \mathbb{R}^{d \times L} \mapsto \mathbb{R}^{d \times L} \mid \tau = f^{h,s,r} \circ \dots \circ f^{h,s,r}\}.$$

Conditional Diffusion Transformer (DiT). Let $f \in \mathcal{T}^{h,s,r}$ be a transformer network, and $(x, y, t) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \times [t_0, T]$ be the input data. We follow the ‘‘in-context conditioning’’ conditional DiT network in (Peebles and Xie, 2023) as in Figure 1. The following reshape layer converts a vector input $x \in \mathbb{R}^{d_x}$ into the sequential matrix input format $Z \in \mathbb{R}^{d \times L}$ for transformer with $d_x = d \cdot L$.

Definition 2.3 (DiT Reshape Layer $R(\cdot)$). Let $R(\cdot) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d \times L}$ be a reshape layer that transforms the d_x -dimensional input into a $d \times L$ matrix. Specifically, for any $d_x = i \times i$ image input, $R(\cdot)$ converts it into a sequence representation with feature dimension $d := p^2$ (where $p \geq 2$) and sequence length $L := (i/p)^2$. Besides, we define the corresponding reverse reshape (flatten) layer $R^{-1}(\cdot) : \mathbb{R}^{d \times L} \rightarrow \mathbb{R}^{d_x}$ as the inverse of $R(\cdot)$. By $d_x = dL$, R, R^{-1} are associative w.r.t. their input.

We define the following transformer network function class with the reshape layer.

Definition 2.4 (Transformer Network Function Class with Reshape Layer $\mathcal{T}_R^{h,s,r}$).

$\mathcal{T}_R^{h,s,r}(C_{\mathcal{T}}, C_Q^{2,\infty}, C_Q, C_K^{2,\infty}, C_K, C_V^{2,\infty}, C_V, C_O^{2,\infty}, C_O, C_E, C_{f_1}^{2,\infty}, C_{f_1}, C_{f_2}^{2,\infty}, C_{f_2}, L_{\mathcal{T}})$ satisfies

- $\mathcal{T}_R^{h,s,r} := \{R^{-1} \circ f_{\mathcal{T}} \circ R : \mathbb{R}^{d_x} \mapsto \mathbb{R}^{d_x} \mid f_{\mathcal{T}} \in \mathcal{T}^{h,s,r}\};$

¹ t_0 is the early stopping time to prevent the score function from blowing up (Fu et al., 2024a; Chen et al., 2023c; Dhariwal and Nichol, 2021; Song et al., 2021).

- Model output bound: $\sup_Z \|f_{\mathcal{T}}(Z)\|_2 \leq C_{\mathcal{T}}$;
- Parameter bound in $f^{(\text{SA})}$: $\|(W_Q)^\top\|_{2,\infty} \leq C_Q^{2,\infty}$, $\|(W_Q)^\top\|_2 \leq C_Q$, $\|W_K\|_{2,\infty} \leq C_K^{2,\infty}$, $\|W_K\|_2 \leq C_K$, $\|W_V\|_{2,\infty} \leq C_V^{2,\infty}$, $\|W_V\|_2 \leq C_V$, $\|W_O\|_{2,\infty} \leq C_O^{2,\infty}$, $\|W_O\|_2 \leq C_O$, $\|E^\top\|_{2,\infty} \leq C_E$;
- Parameter bound in $f^{(\text{FF})}$: $\|W_1\|_{2,\infty} \leq C_{f_1}^{2,\infty}$, $\|W_1\|_2 \leq C_{f_1}$, $\|W_2\|_{2,\infty} \leq C_{f_2}^{2,\infty}$, $\|W_2\|_2 \leq C_{f_2}$;
- Lipschitz of $f_{\mathcal{T}} \in \mathcal{T}^{h,s,r}$: $\|f_{\mathcal{T}}(Z_1) - f_{\mathcal{T}}(Z_2)\|_F \leq L_{\mathcal{T}}\|Z_1 - Z_2\|_F$, for any $Z_1, Z_2 \in \mathbb{R}^{d \times L}$.

These norm bounds are critical to quantify the interplay between model, performance and data.

3 STATISTICAL LIMITS OF CONDITIONAL DiTs

In this section, we present a refined decomposition scheme for the fine-grained analysis of score approximation, score estimation, and distribution estimation in conditional DiT. Our analysis considers two assumptions on initial data distributions: (i) a generic Hölder smooth data assumption (Section 3.1 for approximation, and Section 3.3 for estimation), (ii) a stronger Hölder smooth data assumption (Section 3.2 for approximation, and Section 3.3 for estimation). This new scheme leads to tighter bounds, including the minimax optimality of the unconditional DiT score estimator.

3.1 SCORE APPROXIMATION: GENERIC HÖLDER SMOOTH DATA DISTRIBUTIONS

We present a fine-grained piecewise approximation using transformers to approximate the conditional score function under the Hölder smoothness assumption on the initial data (Fu et al., 2024b). At its core, we introduce a score function decomposition scheme with term-by-term tractability.

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 let $\beta = k_1 + \gamma$ denote the smoothness parameter, where $k_1 = \lfloor \beta \rfloor$ and $\gamma \in [0, 1)$. For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the Hölder space $\mathcal{H}^\beta(\mathbb{R}^d)$ is defined as the set of α -differentiable functions satisfying: $\mathcal{H}^\beta(\mathbb{R}^d) := \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \|f\|_{\mathcal{H}^\beta(\mathbb{R}^d)} < \infty\}$, where the Hölder norm $\|f\|_{\mathcal{H}^\beta(\mathbb{R}^d)}$ satisfies:

$$\|f\|_{\mathcal{H}^\beta(\mathbb{R}^d)} := \max_{\alpha: \|\alpha\|_1 < k_1} \sup_x |\partial^\alpha f(x)| + \max_{\alpha: \|\alpha\|_1 = k_1} \sup_{x \neq x'} \frac{|\partial^\alpha f(x) - \partial^\alpha f(x')|}{\|x - x'\|_2^\gamma}.$$

We also define the Hölder ball of radius B : $\mathcal{H}^\beta(\mathbb{R}^d, B) := \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \|f\|_{\mathcal{H}^\beta(\mathbb{R}^d)} < B\}$.

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 first assumption on the conditional distribution of initial data x_0 .

Assumption 3.1 (Generic 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 $p_0(x_0|y) \in \mathcal{H}^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$ (see Definition 3.1 for precise definition.) Also, for any $y \in [0, 1]^{d_y}$, there exist positive constants C_1, C_2 such that $p_0(x_0|y) \leq C_1 \exp(-C_2\|x_0\|_2^2/2)$.

Remark 3.1. The Hölder continuity assumption captures various smoothness levels in the conditional density function. The light-tail condition relaxes the bounded support assumption in (Oko et al., 2023). Moreover, Assumption 3.1 only applies to the initial conditional distribution and imposes no constraints on the induced conditional score function. This is far less restrictive than the Lipschitz score condition in prior works (Yuan et al., 2024; Lee et al., 2023; Chen et al., 2022).

In our work, we aim to approximate the conditional score function $\nabla \log p_t(x_t|y)$ using transformer architectures. Hu et al. (2024) analyze the unconditional DiTs based on the established universality of transformers (Yun et al., 2020). These theories discretize the input and output domains into infinitesimal grids and employ piecewise constant approximations to construct universal approximators with controllable errors. However, such methods do not yield tight bounds for DiT architectures (Hu et al., 2024). To combat this, we build on the key observation by Fu et al. (2024a)²:

$$p_t(x_t|y) = \int_{\mathbb{R}^{d_x}} \frac{dx_0}{\sigma_t^{d_x} (2\pi)^{d_x/2}} \cdot \underbrace{p_0(x_0|y)}_{\approx k_1\text{-order Taylor polynomial}} \cdot \underbrace{\exp\left(-\frac{\|\alpha_t x_0 - x_t\|^2}{2\sigma_t^2}\right)}_{\approx k_2\text{-order Taylor polynomial}}. \quad (3.1)$$

²Recall that $p_t(x_t|y) = \int_{\mathbb{R}^{d_x}} p(x_0|y)p_t(x_t|x_0) dx_0$ with $P_t(\cdot|y) \sim N(\alpha_t x_0, \sigma_t I_{d_x})$. In below, when the context is clear, we suppress the notation dependence of x_t on the time step t .

A term-by-term Taylor expansion of the above conditional distribution under [Assumption 3.1](#) enables a more fine-grained analysis (e.g., [Lemma I.2](#)). As a result, we propose a *fine-grained version of piecewise constant approximation* for conditional DiTs, allowing transformers to approximate the conditional score function with tighter error bounds. In particular, we utilize a refined transformer universal approximation modified from ([Kajitsuka and Sato, 2024](#)) (see [Appendix H.1](#) for details).

Our score approximation procedure has two stages: first, we approximate p_t and ∇p_t using a Taylor expansion, then use transformers to approximate p_t , ∇p_t , and the required algebraic operators to construct $\nabla \log p_t(x|y) = \frac{\nabla p_t(x|y)}{p_t(x|y)}$. These lead to provably tight estimation results in [Section 3.3](#).

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](#)). Assume [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

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

Notably, for $\epsilon = \mathcal{O}(N^{-\beta})$, the approximation error has the upper bound $\mathcal{O}((\log(\frac{1}{\epsilon}))^{d_x}/\sigma_t^4)$.

The parameter bounds for the transformer network class are as follows:

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

Remark 3.2. N is the resolution of the input domain discretization (see [Lemma I.2](#)). We remark that domain discretization is essential for utilizing the local smoothness of functions under Hölder assumptions. C_σ and C_α control the stability cutoff and early stopping time, respectively.

Proof Sketch. Recall that $\nabla \log p_t(x|y) = \frac{\nabla p_t(x|y)}{p_t(x|y)}$. We employ the following strategy: discretize the domains, apply a term-by-term Taylor approximation to the decomposed conditional distribution (3.1), decompose the conditional score function $\nabla \log p_t(x|y) = \frac{\nabla p_t(x|y)}{p_t(x|y)}$ into two fundamental functions and a parsimonious set of algebraic operators, and then approximate the fundamental functions and operators with transformer networks. The resulting joint error of this strategy is controllable under [Assumption 3.1](#). Our proof follows three steps:

Step 1. Input Domains Discretization. For any $x \in \mathbb{R}^{d_x}$, we construct a bounded domain $B_{x,N}$ to approximate polynomial functions evaluated at x on \mathbb{R}^{d_x} with the same functions on $B_{x,N}$ to arbitrary precision $1/N$ ([Lemma I.1](#)). Then, we discretize $B_{x,N} \times [0, 1]^{d_y}$ into $N^{d_x+d_y}$ hypercubes ([Lemma I.2](#)). This technique confines the approximation to a compact domain by controlling error outside this domain under [Assumption 3.1](#). Each hypercube is now compact and local, enabling a well-behaved Taylor expansion at x . This confinement reduces approximation error in **Step 2**.

Step 2. Local, Term-by-Term Taylor Expansion for $\nabla \log p_t$. To approximate $\nabla \log p_t$, we expand $p_t(x|y)$ and $\nabla p_t(x|y)$ with Taylor polynomials on each local grid on $B_{x,N}$, following the term-by-term expansion (3.1). Specifically, we approximate $p_t(x|y)$ with a scalar polynomial function $f_1(x, y, t) \in \mathbb{R}$ ([Lemma I.3](#)) and $\nabla p_t(x|y)$ with a vector-valued polynomial function $f_2(x, y, t) \in \mathbb{R}^{d_x}$ ([Lemma I.4](#)). Together with a parsimonious set of algebraic operators (inverse, product), the obtained f_1, f_2 resemble $\nabla \log p_t$ with a bounded error $\text{Error}_{\text{Taylor}}$.

Step 3. Term-by-Term Approximations with Transformers. We utilize a refined universal approximation theorem for transformers ([Appendix H.1](#)) to approximate all Taylor-expanded terms: f_1, f_2 , and the set of algebraic operators. Specifically, we approximate $f_1(x, y, t)$ and $f_2(x, y, t)$ with transformer models \mathcal{T}_{f_1} ([Lemma I.5](#)) and \mathcal{T}_{f_2} ([Lemma I.6](#)). For the operators, we also approximate each of them with a corresponding transformer \mathcal{T}_μ with $\mu = \{\text{inverse, square} \dots\}$ ([Lemmas I.8](#)

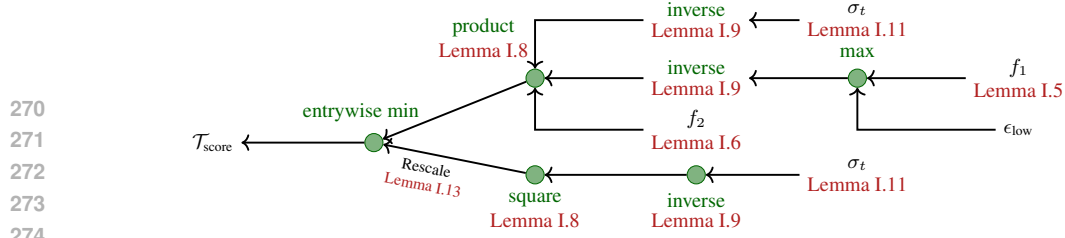


Figure 2: **Approximate Score Function with Transformer $\mathcal{T}_{\text{score}}$ under Assumption 3.1.** The construction consists of the transformers to approximate local polynomials f_1 and f_2 , and the algebraic operators. We highlight the overall term-by-term approximations and their corresponding lemmas to ensemble the transformers to I.9 and I.11). All approximations have precision guarantees. Finally, we combine the transformer approximations \mathcal{T}_{f_1} , \mathcal{T}_{f_2} and \mathcal{T}_μ for the set of algebraic operators, resulting in a joint approximation for $\nabla \log p_t$ (see Figure 2) with arbitrary small error $\text{Error}_{\mathcal{T}}$.

Error Matching. The overall error includes $\text{Error}_{\text{Taylor}}$ and $\text{Error}_{\mathcal{T}}$. Given a fixed discretization resolution N , $\text{Error}_{\text{Taylor}}$ remains fixed. However, the approximation error bound of the transformer can be an arbitrary value. We align $\text{Error}_{\mathcal{T}}$ and $\text{Error}_{\text{Taylor}}$ to optimize the final results.

Please see Appendix I for a detailed proof. \square

Remark 3.3 (Approximation Rate). Given a fixed resolution N , the approximation error scales inversely with the smoothness β . As the smoothness increases, we get a tighter approximation error.

Remark 3.4 (Comparing with Existing Works). Fu et al. (2024a) provide approximation rates for conditional diffusion models using ReLU networks. We are the first to establish approximation error bounds with transformer networks. Additionally, Oko et al. (2023) establish approximation rates under a compactness condition on the input data. We mitigate this compactness requirement by applying a Hölder smoothness assumption to control approximation error outside a compact domain.

3.2 SCORE APPROXIMATION: STRONGER HÖLDER SMOOTH DATA DISTRIBUTIONS

Next, we study the conditional DiT score approximation problem using our score decomposition scheme under the stronger Hölder smoothness assumption from Fu et al. (2024b, Assumption 3.3).

Assumption 3.2 (Stronger Hölder Smooth Data). Let function $f \in \mathcal{H}^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$. Given a constant radius B , positive constants C and C_2 , we assume the conditional density function $p(x_0|y) = \exp\left(-C_2\|x_0\|_2^2/2\right) \cdot f(x_0, y)$ and $f(x_0, y) \geq C$ for all $(x_0, y) \in \mathbb{R}^{d_x} \times [0, 1]^{d_y}$.

Assumption 3.2 imposes stronger assumption than Assumption 3.1 and induces a refined conditional score function decomposition. Explicitly, by Lemma J.1, $\nabla \log p_t(x|y)$ becomes:

$$\nabla \log p_t(x|y) = \frac{-C_2 x}{\alpha_t^2 + C_2 \sigma_t^2} + \frac{\nabla h(x, y, t)}{h(x, y, t)}, \quad (3.2)$$

where $h(x, y, t) := \int_{\mathbb{R}^{d_x}} \frac{f(x_0, y)}{\hat{\sigma}_t^{d_x} (2\pi)^{d_x/2}} \exp\left(-\frac{\|x_0 - \hat{\alpha}_t x\|_2^2}{2\hat{\sigma}_t^2}\right) dx_0$, $\hat{\sigma}_t = \frac{\sigma_t}{\sqrt{\alpha_t^2 + C_2 \sigma_t^2}}$, and $\hat{\alpha}_t = \frac{\alpha_t}{\alpha_t^2 + C_2 \sigma_t^2}$.

We highlight that (3.2) leads to a tighter approximation error compared with Theorem 3.1. Intuitively, Assumption 3.2 imposes a lower bound on the conditional density function and hence implies in better regularity of the score function. In contrast, under Assumption 3.1, the score function lacks such regularity and may explode when p_t is small. These low-density regions act as holes in the data support. They cause the score function to diverge near the boundary of these holes. To combat this, an implication of (3.2) is handy — h is bounded from zero, ensuring that the score function remains well-behaved across the entire data domain. To elaborate more, two technical remarks are in order.

Remark 3.5 (Linearity). The first term on the RHS of (3.2) is linear in x . This makes part of $\nabla \log p_t(x|y)$ a linear function of x , enabling easy approximation with a tighter bound.

Remark 3.6 (Tightened Approximation Induced by h 's Lower Bound). Moreover, the introduction of h tightens the approximation error due to the lower bound imposed by Assumption 3.2 (i.e., $f(x, y) \geq C$). The second term on the RHS of (3.2) mirrors the form $\nabla \log p_t(x|y) = \frac{\nabla p_t(x|y)}{p_t(x|y)}$ by replacing p with h . In the analysis of Section 3.1, especially in Step 2 of the proof (resembling f_1 , f_2 to approximate $\nabla p_t(x|y)$), we have to impose a threshold on the denominator of $\frac{\nabla p_t(x|y)}{p_t(x|y)}$ to prevent score explosion under Assumption 3.1. This threshold introduces additional approximation error (Lemma I.13). Assumption 3.2 remedies this by ensuring a lower bound on $p_t(x|y)$ through the minimum values of $f(x, y)$ and $\exp(-C_2\|x\|_2^2/2)$ within the compact domain after discretization. Setting this lower bound eliminates the need for a threshold and improves the approximation.

Consequently, decomposition (3.2) improves our approximation result from Section 3.1. We state our main result of score approximation using transformers under Assumption 3.2 as follows:

Theorem 3.2 (Conditional Score Approximation under Assumption 3.2 (Informal Version of Theorem J.1)). Assume Assumption 3.2. 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

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

Notably, for $\epsilon = \mathcal{O}(N^{-\beta})$, the approximation error has the upper bound $(\log(\frac{1}{\epsilon}))^{\mathcal{O}(1)}/\sigma_t^2$.

Proof Sketch. Our proof follows Theorem 3.1, but uses a different conditional score function decomposition in the form of (3.2). We highlight key differences in each corresponding step:

Step 0: Score Decomposition and Bounds on h and ∇h . We decompose ∇p_t to the form of (3.2) by Lemma J.1. Different from Section 3.2, we derive a lower bound on h in Lemma J.2.

Step 1: Input Domains Discretization. This step remains the same as Section 3.1, except the approximation target changes from $p, \nabla p$ to $h, \nabla h$. We confine and discretize input domains $\mathbb{R}^{d_x} \times [0, 1]^{d_y}$ into $N^{d_x+d_y}$ hypercubes (Lemma I.2), each supporting well-behaved Taylor expansions.

Step 2: Local, Term-by-Term Taylor Expansion for h and ∇h . Similar to Section 3.1, we utilize Taylor polynomials f_1 and f_2 to approximate h and ∇h on obtained hypercubes. The approximation on h and ∇h differs from approximation on p_t and ∇p_t , as their boundedness eliminates the need for a threshold to prevent score function blow-up. This leads to a faster approximation rate.

Step 3: Transformer Network Approximation. Similar to Section 3.1, we approximate polynomial functions f_1, f_2 and all necessary algebraic operators to construct an approximator f_3 for ∇p_t :

$$f_3(x, y, t) = -\frac{C_2 x}{\alpha_t^2 + C_2 \sigma_t^2} + \frac{\hat{\alpha}_t}{\hat{\sigma}_t} \cdot \frac{f_2(x, y, t)}{f_1(x, y, t)}, \quad (3.3)$$

following (3.2). Differed from Section 3.1, (3.2) requires transformers to approximate two additional operators, $\hat{\sigma}_t$ and $\hat{\alpha}_t$. All approximations have precision guarantees. Finally, we combine all transformer approximations required in (3.3) and obtain a joint approximation error for $\nabla \log p_t$ (see Figure 5) with arbitrary precision. We complete the proof by matching the approximation errors of the Taylor polynomial and transformer. Importantly, second term on the RHS of (3.3) manifests a tighter bound than that of $\frac{\nabla p_t(x|y)}{p_t(x|y)}$. The first linear-in- x term achieves a even tighter bound due to its linearity. Combined, we obtain a smaller overall joint approximation error than Theorem 3.1.

Please see Appendix J for a detailed proof, and see Theorem J.1 for the formal version. \square

Remark 3.7 (Comparing with Theorem 3.1). Let $\tilde{\mathcal{O}}(\cdot)$ hide the terms about $t_0, \log t_0, \log n$. In Theorem 3.2, the approximation rate $\tilde{\mathcal{O}}(N^{-\frac{2\beta}{d_x+d_y}})$ is faster than that of Theorem 3.1, i.e., $\tilde{\mathcal{O}}(N^{-\frac{\beta}{d_x+d_y}})$.

3.3 SCORE ESTIMATION AND DISTRIBUTION ESTIMATION OF CONDITIONAL DiTs

Next, we study score and distribution estimations based on the two score approximation results for two different data assumptions: Theorems 3.1 and 3.2. Let \hat{s} denote the trained score estimator.

Score Estimation. Building on our approximation results from Sections 3.1 and 3.2, the next objective is to evaluate the performance of the score estimator \hat{s} trained with a set of finite samples by optimizing the empirical loss (2.1). 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 risk as the expectation of the squared ℓ_2 difference between the score estimator and the ground truth with respect to (x_t, y, t) :

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

Given a set of i.i.d sample $\{x_i, y_i\}_{i \in [n]}$, direct computation of $\mathbb{E}_{\{x_i, y_i\}_{i \in [n]}}[\mathcal{R}(\hat{s})]$ is infeasible due to the absence of access to the joint distribution $P(x_t, y)$. To address this, we: (i) Decompose the risk into estimation and approximation errors, (ii) Bound the estimation error using the covering number of transformers, and (iii) Bound the approximation error using Theorem 3.1 and Theorem 3.2.

Theorem 3.3 (Conditional Score Estimation with Transformer). Assume $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$.

- Under [Assumption 3.1](#), by taking $N = n^{\frac{1}{\nu_1} \cdot \frac{d_x + d_y}{\beta + d_x + d_y}}$, $t_0 = N^{-C_\sigma} < 1$ and $T = C_\alpha \log n$, it holds

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

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

- Under [Assumption 3.2](#), by taking $N = n^{\frac{1}{\nu_3} \cdot \frac{d_x + d_y}{2\beta + d_x + d_y}}$, $t_0 = N^{-C_\sigma} < 1$ and $T = C_\alpha \log n$, it holds

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] = \mathcal{O}\left(\log \frac{1}{t_0} n^{-\frac{1}{\nu_3} \cdot \frac{2\beta}{d_x + d_y + 2\beta}} (\log n)^{\max(10, \beta + 1)}\right),$$

where $\nu_3 = \frac{4(12\beta d_x + 31\beta d + 6\beta)}{d(d_x + d_y)} + \frac{12(12C_\alpha d_x + 25C_\alpha \cdot d + 6C_\alpha)}{d} + 72C_\sigma$.

Corollary 3.3.1 (Low-Dimensional Input Region). Assume $d_x = o\left(\frac{\log N}{\log \log N}\right)$, i.e., $d_x \ll n$. Under [Assumption 3.1](#), by setting N, t_0, T as specified in [Theorem 3.3](#), we have $\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] = \mathcal{O}\left(\frac{1}{t_0} n^{-\frac{1}{\nu_4} \cdot \frac{\beta}{d_x + d_y + \beta}}\right)$, where $\nu_4 = \frac{72\beta(2d_x + 5d + 1)}{d(d_x + d_y)} + \frac{48C_\sigma(2d_x + 5d + 1)}{d} - 4\beta$.

Proof. Please see [Appendix K.2](#) and [Appendix K.4](#) for detailed proofs. \square

Remark 3.8 (Sample Complexity Bounds). To obtain ϵ -error in terms of score estimation, we have the sample complexity $\tilde{\mathcal{O}}\left(\epsilon^{-\nu_1(d_x + d_y + \beta)/\beta}\right)$ under [Assumption 3.1](#) and $\tilde{\mathcal{O}}\left(\epsilon^{-\nu_3(d_x + d_y + 2\beta)/2\beta}\right)$ under [Assumption 3.2](#). Here $\tilde{\mathcal{O}}(\cdot)$ ignores the terms about t_0 , $\log t_0$ and $\log n$. The Hölder data smoothness degree β affects the sample complexity. This indicates that the regularity of the initial data distribution determines the complexity of score estimation.

Distribution Estimation. Next, we study the distributional estimation capability of the trained conditional score network $s(x, y, t)$ by analyzing the total variation distance between the estimated and true distributions. Our strategy uses a three-part decomposition: (i) the total variation between the true distributions at timestamps 0 and t_0 , (ii) the total variation between the true distribution at t_0 and the reverse process distribution using the true score function, and (iii) the total variation between the reverse process distributions using the true and estimated score functions at t_0 .

Theorem 3.4 (Conditional Distribution Estimation). Assume $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$. For $y \in [0, 1]^{d_y}$, let $\hat{P}_{t_0}(\cdot|y)$ denote *estimated* conditional distributions at t_0 . Recall that $P_0(\cdot|y)$ is the conditional distribution of initial data x_0 given y . Assume $\text{KL}(P_0(\cdot|y) | N(0, I)) \leq c$ for some constant $c < \infty$.

- Under [Assumption 3.1](#), by taking the early-stopping time $t_0 = n^{-\frac{\beta}{d_x + d_y + \beta}}$ and terminal time $T = \frac{2\beta}{d_x + d_y + 2\beta} \log n$, it holds

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\hat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O}\left(n^{-\frac{\beta}{2(\nu_1 - 1)(d_x + d_y + \beta)}} (\log n)^{\frac{\nu_2}{2} + \frac{3}{2}}\right),$$

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

- Under [Assumption 3.2](#), by taking $t_0 = n^{-\frac{4\beta}{d_x + d_y + 2\beta} - 1}$, it holds

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\hat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O}\left(n^{-\frac{1}{2\nu_3} \frac{\beta}{d_x + d_y + 2\beta}} (\log n)^{\max(6, \frac{\beta}{2} + \frac{3}{2})}\right),$$

where $\nu_3 = \frac{4(12\beta d_x + 31\beta d + 6\beta)}{d(d_x + d_y)} + \frac{12(12C_\alpha d_x + 25C_\alpha \cdot d + 6C_\alpha)}{d} + 72C_\sigma$ and $C_\alpha = \frac{2\beta}{d_x + d_y + 2\beta}$.

We remark that the choice of t_0, T (i.e., C_σ, C_α) leads to the tightest rates in our analysis.

Corollary 3.4.1 (Low-Dimensional Input Region). Assume $d_x = o\left(\frac{\log N}{\log \log N}\right)$, i.e., $d_x \ll n$. Under [Assumption 3.1](#), by setting t_0, T as specified in [Theorem 3.4](#), we have

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\hat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O}\left(n^{-\frac{\beta}{2(\nu_4 + 1)(d_x + d_y + \beta)}}\right),$$

where $\nu_4 = \frac{72\beta(2d_x + 5d + 1)}{d(d_x + d_y)} + \frac{48C_\sigma(2d_x + 5d + 1)}{d} - 4\beta$.

Proof. Please see [Appendix K.6](#) for a detailed proof. \square

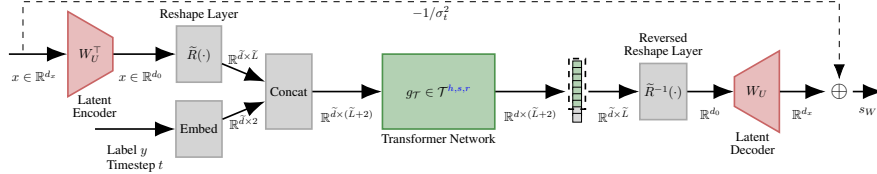


Figure 3: **Network Architecture of Latent Conditional DiT.** The overall architecture consists of linear layer of encoder and decoder W_U^\top and W_U that transform input $x \in \mathbb{R}^{d_x}$ into linear latent space \mathbb{R}^{d_0} , reshaping layer $\tilde{R}(\cdot)$ and $\tilde{R}^{-1}(\cdot)$, embedding layer for label y and timestep t . The embedding concatenates with input sequences and processes by the adapted transformer network $\mathcal{T}_{\tilde{R}}^{h,s,r} = \tilde{R}^{-1} \circ g_{\mathcal{T}} \circ f^{(\text{FF})} \circ \tilde{R}$.

3.4 MINIMAX OPTIMAL ESTIMATION OF UNCONDITIONAL DiTs

In this section, we show the minimax optimality of the unconditional DiT architecture under [Assumption 3.2](#). Specifically, we obtain the distribution estimation error of unconditional DiTs by removing the condition y and let $d_y = 0$ in [Theorem 3.4](#). Then the distribution estimation error becomes $\tilde{\mathcal{O}}(\epsilon^{-\frac{1}{2\nu_3} \frac{\beta}{d_x + 2\beta}})$ under [Assumption 3.2](#). Here $\tilde{\mathcal{O}}(\cdot)$ ignores the term about $\log n$. By setting $2\nu_3 = 1$, we show that the unconditional DiT is the minimax optimal distribution estimator.

Corollary 3.4.2 (Proposition 4.3 of [Fu et al. \(2024b\)](#)). For a fixed constant C_2 and a Hölder index $\beta > 0$. We consider the task of estimating a probability distribution $P(x)$ with its density function defined within the following function space

$$\mathcal{P} = \left\{ p(x) = f(x) \exp\left(-C_2 \|x\|_2^2\right) : f(x) \in \mathcal{H}^\beta(\mathbb{R}^{d_x}, B), f(x) \geq C \geq 0 \right\},$$

Given n i.i.d data $\{x_i\}_{i=1}^n$, we have $\inf_{\hat{\mu}} \sup_{P \in \mathcal{P}} \mathbb{E} \{ \text{TV}(\hat{\mu}, P) \} \geq \Omega(n^{-\frac{\beta}{d_x + 2\beta}})$. Here, the estimator $\hat{\mu}$ ranges over all possible estimators constructed from the data.

Remark 3.9 (Comparing with Existing Works). [Oko et al. \(2023\)](#) analyze the ReLU network and provide the near minimax optimal estimation rates in both the total variation distance and Wasserstein distance of order one. [Fu et al. \(2024b\)](#) also uses the ReLU network and provides the minimax optimality for distribution in total variation. Our results offer the first and exact minimax optimal guarantee for unconditional DiTs in distribution estimation.

4 LATENT CONDITIONAL DiTs

In this section, we extend the results from [Section 3](#) by considering the latent conditional DiTs. Specifically, we assume the raw input $x \in \mathbb{R}^{d_x}$ has an intrinsic lower-dimensional representation.

Assumption 4.1 (Low-Dimensional Linear Latent Space). Initial data x has a latent representation via $x = Uh$, where $U \in \mathbb{R}^{d_x \times d_0}$ is an unknown matrix with orthonormal columns. The latent variable $h \in \mathbb{R}^{d_0}$ follows the distribution P_h with a density function p_h .

Remark 4.1. “Linear Latent Space” means that each entry of a given latent vector is a linear combination of the corresponding input, i.e., $x = Uh$. This is also known as the “low-dimensional data” assumption in literature ([Hu et al., 2024](#); [Chen et al., 2023c](#)). This assumption is fundamental in dimensionality reduction techniques for capturing the intrinsic lower-dimensional structure of data.

Score Decomposition and Model Architecture. To derive approximation and estimation results, we extend the techniques and network architecture presented in [Section 3](#) to latent diffusion by considering the “low-dimensional linear subspace”. Under [Assumption 4.1](#), we decompose the score:

$$\nabla \log p_t(x|y) = U \left(\underbrace{\sigma_t^2 \nabla \log p_t^h(U^\top x|y) + U^\top x}_{:=q(U^\top x, y, t): \mathbb{R}^{d_0} \times \mathbb{R}^{d_y} \times [t_0, T] \rightarrow \mathbb{R}^{d_0}} \right) / \sigma_t^2 - \underbrace{x / \sigma_t^2}_{\text{residual connection}}, \quad (4.1)$$

following [Hu et al. \(2024\)](#); [Chen et al. \(2023c\)](#) (see [Lemma E.1](#)). Based on this decomposition, we construct the model architecture in [Figure 3](#). The network detail for approximate (4.1) are as follow: a transformer $g_{\mathcal{T}}(W_U^\top x, y, t) \in \mathcal{T}_{\tilde{R}}^{h,s,r}$ to approximate $q(U^\top x, y, t)$, a latent encoder $W_U^\top \in \mathbb{R}^{d_0 \times d_x}$ and decoder $W_U \in \mathbb{R}^{d_x \times d_0}$ to approximate $U^\top \in \mathbb{R}^{d_0 \times d_x}$ and $U \in \mathbb{R}^{d_x \times d_0}$, and a residual connection to approximate $-x/\sigma_t^2$. Importantly, d_0 is the latent dimension.

For latent diffusion, we follow the standard setting by [Peebles and Xie \(2023\)](#). For each input $x \in \mathbb{R}^{d_x}$ and corresponding label $y \in \mathbb{R}^{d_y}$, we use a transformer network to obtain a score estimator $s_W \in \mathbb{R}^{d_x}$. The key differences from [Section 3](#) are as follows: First, we apply a latent encoder $W_U^\top \in \mathbb{R}^{d_0 \times d_x}$ to map the raw data $x \in \mathbb{R}^{d_x}$ into a low-dimensional representation $h := W_U^\top x \in \mathbb{R}^{d_0}$, where $d_0 \leq d_x$. Second, we reshape $h \in \mathbb{R}^{d_0}$ into a sequence $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$ using a layer $\tilde{R}(\cdot): \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}$, with $d_0 = \tilde{d} \cdot \tilde{L}$. Note that, by $d_0 \leq d_x$, $\tilde{d} \leq d$, and $\tilde{L} \leq L$. Third, we pass $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$ through the

transformer $g_{\mathcal{T}}$. Lastly, We then obtain the predicted score $s_W \in \mathbb{R}^{d_x}$ by applying the inverse reshape layer $\tilde{R}^{-1}(\cdot) : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{d_0}$, followed by the latent decoder $W_U : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_x}$.

For our analysis, we study the cases under both the generic and strong Hölder smoothness assumptions on latent representation $z \in \mathbb{R}^{d_0}$. Specifically, we assume the ‘‘latent’’ data is β_0 -Hölder smooth with radius B_0 following [Assumptions 3.1](#) and [3.2](#). We extend both approximation and estimation results from [Section 3](#) to latent diffusion and establish the minimax optimality of latent conditional DiTs.

Score Approximation. We now present the approximation rates for latent score function under both generic and stronger Hölder data assumptions. Let $h := W_U^\top x \in \mathbb{R}^{d_0}$ and $\bar{h} := U^\top x \in \mathbb{R}^{d_0}$ be the estimated and ground truth (according to [Assumption 4.1](#)) latent representations, respectively.

Theorem 4.1 (Score Approximation of Latent Conditional DiTs (Informal Version of [Theorems E.1](#) and [E.2](#))). *Assume $d_x = \Omega(\frac{\log N}{\log \log N})$. For any precision $0 < \epsilon < 1$ and smoothness $\beta_0 > 0$, let $\epsilon \leq \mathcal{O}(N^{-\beta_0})$ 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 $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_{\tilde{R}}^{h, s, r}$ such that*

- Under [Assumption 3.1](#), we have

$$\int_{\mathbb{R}^{d_0}} \|\mathcal{T}_{\text{score}}(\bar{h}, y, t) - \nabla \log p_t^h(\bar{h}|y)\|_2^2 \cdot p_t^h(\bar{h}|y) d\bar{h} = \mathcal{O}\left(\frac{B_0^2}{\sigma_t^4} \cdot N^{-\frac{\beta_0}{d_0+d_y}} \cdot (\log N)^{d_0+\frac{\beta_0}{2}+1}\right).$$

Notably, for $\epsilon = \mathcal{O}(N^{-\beta_0})$, the approximation error has the upper bound $\mathcal{O}((\log(\frac{1}{\epsilon}))^{d_0}/\sigma_t^4)$.

- Under [Assumption 3.2](#), we have

$$\int_{\mathbb{R}^{d_0}} \|\mathcal{T}_{\text{score}}(x, y, t)(\bar{h}, y, t) - \nabla \log p_t^h(\bar{h}|y)\|_2^2 \cdot p_t^h(\bar{h}|y) d\bar{h} = \mathcal{O}\left(\frac{B_0^2}{\sigma_t^2} \cdot N^{-\frac{2\beta_0}{d_0+d_y}} \cdot (\log N)^{\beta_0+1}\right).$$

Notably, for $\epsilon = \mathcal{O}(N^{-\beta_0})$, the approximation error has the upper bound $(\log(\frac{1}{\epsilon}))^{\mathcal{O}(1)}/\sigma_t^2$.

Proof. See [Theorems E.1](#) and [E.2](#) for the formal versions and [Appendices I](#) and [J](#) for proofs. \square

Remark 4.2 (Comparing with [Theorems 3.1](#) and [3.2](#)). Recall $d_x \geq d_0$, and the approximation error bounds are $\tilde{\mathcal{O}}(\epsilon^{1/(d_x+d_y)}/\sigma_t^2)$ in [Theorem 3.1](#) and $\tilde{\mathcal{O}}(\epsilon^{2/(d_x+d_y)}/\sigma_t^2)$ in [Theorem 3.2](#). These results show that the latent conditional DiT achieves better approximation and has the potential to bypass the challenges associated with the high dimensionality of initial data.

Score and Distribution Estimation. Based on [Theorem 4.1](#), we derive the score estimation bounds in [Theorem E.3](#), and report the results for distribution estimation in next theorem.

Theorem 4.2 (Distribution Estimation of Latent Conditional DiTs). *Assume $d_0 = \Omega(\frac{\log N}{\log \log N})$. For $y \in [0, 1]^{d_y}$, let $\hat{P}_{t_0}(\cdot|y)$ denote *estimated* conditional distributions at t_0 . Recall that $P_0(\cdot|y)$ is the conditional distribution of initial data x_0 given y . Assume $\text{KL}(P_0(\cdot|y) | N(0, I)) \leq c$ for some constant $c < \infty$.*

- Under [Assumption 3.1](#), taking $t_0 = n^{-\frac{\beta_0}{(d_0+d_y+\beta_0)}}$ and $T = \frac{2\beta_0}{d_0+d_y+2\beta_0} \log n$, it holds

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\hat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{\beta_0}{2(\tilde{\nu}_1-1)(d_0+d_y+\beta_0)}} (\log n)^{\frac{\tilde{\nu}_2}{2} + \frac{3}{2}} \right),$$

where $\tilde{\nu}_1 = \frac{68\beta_0}{(d_0+d_y)} + 104C_\sigma$, $\tilde{\nu}_2 = 12d_0 + 12\beta_0 + 2$ and $C_\sigma = \frac{\beta_0}{d_0+d_y+\beta_0}$.

- Under [Assumption 3.2](#), taking $t_0 = n^{-\frac{\beta_0}{4(d_0+d_y+\beta_0)}}$ and $T = \frac{2\beta_0}{d_0+d_y+2\beta_0} \log n$, it holds

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\hat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{1}{2\tilde{\nu}_3} \frac{\beta_0}{d_0+d_y+2\beta_0}} (\log n)^{\max(6, \frac{\beta_0}{2} + \frac{3}{2})} \right),$$

where $\tilde{\nu}_3 = \frac{4(12\beta_0 d_0 + 31\beta_0 \tilde{d} + 6\beta_0)}{\tilde{d}(d_0+d_y)} + \frac{12(12C_\alpha d_0 + 25C_\alpha \tilde{d} + 6C_\alpha)}{\tilde{d}} + 72C_\sigma$ and $C_\alpha = \frac{2\beta_0}{d_0+d_y+2\beta_0}$.

Proof. Please see [Appendix K.6](#) for a detailed proof. \square

Remark 4.3 (Minimax Optimal Estimation). Following the same idea in [Section 3.4](#), we show that the estimation error bound in [Theorem 4.2](#) is the optimal tight bound for the latent unconditional DiT. Specifically, by applying [Corollary 3.4.2](#) and substituting $p(x|y)$ and d_x by $p_t^h(\bar{h}|y)$ and d_0 respectively in [Assumption 3.2](#), we establish a distribution estimation lower bound of $\mathcal{O}(n^{-\beta_0/(d_0+2\beta_0)})$. Setting $2\tilde{\nu}_3 = 1$, we obtain the minimax optimality of latent unconditional DiT.

Concluding Remarks. We defer the discussion of our results and concluding remarks to [Appendix A](#). We extend our analysis to the setting of ([Hu et al., 2024](#)) and improve their results in [Appendix F](#). Importantly, our bounds avoid the gigantic $2^{(1/\epsilon)^{2L}}$ term reported by [Hu et al. \(2024\)](#).

REFERENCES

- 540
541
542 Silas Alberti, Niclas Dern, Laura Thesing, and Gitta Kutyniok. Sumformer: Universal approximation
543 for efficient transformers. In *Topological, Algebraic and Geometric Learning Workshops 2023*,
544 pages 72–86. PMLR, 2023.
- 545
546 Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth
547 words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on
548 computer vision and pattern recognition*, pages 22669–22679, 2023.
- 549
550 Georgios Batzolis, Jan Stanczuk, Carola-Bibiane Schönlieb, and Christian Etmann. Conditional
551 image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*, 2021.
- 552
553 Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly d -linear con-
554 vergence bounds for diffusion models via stochastic localization. In *The Twelfth International
555 Conference on Learning Representations*, 2024.
- 556
557 Clément L Canonne. A short note on an inequality between kl and tv. *arXiv preprint
558 arXiv:2202.07198*, 2022.
- 559
560 Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang. PLOT:
561 Prompt learning with optimal transport for vision-language models. In *The Eleventh International
562 Conference on Learning Representations (ICLR)*, 2023a.
- 563
564 Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling:
565 User-friendly bounds under minimal smoothness assumptions. In *International Conference on
566 Machine Learning*, pages 4735–4763. PMLR, 2023b.
- 567
568 Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and
569 distribution recovery of diffusion models on low-dimensional data. In *International Conference on
570 Machine Learning*, pages 4672–4712. PMLR, 2023c.
- 571
572 Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models: Ap-
573 plications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*,
574 2024a.
- 575
576 Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy
577 as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint
578 arXiv:2209.11215*, 2022.
- 579
580 Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability
581 flow ode is provably fast. *Advances in Neural Information Processing Systems*, 36, 2024b.
- 582
583 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances
584 in neural information processing systems*, 34:8780–8794, 2021.
- 585
586 Anh-Dung Dinh, Daochang Liu, and Chang Xu. Rethinking conditional diffusion sampling with
587 progressive guidance. *Advances in Neural Information Processing Systems*, 36, 2023.
- 588
589 Zehao Dou, Minshuo Chen, Mengdi Wang, and Zhuoran Yang. Theory of consistency diffusion
590 models: Distribution estimation meets fast sampling. In *Forty-first International Conference on
591 Machine Learning*, 2024a.
- 592
593 Zehao Dou, Subhodh Kotekal, Zehao Xu, and Harrison H Zhou. From optimal score matching to
594 optimal sampling. *arXiv preprint arXiv:2409.07032*, 2024b.
- 595
596 Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable
597 creation in self-attention mechanisms. In *International Conference on Machine Learning (ICML)*,
598 pages 5793–5831. PMLR, 2022.
- 599
600 Hengyu Fu, Zehao Dou, Jiawei Guo, Mengdi Wang, and Minshuo Chen. Diffusion transformer
601 captures spatial-temporal dependencies: A theory for gaussian process data. *arXiv preprint
602 arXiv:2407.16134*, 2024a.

- 594 Hengyu Fu, Zhuoran Yang, Mengdi Wang, and Minshuo Chen. Unveil conditional diffusion models
595 with classifier-free guidance: A sharp statistical theory. *arXiv preprint arXiv:2403.11968*, 2024b.
596
- 597 Khashayar Gatmiry, Jonathan Kelner, and Holden Lee. Learning mixtures of gaussians using diffusion
598 models. *arXiv preprint arXiv:2404.18869*, 2024.
599
- 600 Jiuxiang Gu, Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the frontiers
601 of softmax: Provable optimization, applications in diffusion model, and beyond. *arXiv preprint
602 arXiv:2405.03251*, 2024.
603
- 604 Hanzhong Guo, Cheng Lu, Fan Bao, Tianyu Pang, Shuicheng Yan, Chao Du, and Chongxuan Li.
605 Gaussian mixture solvers for diffusion models. *Advances in Neural Information Processing
606 Systems (NeurIPS)*, 37, 2023.
607
- 608 Yingqing Guo, Hui Yuan, Yukang Yang, Minshuo Chen, and Mengdi Wang. Gradient guidance for
609 diffusion models: An optimization perspective. *arXiv preprint arXiv:2404.14743*, 2024.
610
- 611 Alexander Havrilla and Wenjing Liao. Understanding scaling laws with statistical and approximation
612 theory for transformer neural networks on intrinsically low-dimensional data. In *The Thirty-eighth
613 Annual Conference on Neural Information Processing Systems*, 2024.
614
- 615 Satoshi Hayakawa and Taiji Suzuki. On the minimax optimality and superiority of deep neural
616 network learning over sparse parameter spaces. *Neural Networks*, 123:343–361, March 2020.
617 ISSN 0893-6080. doi: 10.1016/j.neunet.2019.12.014.
618
- 619 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
620 2022.
621
- 622 Jerry Yao-Chieh Hu, Weimin Wu, Zhuoru Li, Sophia Pi, , Zhao Song, and Han Liu. On statistical rates
623 and provably efficient criteria of latent diffusion transformers (dits). In *Thirty-eighth Conference
624 on Neural Information Processing Systems (NeurIPS)*, 2024.
625
- 626 Yuling Jiao, Lican Kang, Huazhen Lin, Jin Liu, and Heng Zuo. Latent schr $\{ \backslash " o \}$ dinger bridge
627 diffusion model for generative learning. *arXiv preprint arXiv:2404.13309*, 2024a.
628
- 629 Yuling Jiao, Yanming Lai, Yang Wang, and Bokai Yan. Convergence analysis of flow matching in
630 latent space with transformers. *arXiv preprint arXiv:2404.02538*, 2024b.
631
- 632 Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank
633 weight matrices universal approximators? In *The Twelfth International Conference on Learning
634 Representations (ICLR)*, 2024.
635
- 636 Junghwan Kim, Michelle Kim, and Barzan Mozafari. Provable memorization capacity of transformers.
637 In *The Eleventh International Conference on Learning Representations*, 2022.
638
- 639 Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with
640 polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882,
641 2022.
642
- 643 Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general
644 data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985.
645 PMLR, 2023.
646
- 647 Gen Li, Yu Huang, Timofey Efimov, Yuting Wei, Yuejie Chi, and Yuxin Chen. Accelerating
648 convergence of score-based diffusion models, provably. *arXiv preprint arXiv:2403.03852*, 2024a.
649
- 650 Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards non-asymptotic convergence for diffusion-
651 based generative models. In *The Twelfth International Conference on Learning Representations*,
652 2024b.
653
- 654 Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. A sharp convergence theory for the probability
655 flow odes of diffusion models. *arXiv preprint arXiv:2408.02320*, 2024c.
656
- 657 Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Unraveling the smoothness properties of
658 diffusion models: A gaussian mixture perspective. *arXiv preprint arXiv:2405.16418*, 2024a.

- 648 Yuchen Liang, Peizhong Ju, Yingbin Liang, and Ness Shroff. Non-asymptotic convergence of discrete-
649 time diffusion models: New approach and improved rate. *arXiv preprint arXiv:2402.13901*, 2024b.
650
- 651 Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang,
652 Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and
653 opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.
- 654 Haomiao Ni, Changhao Shi, Kai Li, Sharon X Huang, and Martin Renqiang Min. Conditional
655 image-to-video generation with latent flow diffusion models. In *Proceedings of the IEEE/CVF*
656 *conference on computer vision and pattern recognition*, pages 18444–18455, 2023.
657
- 658 Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution
659 estimators. In *International Conference on Machine Learning*, pages 26517–26582. PMLR, 2023.
660
- 661 Sejun Park, Jaeho Lee, Chulhee Yun, and Jinwoo Shin. Provable memorization via deep neural
662 networks using sub-linear parameters. In *Conference on Learning Theory (COLT)*, pages 3627–
663 3661. PMLR, 2021.
- 664 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
665 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4195–4205, 2023.
666
- 667 William S Peebles and Saining Xie. Scalable diffusion models with transformers. 2023 ieee. In *CVF*
668 *International Conference on Computer Vision (ICCV)*, volume 4172, 2022.
- 669 Peter Potaptchik, Iskander Azangulov, and George Deligiannidis. Linear convergence of diffusion
670 models under the manifold hypothesis. *arXiv preprint arXiv:2410.09046*, 2024.
671
- 672 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
673 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
674 *ence on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 675 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
676 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
677 text-to-image diffusion models with deep language understanding. *Advances in neural information*
678 *processing systems*, 35:36479–36494, 2022.
679
- 680 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation
681 function. *The Annals of Statistics*, 2020, 2020.
- 682 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
683 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
684 pages 2256–2265. PMLR, 2015.
685
- 686 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
687 Poole. Score-based generative modeling through stochastic differential equations. In *International*
688 *Conference on Learning Representations*, 2021.
- 689 Wenpin Tang and Hanyang Zhao. Score-based diffusion models via stochastic differential equations—a
690 technical tutorial. *arXiv preprint arXiv:2402.07487*, 2024.
691
- 692 Matus Telgarsky. Neural networks and rational functions. In *International Conference on Machine*
693 *Learning*, pages 3387–3393. PMLR, 2017.
694
- 695 Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space.
696 *Advances in neural information processing systems*, 34:11287–11302, 2021.
- 697 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computa-*
698 *tion*, 23(7):1661–1674, 2011.
699
- 700 Vikram Voleti, Alexia Jolicoeur-Martineau, and Chris Pal. Mcvd-masked conditional video diffusion
701 for prediction, generation, and interpolation. *Advances in neural information processing systems*,
35:23371–23385, 2022.

702 Junde Wu, Wei Ji, Huazhu Fu, Min Xu, Yueming Jin, and Yanwu Xu. Medsegdiff-v2: Diffusion-based
703 medical image segmentation with transformer. In *Proceedings of the AAAI Conference on Artificial*
704 *Intelligence*, volume 38, pages 6030–6038, 2024a.

705
706 Yuchen Wu, Minshuo Chen, Zihao Li, Mengdi Wang, and Yuting Wei. Theoretical insights for
707 diffusion guidance: A case study for gaussian mixture models. In *Forty-first International*
708 *Conference on Machine Learning*, 2024b.

709 Yuchen Wu, Yuxin Chen, and Yuting Wei. Stochastic runge-kutta methods: Provable acceleration of
710 diffusion models. *arXiv preprint arXiv:2410.04760*, 2024c.

711
712 Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Zou,
713 and Stefano Ermon. Tfg: Unified training-free guidance for diffusion models. *arXiv preprint*
714 *arXiv:2409.15761*, 2024.

715 Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed
716 conditional diffusion: Provable distribution estimation and reward improvement. *Advances in*
717 *Neural Information Processing Systems*, 36, 2023.

718 Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed
719 conditional diffusion: Provable distribution estimation and reward improvement. *Advances in*
720 *Neural Information Processing Systems*, 36, 2024.

721
722 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are trans-
723 formers universal approximators of sequence-to-sequence functions? In *International Conference*
724 *on Learning Representations (ICLR)*, 2020.

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

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Appendix

A Discussion and Conclusion	16
B Notation Table	17
C Related Works, Broader Impact and Limitations	18
C.1 Related Works	18
C.2 Broader Impact	19
C.3 Limitations	19
D Proof-of-Concept Experiments	20
D.1 Experimental Results	20
E Latent Conditional DiT with Hölder Assumption	22
E.1 Score Approximation	24
E.2 Score Estimation	26
E.3 Distribution Estimation	27
F Latent Conditional DiT with Lipschitz Assumption	28
F.1 Score Approximation	30
F.2 Score Estimation	31
F.3 Distribution Estimation	31
F.4 Proof of Score Approximation (Theorem F.1)	32
F.5 Proof of Score Estimation (Theorem F.2)	38
F.6 Proof of Distribution Estimation (Theorem F.3)	43
G Supplementary Theoretical Background	47
G.1 Conditional Diffusion Process	47
G.2 Classifier-free Guidance	48
H Universal Approximation of Transformers	49
H.1 Transformers as Universal Approximators	49
H.2 Parameter Norm Bounds for Transformer Approximation	58
I Proof of Theorem 3.1	62
I.1 Auxiliary Lemmas	62
I.2 Main Proof of Theorem 3.1	83
J Proof of Theorem 3.2	88
J.1 Auxiliary Lemmas	88
J.2 Main Proof of Theorem 3.2	100
K Proof of the Estimation Results for Conditional DiTs	102
K.1 Auxiliary Lemmas for Theorem 3.3	102
K.2 Proof of Theorem 3.3	111
K.3 Dominance Transition between N and $\log N$ for All Norm Bounds under Assumption 3.1	114
K.4 Proof of Corollary 3.3.1	119
K.5 Auxiliary Lemmas for Theorem 3.4	123
K.6 Main Proof of Theorem 3.4	123
K.7 Proof of Corollary 3.4.1	125

A DISCUSSION AND CONCLUSION

We investigate the approximation and estimation rates of conditional DiT and its latent setting. We focus on the “in-context” conditional DiT setting presented by Peebles and Xie (2023), and conduct a comprehensive analysis under various common data conditions (Section 3 for generic and strong Hölder smooth data, Section 4 for data with intrinsic latent subspace).

Interestingly, we establish the minimax optimality of the unconditional DiTs’ estimation by reducing our analysis of conditional DiTs to the unconditional setting (Section 3.4 and Remark 4.3). Our key techniques include a well-designed score decomposition scheme (Section 3.1). These enable a finer use of transformers’ universal approximation, compared to the prior statistical rates of DiTs derived from the universal approximation results in (Yun et al., 2020) by Hu et al. (2024).

Consequently, we provide two extensions in the appendix:

- In Appendix E, we expand Section 4 and extend our well-designed score decomposition scheme from Section 3 to the latent conditional DiT. Notably, we also obtain provably tight rate, i.e., for distribution estimation under Assumption 3.2 (Remark 4.3).
- In Appendix F, we extend the analysis of (Hu et al., 2024) to the conditional DiT setting and provide an improved version. In particular, we analyze conditional latent DiTs under the following three assumptions from (Hu et al., 2024) and obtained sharper rates:
 - Low-Dimensional Linear Latent Space Data (Assumption 4.1)
 - Lipschitz Score Function (Assumption F.2)
 - Light Tail Data Distribution (Assumption F.3)

In detail, we use a modified universal approximation of the single-layer self-attention transformers (modified from (Kajitsuka and Sato, 2024)) to avoid the need for dense layers required in (Yun et al., 2020). This refinement results in tighter error bounds for both score and distribution estimation. Consequently, our sample complexity error bounds avoid the gigantic double exponential term $2^{(1/\epsilon)^{2L}}$ reported by Hu et al. (2024), and obtain sharper rates than those of (Hu et al., 2024).

B NOTATION TABLE

We summarize our notations in the following table for easy reference.

Table 2: Mathematical Notations and Symbols

Symbol	Description
$[I]$	The index set $\{1, \dots, I\}$, where $I \in \mathbb{N}^+$
$a[i]$	The i -th component of vector a
A_{ij}	The (i, j) -th entry of matrix A
$\ x\ $	Euclidean norm of vector x
$\ x\ _1$	1-norm of vector x
$\ x\ _2$	2-norm of vector x
$\ x\ _\infty$	Infinite norm of vector x
$\ W\ _2$	Spectral norm of matrix W
$\ W\ _F$	Frobenius norm of matrix W
$\ W\ _{p,q}$	(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
$d_p(f, g)$	p -norm of the difference between functions f and g defined as $d_p(f, g) = (\int f(x) - g(x) ^p dx)^{1/p}$
$f_\# P$	Pushforward measure, where f is a function and P is a distribution
$\text{KL}(P, Q)$	Kullback-Leibler (KL) divergence between distributions P and Q
$\text{TV}(P, Q)$	Total variation (TV) distance between distributions P and Q
$N(\mu, \sigma^2)$	Normal distribution with mean μ and variance σ^2
$a \lesssim b$	There exist constants $C > 0$ such that $a \leq Cb$
n	Sample size
x	Data point in original data space, $x \in \mathbb{R}^{d_x}$
y	Conditioning Label, $x \in \mathbb{R}^{d_y}$
h	Latent variable in low-dimensional subspace, $h \in \mathbb{R}^{d_0}$
\bar{h}	$\bar{h} = U^\top x$
p_h	The density function of h
U	The matrix with orthonormal columns to transform h to x , where $U \in \mathbb{R}^{d \times d_0}$
B	Radius of Hölder ball for conditional density function $p(x y)$
B_0	Radius of Hölder ball for latent conditional density function $p(\bar{h} y)$
β	Hölder index for conditional density function $p(x y)$
β_0	Hölder index for latent conditional density function $p(\bar{h} y)$
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 ℓ^2 difference between score estimator and ground truth)
$\mathcal{N}(\epsilon, \mathcal{F}, \ \cdot\)$	Covering number of collection \mathcal{F} (see Definition K.5)
T	Stopping time in the forward process of diffusion model
t_0	Stopping time in the backward process of diffusion model
μ	Discretized step size in backward process
$p_t(\cdot)$	The density function of x at time t
$p_t^h(\cdot)$	The density function of \bar{h} at time t
ψ	(Conditional) Gaussian density function
$\mathcal{T}^{h,s,r}$	Transformer network function class (see Definition 2.2)
$f^{h,s,r}$	Transformer block of h -head, s -hidden size, r -MLP dimension (see Definition 2.1)
d	Input dimension of each token in the transformer network of DiT
L	Token length in the transformer network of DiT
\tilde{d}	Latent data input dimension of each token in the transformer network of DiT
\tilde{L}	Latent data token length in the transformer network of DiT
X	Sequence input of transformer network in DiT, where $X \in \mathbb{R}^{d \times L}$
H	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, $R(\cdot) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d \times L}$
$\tilde{R}(\cdot)$	Reshape layer in DiT, $\tilde{R}(\cdot) : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}$
$R^{-1}(\cdot)$	Reverse reshape layer in DiT, $R^{-1}(\cdot) : \mathbb{R}^{d \times L} \rightarrow \mathbb{R}^{d_x}$
$\tilde{R}^{-1}(\cdot)$	Reverse reshape layer in DiT, $\tilde{R}^{-1}(\cdot) : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{d_0}$
W_U	The orthonormal matrix to approximate U , where $W_U \in \mathbb{R}^{d_x \times d_0}$

C RELATED WORKS, BROADER IMPACT AND LIMITATIONS

C.1 RELATED WORKS

In the following, we discuss the recent success of the techniques used in our work. We first give the universality (universal approximation) of the transformer. Then, we discuss recent theoretical developments (approximation and estimation) in diffusion generative models.

Universality of Transformers. The universality of transformers refers to their ability to approximate any sequence-to-sequence function with arbitrary precision. Yun et al. (2020) establish this by showing that transformers is capable of universally approximate sequence-to-sequence functions using deep stacks of feed-forward and self-attention layers. Additionally, Alberti et al. (2023) demonstrate universal approximation for architectures employing non-standard attention mechanisms. Recently, Kajitsuka and Sato (2024) show that even a single-layer transformer with self-attention suffices for universal approximation assuming all attention weights are rank-1. Moreover, Hu et al. (2024) leverage Yun et al. (2020) universality results to analyze the approximation and estimation capabilities of DiT.

Our paper is motivated by and builds upon the works of Hu et al. (2024); Kajitsuka and Sato (2024); Yun et al. (2020). Specifically, we utilize and extend the transformer universality result from Kajitsuka and Sato (2024). We employ a relaxed contextual mapping property in Kajitsuka and Sato (2024) (see Appendix H.1). This generalization allows us to avoid the “double exponential” sample complexity bounds in previous DiT analyses (Hu et al., 2024, Remark 3.4) and establish transformer approximation in the simplest configuration — a single-layer, single-head attention model.

Approximation and Estimation Theories of Diffusion Models. The theories of DiTs revolve around two main frontiers: score function approximation and statistical estimation (Chen et al., 2024a; Tang and Zhao, 2024). Score function approximation refers to the ability of the score network to approximate the score function. It leverages the universal approximation ability of the neural network in L^p norms (Hayakawa and Suzuki, 2020), the approximation characterized as Taylor polynomial (Fu et al., 2024a) or B-Spline (Oko et al., 2023). Chen et al. (2023c) and Fu et al. (2024a) investigate score approximation under specific conditions, such as low-dimensional linear subspaces and Hölder smooth data assumptions, using ReLU-based models. Furthermore, Hu et al. (2024) presents the first characterization of score approximation in diffusion transformers (DiTs).

The statistical estimation includes score function and distribution estimation (Wu et al., 2024b; Dou et al., 2024a; Guo et al., 2024; Chen et al., 2023c). Under a L_2 accurate score estimation, several works have provided the convergence bounds under either smoothness assumptions (Benton et al., 2024; Chen et al., 2022) or bounded second-order moment assumptions (Chen et al., 2023b; Lee et al., 2023). Chen et al. (2023c) provide the first complete estimation theory using ReLU networks without precise estimators. Oko et al. (2023) achieve nearly minimax optimal estimation rates for total variation and Wasserstein distances. Meanwhile, Dou et al. (2024b) define exact minimax optimality using kernel functions without characterizing the network architectures. In the realm of diffusion transformers, Hu et al. (2024) introduces the first complete estimation theory. Jiao et al. (2024a;b) demonstrate theoretical convergence for latent DiTs using ODE-based and Schrödinger bridge diffusion models.³

Our paper advances the foundational works of Fu et al. (2024b); Oko et al. (2023); Hu et al. (2024). We adopt the Hölder smooth data distribution assumption⁴, a more practical approach than the bounded support assumption in Oko et al. (2023). Unlike the simple ReLU networks in Fu et al. (2024b), we provide a complete approximation and estimation analysis for conditional DiTs and establish their exact minimax optimality. Furthermore, while Hu et al. (2024) analyze DiTs, their estimation upper bounds are suboptimal. We refine this by avoiding the substantial double exponential

³Of independent interest, many works investigate the convergence rates of diffusion models under various score and data smoothness assumptions or with different samplers. Please see (Li et al., 2024a;b;c; Potapchik et al., 2024; Wu et al., 2024c; Liang et al., 2024b;a; Garmiry et al., 2024; Gu et al., 2024; Guo et al., 2023; Chen et al., 2024b; 2023b; 2022; Lee et al., 2023; 2022) and references therein.

⁴Recent work by Havrilla and Liao (2024) examines the generalization and approximation of transformers under Hölder smoothness and low-dimensional subspace assumptions.

972 term $2^{(1/\epsilon)^{2L}}$ reported by [Hu et al. \(2024, Remark 3.4\)](#) and present a provably tight, minimax optimal
973 estimation.

974

975 C.2 BROADER IMPACT

976

977 This theoretical work aims to shed light on the foundations of generative diffusion models and is not
978 expected to have negative social impacts.

979

980 C.3 LIMITATIONS

981

982 Although our study provides a complete theoretical analysis of the conditional DiTs and establishes
983 the minimax optimality of the unconditional DiT, we acknowledge three main limitations:

984

- 985 • The minimax optimality of conditional DiT remains not clear.
- 986 • We did not explore other architectures such as “adaptive layer norm” and “cross-attention” DiT. A
987 potential direction is by establishing the universal approximation capacity of the transformer with
988 cross-attention mechanisms.
- 989 • Although we achieve a better bound for the latent conditional DiT under the Lipschitz assumption
990 than under the Hölder assumption, we do not show the minimax optimality under the Lipschitz
991 assumption.

992

993 We leave these for future work.

994

995 Furthermore, there are limitations regarding the Hölder smooth data assumptions in [Assumption 3.1](#)
996 and [Assumption 3.2](#). Our results in [Section 3](#) and [Section 4](#) depend on the Hölder smooth data
997 assumptions. However, it is challenging to measure the smoothness of a given dataset (e.g., CIFAR10),
998 because it requires knowledge of the dataset’s exact distribution. Conversely, it is feasible to create a
999 dataset with a predefined level of smoothness. To illustrate this, we provide two examples.

1000

- 1001 • **Diffusion Models in Image Generation:** When modeling conditional distributions of images given
1002 attributes (e.g., generating images based on class labels), these assumptions hold if the data
1003 distribution around these attributes is smooth and decays. In diffusion-based generative models,
1004 the data distribution often decays smoothly in high-dimensional space. The assumption that the
1005 density function decays exponentially reflects the natural behavior of image data, where pixels or
1006 features far from a central region or manifold are less likely. This is commonly observed in images
1007 with blank boundaries.

1008

- 1009 • **Physical Systems with Gaussian-Like Decay:** This applies to cases where the spatial distribution
1010 of a physical quantity, such as temperature, is smooth and governed by diffusion equations with
1011 exponential decay. In physics-based diffusion models, like those simulating the spread of particles
1012 or heat in a medium (e.g., stars in galaxies for astrophysics applications), the conditional density
1013 typically decays exponentially with distance from a central region.

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

D PROOF-OF-CONCEPT EXPERIMENTS

Experimental Objectives. We train a conditional diffusion transformer model on the CIFAR10 dataset to validate the following three parts:

- **Objective 1.** Validating the influence of input data dimension d_x on the testing loss (score estimation error) in [Theorem 3.3](#).
- **Objective 2.** Validating the influence of input data dimension d_x on the parameter norm bounds ($\|W_O\|_{2,\infty}$ and $\|W_V\|_{2,\infty}$) in [Theorem 3.1](#).
- **Objective 3.** Validating the influence of backward timestamp t_0 on the testing loss (score estimation error) in [Theorem 3.3](#).

Experimental Details. We train the model on the CIFAR10 training dataset for 10 epochs. The dataset consists of 50,000 images across 10 classes. We set the forward process termination step to $T = 1000$. Then, we evaluate the model’s performance using the CIFAR10 testing dataset of 10,000 images from 10 classes. We use the testing loss as the measurement.

- To validate objectives 1 and 2, we test various values of d_x at backward timestamp $t_0 = 5$, including $32 \cdot 32 = 1,024$, $48 \cdot 48 = 2,304$, $64 \cdot 64 = 4,096$, and $80 \cdot 80 = 6,400$.
- To validate objective 3, we test different backward timestamps t_0 , including 5, 4, 3, 2 and 1 for both $d_x = 32 \cdot 32 = 1,024$ and $d_x = 48 \cdot 48 = 2,304$.

Model Setup. The conditional diffusion transformer model has 12 transformer blocks. The number of attention heads is $h = 6$, and the hidden dimension is $s = 384$. We set the MLP dimension to $r = 1536$. We fix $d = 4$ in the DiT reshape layer ([Definition 2.3](#)).

Computational Resource. We conduct all experiments using 1 NVIDIA A100 GPU with 80GB of memory. Our code is based on the PyTorch implementation of the diffusion transformer ([Peebles and Xie, 2023](#)) at <https://github.com/chuanyangjin/fast-DiT>.

D.1 EXPERIMENTAL RESULTS

Results for Objectives 1 and 2. We report the numerical results of objectives 1 and 2 in [Table 3](#).

We observe an increase in the loss value with increasing d_x . This is consistent with the score estimation result in [Theorem 3.3](#).

Additionally, we note an increase in the parameter norm bounds ($\|W_O\|_{2,\infty}$ and $\|W_V\|_{2,\infty}$) with increasing d_x . These align with the parameter norm bound results in [Theorem 3.1](#).

Table 3: **Influence of Input Data Dimension d_x on the Testing Loss and Parameter Norm Bounds at Backward Timestamp $t_0 = 5$:** The testing loss and parameter norm bounds ($\|W_O\|_{2,\infty}$ and $\|W_V\|_{2,\infty}$) increase with an increasing d_x . These results are consistent with the results in [Theorem 3.3](#) and [Theorem 3.1](#).

Input Data Dim. d_x	$32 \cdot 32 = 1,024$	$48 \cdot 48 = 2,304$	$64 \cdot 64 = 4,096$	$80 \cdot 80 = 6,400$
Testing loss	0.9321	0.9356	0.9364	0.9476
$\ W_O\ _{2,\infty}$	1.6074	1.6332	1.6789	1.6886
$\ W_V\ _{2,\infty}$	2.1513	2.1767	2.1858	2.1994

Results for Objective 3. We report numerical results of objectives 3 for $d_x = 32 \cdot 32 = 1,024$ and $d_x = 48 \cdot 48 = 2,304$ in [Table 4](#). We observe an increase in the loss value as t_0 decreases. This is consistent with the score estimation result in [Theorem 3.3](#).

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Table 4: **Influence of Backward Timestamp t_0 on the Testing Loss:** The testing loss increases with increasing t_0 . This is consistent with the result in [Theorem 3.3](#).

Testing loss	$t_0 = 5$	$t_0 = 4$	$t_0 = 3$	$t_0 = 2$	$t_0 = 1$
$32 \cdot 32 = 1,024$	0.9321	0.9329	0.9335	0.9350	0.9361
$48 \cdot 48 = 2,304$	0.9356	0.9357	0.9360	0.9363	0.9367

E LATENT CONDITIONAL DiT WITH HÖLDER ASSUMPTION

In this section, we extend the results on approximation and estimation of DiT from Section 3 by considering the latent conditional DiTs. Latent DiTs enables efficient data generation from latent space and therefore scales better in terms of spatial dimensionality (Rombach et al., 2022). Specifically, we assume the raw input $x \in \mathbb{R}^{d_x}$ has an intrinsic lower-dimensional representation in a d_0 -dimensional subspace, where $d_0 \leq d_x$. This setting is common in both empirical (Peebles and Xie, 2022; Rombach et al., 2022) and theoretical studies (Hu et al., 2024; Chen et al., 2023c).

Organization. We present the statistical results under Hölder data smooth Assumptions 3.1 and 3.2 and state the results in Theorem E.1, Theorem E.2, Theorem E.3, and Theorem E.4, respectively. Appendix E.1 discusses score approximation. Appendix E.2 discusses score estimation. Appendix E.3 discusses distribution estimation. The proofs in this section primarily follow Appendices I and J.

Let d_0 denote the latent dimension. We summarize the key points of this section as follows:

K1. Low-Dimensional Subspace Space Data Assumption. We consider the setting that latent representation lives in a “Low-Dimensional Subspace” under Assumption 4.1, following (Hu et al., 2024; Chen et al., 2023c).

Assumption E.1 (Low-Dimensional Linear Latent Space (Assumption 4.1 Restated)). Data point $x = Uh$, where $U \in \mathbb{R}^{d_x \times d_0}$ is an unknown matrix with orthonormal columns. The latent variable $h \in \mathbb{R}^{d_0}$ follows a distribution P_h with a density function p_h .

For raw data $x \in \mathbb{R}^{d_x}$, we utilize linear encoder $W_U^\top \in \mathbb{R}^{d_0 \times d_x}$ and decoder $W_U \in \mathbb{R}^{d_x \times d_0}$ to convert the raw $x \in \mathbb{R}^{d_x}$ and latent $h \in \mathbb{R}^{d_0}$ data representations. Importantly, $x = Uh$ with $U \in \mathbb{R}^{d_x \times d_0}$ by Assumption 4.1.

For each input $x \in \mathbb{R}^{d_x}$ and corresponding label $y \in \mathbb{R}^{d_y}$, we use a transformer network to obtain a score estimator $s_W \in \mathbb{R}^{d_x}$. To utilize the transformer network as the score estimator, we introduce reshape layer to convert vector input $h \in \mathbb{R}^{d_0}$ to matrix (sequence) input $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$. Specifically, the reshape layer in the network Figure 3 is defined as $\tilde{R}(\cdot) : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}$ and its reverse $\tilde{R}^{-1}(\cdot) : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{d_0}$, where $d_0 \leq d_x$, $\tilde{d} \leq d$, and $\tilde{L} \leq L$.

We remark that the “low-dimensional data” assumption leads to tighter approximation rates than those of Sections 3.1 and 3.2 and estimation errors due to $d_0 \leq d_x$ (Theorems E.1 and E.2).

K2. Hölder Smooth Assumption. For approximation and estimation results for latent conditional DiTs (Theorems E.1 to E.4), we study the cases under both the generic and strong Hölder smoothness assumptions on latent representation $h \in \mathbb{R}^{d_0}$. Specifically, we assume the “latent” data is β_0 -Hölder smooth with radius B_0 following Assumptions 3.1 and 3.2. We extend both approximation and estimation results from Section 3 to latent diffusion and establish the minimax optimality of latent conditional DiTs.

Assumption E.2 (Generic Hölder Smooth Data (Assumption 3.1 Restated)). The conditional density function $p_0^h(h_0|y)$ is defined on the domain $\mathbb{R}^{d_0} \times [0, 1]^{d_y}$ and belongs to Hölder ball of radius $B_0 > 0$ for Hölder index $\beta_0 > 0$, denoted by $p_0^h(h_0|y) \in \mathcal{H}^{\beta_0}(\mathbb{R}^{d_0} \times [0, 1]^{d_y}, B_0)$ (see Definition 3.1 for precise definition.) Also, for any $y \in [0, 1]^{d_y}$, there exist positive constants C_1, C_2 such that $p_0^h(h_0|y) \leq C_1 \exp(-C_2 \|h_0\|_2^2/2)$.

Assumption E.3 (Stronger Hölder Smooth Data (Assumption 3.2 Restated)). Let function $f \in \mathcal{H}^{\beta_0}(\mathbb{R}^{d_0} \times [0, 1]^{d_y}, B_0)$. Given a constant radius B_0 , positive constants C and C_2 , we assume the conditional density function $p(h_0|y) = \exp(-C_2 \|h_0\|_2^2/2) \cdot f(h_0, y)$ and $f(h_0, y) \geq C$ for all $(h_0, y) \in \mathbb{R}^{d_0} \times [0, 1]^{d_y}$.

K3. Latent Score Network. Under low-dimensional data assumption, we decompose the score function following (Hu et al., 2024; Chen et al., 2023c) (see Lemma E.1):

$$\nabla \log p_t(x|y) = \underbrace{U(\sigma_t^2 \nabla \log p_t^h(U^\top x|y) + U^\top x)/\sigma_t^2}_{:=q(U^\top x, y, t): \mathbb{R}^{d_0} \times [t_0, T] \rightarrow \mathbb{R}^{d_0}} - \underbrace{x/\sigma_t^2}_{\text{residual connection}}. \quad (\text{E.1})$$

Based on this decomposition, we construct the model architecture in [Figure 3](#). The network detail for approximate [\(E.1\)](#) are as follow: a transformer $g_{\mathcal{T}}(W_U^\top x, y, t) \in \mathcal{T}^{h,s,r}$ to approximate $q(U^\top x, y, t)$, a latent encoder $W_U^\top \in \mathbb{R}^{d_0 \times d_x}$ and decoder $W_U \in \mathbb{R}^{d_x \times d_0}$ to approximate $U^\top \in \mathbb{R}^{d_0 \times d_x}$ and $U \in \mathbb{R}^{d_x \times d_0}$, and a residual connection to approximate $-x/\sigma_t^2$.

We adopt the following transformer network class of one-layer single-head self-attention

$$\mathcal{T}_{\tilde{R}}^{h,s,r} = \left\{ s_W(x, y, t) = \frac{1}{\sigma_t^2} W_U g_{\mathcal{T}}(W_U^\top x, y, t) - \underbrace{\frac{1}{\sigma_t^2} x}_{\text{residual connection}} \right\}, \quad (\text{E.2})$$

where $g_{\mathcal{T}} \in \mathcal{T}^{h,s,r} = \{f_2^{\text{FF}} \circ f^{h,s,r} : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}\}$.

Let $h := W_U^\top x \in \mathbb{R}^{d_0}$ and $\bar{h} := U^\top x \in \mathbb{R}^{d_0}$ be the estimated and ground truth (according to [Assumption 4.1](#)) latent representations, respectively. Here we construct a network $s_W(x, y, t)$ to approximate the score function in [\(E.1\)](#) (see [Figure 3](#) for network illustration).

In [Section 3](#), we derive the approximation theory of conditional DiTs using a one-layer, single-head self-attention transformer to approximate the score function $\nabla \log p_t(x|y)$. Here, we use the similar transformer architecture to approximate latent score function $\nabla \log p_t^h(\bar{h}|y)$, where $p_t^h(\bar{h}|y) = \int \psi_t(\bar{h}|h) p_h(h|y) dh$, $\psi_t(\cdot|h)$ is the Gaussian density function of $N(\beta_t h, \sigma_t^2 I_{d_0})$, $\beta_t = e^{-t/2}$, and $\sigma_t^2 = 1 - e^{-t}$.

Base on the latent network construction in [\(K3\)](#), we employ the same techniques presented in [Section 3](#) for score function approximation and estimation. We restate for completeness. First, we decompose the conditional score function $\nabla \log p_t^h(\bar{h}|y)$ as following:

$$\nabla \log p_t^h(\bar{h}|y) = \frac{\nabla p_t^h(\bar{h}|y)}{p_t^h(\bar{h}|y)}. \quad (\text{E.3})$$

By the definition of Gaussian kernel, we have

$$p_t^h(\bar{h}|y) = \int_{\mathbb{R}^{d_0}} (2\pi\sigma_t^2)^{-d_x/2} \underbrace{p_h(h|y)}_{\approx k_1\text{-order Taylor polynomial}} \underbrace{\exp\left(-\frac{\|\beta_t h - \bar{h}\|_2^2}{2\sigma_t^2}\right)}_{\approx k_2\text{-order Taylor polynomial}} dh.$$

Similar to [Section 3](#), our strategy is to expand above term-by-term with k_1 - and k_2 -order Taylor polynomials for fine-grained characterizations.

Remark E.1. Here in the latent density function, we have $(2\pi\sigma_t^2)^{-d_x/2}$ instead of $(2\pi\sigma_t^2)^{-d_0/2}$. However, the additional $(2\pi\sigma_t^2)^{-(d_x-d_0)/2}$ term does not affect the application of [Section 3](#) into latent diffusion approximation.

Based on the low-dimensional data structure assumption, we have the following score decomposition terms: on-support score $s_+(U^\top x, y, t)$ and orthogonal score $s_-(x, y, t)$.

Lemma E.1 (Score Decomposition, Lemma 1 of [\(Chen et al., 2023c\)](#)). Let data $x = Uh$ follow [Assumption 4.1](#). The decomposition of score function $\nabla \log p_t(x)$ is

$$\nabla \log p_t(x) = \underbrace{U \nabla \log p_t^h(\bar{h}|y)}_{s_+(\bar{h}, y, t)} - \underbrace{(I_D - UU^\top) x / \sigma_t^2}_{s_-(x, t)}, \quad \bar{h} = U^\top x, \quad (\text{E.4})$$

where $p_t^h(\bar{h}|y) := \int \psi_t(\bar{h}|h) p_h(h|y) dh$, $\psi_t(\cdot|h)$ is the Gaussian density function of $N(\beta_t h, \sigma_t^2 I_{d_0})$, $\beta_t = e^{-t/2}$ and $\sigma_t^2 = 1 - e^{-t}$.

Following the proof strategy of conditional DiTs in [Appendices I and J](#) with differences highlighted in [\(K1\)](#), [\(K2\)](#), and the latent network in [\(K3\)](#). To derive the approximation and estimation under generic

and stronger Hölder assumptions results in [Theorems 3.1 to 3.4](#) for data under low-dimensional data assumption, we just need to replace the input dimension d, L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 , and consider the β_0 -Hölder smoothness assumption on latent data.

To begin, we clarify the relation between initial data admits to $p(x|y) \in \mathcal{H}^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$, and under linear transformed data [Assumption 4.1](#) admits to $p(\bar{h}|y) \in \mathcal{H}^{\beta_0}(\mathbb{R}^{d_0} \times [0, 1]^{d_y}, B_0)$ where $\beta_0 = \beta$ and $B_0 \leq \tilde{C}B$ by [Lemma E.2](#).

Lemma E.2 (Transformation of Stronger Hölder Smooth Data Distribution under Linear Mapping). Let $f \in H^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$ satisfy $f(x, y) \geq C > 0$ for all $(x, y) \in \mathbb{R}^{d_x} \times [0, 1]^{d_y}$. Consider the conditional density function:

$$p(x|y) = f(x, y) \exp\left(-\frac{C_2}{2}\|x\|_2^2\right).$$

Suppose the data undergo the linear transformation $x = Uh$, where $U \in \mathbb{R}^{d_x \times d_0}$ has orthonormal columns ($U^\top U = I_{d_0}$) and $f_0(h|y) = f(Uh|y)$. The transformed density $p(h|y)$ becomes:

$$p(h|y) = f(Uh, y) \exp\left(-\frac{C_2}{2}\|h\|_2^2\right).$$

The following condition holds for Hölder smooth data undergo linear transformation: $f_0 \in H^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B_0)$ with $B_0 \leq \tilde{C}B$, where $\tilde{C} = \max\{C', C''\}$.

Proof. First, we compute the partial derivative of the transformed function $f_0(h|y) := f(Uh|y)$. From the definition of Hölder space [Definition 3.1](#), and let $\alpha = (\alpha_h, \alpha_y)$ where $\alpha_h + \alpha_y \leq k_1$. We compute the partial derivative up to the order of k_1 and show that it is bounded by some C' , that is

$$\begin{aligned} \partial_h^{\alpha_h} \partial_y^{\alpha_y} p(h|y) &= \partial_h^{\alpha_h} \partial_y^{\alpha_y} \left[f(Uh, y) \exp\left(-\frac{C_2}{2}\|h\|_2^2\right) \right] \\ &= \sum_{\alpha \leq \nu} \binom{\alpha}{\mu} (\partial_h^{\alpha_h} f(Uh, y)) \left(\partial_h^{\alpha_h - \mu} \exp\left(-\frac{C_2}{2}\|h\|_2^2\right) \right). \quad (\text{By product rule}) \end{aligned}$$

From the relation $\partial_h^{\alpha_h} f(Uh, y) = U^{\alpha_h} \partial_x^{\alpha_h} f(Uh, y)$ where U^{α_h} is the product of U entries correspond to α_h . Therefore, $\|\partial_h^{\alpha_h} \partial_y^{\alpha_y} f_0(h|y)\| \leq C'B$ for some C' depends on U and α_h . Since f satisfied Hölder condition and the mapping $h \mapsto Uh$ is linear, for Hölder condition $|\alpha_h| + |\alpha_y| = k_1$ there exist C'' such that

$$\frac{|\partial_h^{\alpha_h} \partial_y^{\alpha_y} f_0(h|y) - \partial_h^{\alpha_h} \partial_y^{\alpha_y} f_0(h'|y')|}{\|(h, y) - (h', y')\|_\infty^{\beta_0}} \leq C''B.$$

The bounded partial derivate up to order k_1 satisfied Hölder condition.

This completes the proof. \square

E.1 SCORE APPROXIMATION

We present the approximation rate of latent score function under generic Hölder and stronger Hölder data assumption in [Theorems E.1 and E.2](#), respectively.

Theorem E.1 (Latent Conditional DiT Score Approximation (Formal Version of [Theorem 4.1](#))). Assume [Assumption 3.1](#) and Assume $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$. For any precision $0 < \epsilon < 1$ and smoothness $\beta_0 > 0$, let $\epsilon \leq \mathcal{O}(N^{-\beta_0})$ for some $N \in \mathbb{N}$. For some positive constants $C_\alpha, C_\sigma > 0$, for any

1296 $y \in [0, 1]^{d_y}$ and $t \in [N^{-C_\sigma}, C_\alpha \log N]$, there exists a $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_{\tilde{R}}^{h,s,r}$ such that

$$1299 \int_{\mathbb{R}^{d_0}} \|\mathcal{T}_{\text{score}}(\bar{h}, y, t) - \nabla \log p_t^h(\bar{h}|y)\|_2^2 \cdot p_t^h(\bar{h}|y) d\bar{h} = \mathcal{O}\left(\frac{B_0^2}{\sigma_t^4} \cdot N^{-\frac{\beta_0}{d_0+d_y}} \cdot (\log N)^{d_0+\frac{\beta_0}{2}+1}\right).$$

1302 Notably, for $\epsilon = \mathcal{O}(N^{-\beta_0})$, the approximation error has the upper bound $\mathcal{O}((\log(\frac{1}{\epsilon}))^{d_0}/\sigma_t^4)$.
1303 The parameter bounds for the transformer network class are as follows:

$$1305 \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{\frac{7\beta_0}{d_0+d_y}+6C_\sigma}\right);$$

$$1307 \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{3\beta_0}{d_0+d_y}+6C_\sigma}(\log N)^{3(d_0+\beta_0)}\right);$$

$$1309 \|W_V\|_2 = \mathcal{O}(\sqrt{\tilde{d}}); \quad \|W_V\|_{2,\infty} = \mathcal{O}(\tilde{d});$$

$$1311 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{2\beta_0}{d_0+d_y}+4C_\sigma}\right); \|E^\top\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}}\tilde{L}^{\frac{3}{2}}\right);$$

$$1313 \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{3\beta_0}{d_0+d_y}+2C_\sigma}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t^2\right).$$

1316 *Proof Sketch.* The proof closely follows [Theorem 3.1](#), with differences highlighted in [\(K1\)](#) and [\(K2\)](#).
1317 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](#),
1318 and under the the β_0 -Hölder smoothness assumption on latent data detailed in [\(K2\)](#), the proof is
1319 complete. Please see [Appendix I](#) for a detailed proof. \square

1321 **Theorem E.2** (Latent Conditional DiT Score Approximation under Stronger Hölder Assumption
1322 under Generic Hölder Assumption (Formal Version of [Theorem 4.1](#))). Assume [Assumption 3.2](#)
1323 and Assume $d_x = \Omega(\frac{\log N}{\log \log N})$. For any precision $0 < \epsilon < 1$ and smoothness $\beta_0 > 0$, let
1324 $\epsilon \leq \mathcal{O}(N^{-\beta_0})$ for some $N \in \mathbb{N}$. For some positive constants $C_\alpha, C_\sigma > 0$, for any $y \in [0, 1]^{d_y}$ and
1325 $t \in [N^{-C_\sigma}, C_\alpha \log N]$, there exists a $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_{\tilde{R}}^{h,s,r}$ such that

$$1328 \int_{\mathbb{R}^{d_0}} \|\mathcal{T}_{\text{score}}(x, y, t)(\bar{h}, y, t) - \nabla \log p_t^h(\bar{h}|y)\|_2^2 \cdot p_t^h(\bar{h}|y) d\bar{h} = \mathcal{O}\left(\frac{B_0^2}{\sigma_t^2} \cdot N^{-\frac{2\beta_0}{d_0+d_y}} \cdot (\log N)^{\beta_0+1}\right).$$

1331 Notably, for $\epsilon = \mathcal{O}(N^{-\beta_0})$, the approximation error has the upper bound $(\log(\frac{1}{\epsilon}))^{\mathcal{O}(1)}/\sigma_t^2$.
1332 The parameter bounds in the transformer network class satisfy

$$1334 \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{\frac{3\beta_0(2d_0+4\tilde{d}+1)}{\tilde{d}(d_0+d_y)} + \frac{9C_\alpha(2d_0+4\tilde{d}+1)}{\tilde{d}}}\right);$$

$$1336 \|W_V\|_2 = \mathcal{O}(\sqrt{\tilde{d}}); \|W_V\|_{2,\infty} = \mathcal{O}(\tilde{d}); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{\beta_0}{d_0+d_y}}\right);$$

$$1338 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{4\beta_0}{d_0+d_y}+9C_\sigma+\frac{3C_\alpha}{2}} \cdot \log N\right); \|E^\top\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}}\tilde{L}^{\frac{3}{2}}\right);$$

$$1340 \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{4\beta_0}{d_0+d_y}+9C_\sigma+\frac{3C_\alpha}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t\right).$$

1344 *Proof Sketch.* The proof closely follows [Theorem J.1](#), with differences highlighted in [\(K1\)](#) and [\(K2\)](#).
1345 By replacing the input dimension d, L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 in [Theorem J.1](#),
1346 and under the the β_0 -Hölder smoothness assumption on latent data detailed in [\(K2\)](#), the proof is
1347 complete. Please see [Appendix J](#) for a detailed proof. \square

1349 **Remark E.2** (Score Approximation for Low-Dimensional Linear Latent Space). With the assumption
of low-dimensional latent space [Assumption 4.1](#), [Theorems E.1](#) and [E.2](#) provide better approximation

1350 rates than **Theorems 3.1** and **3.2** under Hölder smooth assumptions in **Assumptions 3.1** and **3.2**,
 1351 respectively. Specifically, from **Lemma E.2** we have $\beta_0 = \beta$ and $B_0 \lesssim B$. Therefore, **Theorems E.1**
 1352 and **E.2** deliver $\mathcal{O}\left(N^{2\beta\left(\frac{d_x-d_0}{(d_0+d_y)(d_x+d_y)}\right)}\right)$ better approximation error over **Theorem 3.1**, where
 1353 $d_0 \leq d_x$.
 1354
 1355

1356 E.2 SCORE ESTIMATION

1357
 1358 In this section, we provide the extended results for **Section 3.3** on score estimation with the estimator
 1359 $\mathcal{T}_{\text{score}}$. We state the main results under Hölder data assumptions in **Theorem E.3**.
 1360

Theorem E.3 (Conditional Score Estimation of Latent DiT). **Assume** $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$. Let \hat{s}
 1361 denote the score estimator trained with a set of finite samples $\{x_i, y_i\}_{i \in [n]}$ by optimizing the empirical
 1362 loss (2.1), and \mathcal{R} denote the conditional score risk defined in **Definition 3.2**.
 1363

- 1364 • Under **Assumption 3.1**, by taking $N = n^{\frac{1}{\tilde{\nu}_1} \cdot \frac{d_0+d_y}{\beta_0+d_0+d_y}}$, $t_0 = N^{-C_\sigma} < 1$ and $T = C_\alpha \log n$, it holds
 1365

$$1366 \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] = \mathcal{O}\left(\frac{1}{t_0} n^{-\frac{\beta_0}{\tilde{\nu}_1(d_0+d_y+\beta_0)}} (\log n)^{\tilde{\nu}_2+2}\right),$$

1367
 1368 where $\tilde{\nu}_1 = 68\beta_0/(d_0 + d_y) + 104C_\sigma$ and $\tilde{\nu}_2 = 12d_0 + 12\beta_0 + 2$.
 1370

- 1371 • Under **Assumption 3.2**, by taking $N = n^{\frac{1}{\tilde{\nu}_3} \cdot \frac{d_0+d_y}{2\beta_0+d_0+d_y}}$, $t_0 = N^{-C_\sigma} < 1$ and $T = C_\alpha \log n$, it
 1372 holds
 1373

$$1374 \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] = \mathcal{O}\left(\log \frac{1}{t_0} n^{-\frac{1}{\tilde{\nu}_3} \frac{2\beta_0}{d_0+d_y+2\beta_0}} (\log n)^{\max(10, \beta_0+1)}\right),$$

1375
 1376 where $\tilde{\nu}_3 = \frac{4(12\beta_0 d_0 + 31\beta_0 \tilde{d} + 6\beta_0)}{\tilde{d}(d_0+d_y)} + \frac{12(12C_\alpha d_0 + 25C_\alpha \tilde{d} + 6C_\alpha)}{\tilde{d}} + 72C_\sigma$.
 1377
 1378

1379
 1380 *Proof Sketch.* The proof closely follows **Theorem 3.3**, with differences highlighted in **(K1)** and **(K2)**.
 1381 By replacing the input dimension d , L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 in **Theorem 3.3**,
 1382 and under the the β_0 -Hölder smoothness assumption on latent data detailed in **(K2)**, the proof is
 1383 complete. Please see **Appendix K.2** for a detailed proof. \square
 1384

1385 Next, we present the score estimation result for low-dimensional input data.
 1386

Corollary E.3.1 (Low-Dimensional Input Region). **Assume** $d_0 = o\left(\frac{\log N}{\log \log N}\right)$, i.e., $d_0 \ll n$. Under
 1387 **Assumption 3.1**, by setting N, t_0, T as specified in **Theorem E.3**, we have $\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] =$
 1388 $\mathcal{O}\left(\frac{1}{t_0} n^{-\frac{1}{\tilde{\nu}_4} \cdot \frac{\beta_0}{d_0+d_y+\beta_0}}\right)$, where $\tilde{\nu}_4 = \frac{72\beta_0(2d_0+5\tilde{d}+1)}{\tilde{d}(d_0+d_y)} + \frac{48C_\sigma(2d_0+5\tilde{d}+1)}{\tilde{d}} - 4\beta_0$.
 1389
 1390
 1391
 1392

1393
 1394 *Proof.* The proof closely follows **Corollary 3.3.1**, with differences highlighted in **(K1)** and **(K2)**. By
 1395 replacing the input dimension d , L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 in **Corollary 3.3.1**,
 1396 and under the the β_0 -Hölder smoothness assumption on latent data detailed in **(K2)**, the proof is
 1397 complete. Please see **Appendix K.2** and **Appendix K.4** for detailed proofs. \square
 1398

Remark E.3 (Comparing Score Estimation in **Theorems 3.3** and **E.3**). Under Hölder data assumption,
 1399 the sample complexity of L_2 estimator for achieving ϵ -error are bound by $\tilde{\mathcal{O}}\left(\epsilon^{-\tilde{\nu}_1(d_0+d_y+\beta_0)/\beta_0}\right)$
 1400 and $\tilde{\mathcal{O}}\left(\epsilon^{-\tilde{\nu}_3(d_0+d_y+2\beta_0)/\beta_0}\right)$ where $\tilde{\mathcal{O}}$ ignores \tilde{d} , \tilde{L} , $\log \tilde{L}$, $\log 1/t_0$, $1/t_0$, and $\log n$. Invoking
 1401 **Lemma E.2** where $\beta_0 = \beta$ and $B_0 \lesssim B$ the sample complexity in **Theorem E.3** improves
 1402 **Theorem 3.3** by $\mathcal{O}\left(\epsilon^{-\zeta(d_x-d_0)}\right)$ where ζ is a positive constant defined by $\zeta = 104C_\sigma/\beta -$
 1403 $68\beta(1/((d_x + d_y)(d_0 + d_y)))$ and $d_0 \leq d_x$.

1404 E.3 DISTRIBUTION ESTIMATION

1405
1406 In this section, we provide the extended results for [Section 3.3](#) on distribution estimation with the
1407 estimator $\mathcal{T}_{\text{score}}$. We state the main results under Hölder data assumptions in [Theorem E.3](#).

1408 **Theorem E.4** (Conditional Distribution Estimation of Latent DiT). Assume $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$. For
1409 all $y \in [0, 1]^{d_y}$, let $\text{KL}(P(\cdot|y)|N(0, I)) \leq c$ for some constant $c < \infty$. Taking the early-stopping
1410 time $t_0 = n^{-\frac{\beta_0}{(d_0+d_y+\beta_0)}}$ and terminal time $T = \frac{2\beta_0}{d_0+d_y+2\beta_0} \log n$.

- 1411 • Under [Assumption 3.1](#), we have

$$1412 \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\widehat{P}_{t_0}(\cdot|y), P(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{\beta}{2(\tilde{\nu}_1-1)(d_0+d_y+\beta_0)}} (\log n)^{\frac{\tilde{\nu}_2}{2} + \frac{3}{2}} \right),$$

1413 where $\tilde{\nu}_1 = 68\beta_0/(d_0 + d_y) + 104C_\sigma$ and $\tilde{\nu}_2 = 12d_0 + 12\beta_0 + 2$.

- 1414 • Under [Assumption 3.2](#), we have

$$1415 \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\widehat{P}_{t_0}(\cdot|y), P(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{1}{2\tilde{\nu}_3} \frac{\beta_0}{d_0+d_y+2\beta_0}} (\log n)^{\max(6, \frac{\beta_0}{2} + \frac{3}{2})} \right),$$

1416 where $\tilde{\nu}_3 = \frac{4(12\beta_0 d_0 + 31\beta_0 \tilde{d} + 6\beta_0)}{\tilde{d}(d_0 + d_y)} + \frac{12(12C_\alpha d_0 + 25C_\alpha \tilde{d} + 6C_\alpha)}{\tilde{d}} + 72C_\sigma$.

1417 *Proof.* The proof closely follows [Theorem 3.4](#), with differences highlighted in [\(K1\)](#) and [\(K2\)](#). By
1418 replacing the input dimension d, L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 in [Theorem 3.4](#),
1419 and under the the β_0 -Hölder smoothness assumption on latent data detailed in [\(K2\)](#), the proof is
1420 complete. Please see [Appendix K.6](#) for a detailed proof. \square

1421 Next, we present the distribution estimation result for low-dimensional input data.

1422 **Corollary E.4.1** (Low-Dimensional Input Region). Assume $d_0 = o\left(\frac{\log N}{\log \log N}\right)$, i.e., $d_0 \ll n$. Under
1423 [Assumption 3.1](#), by setting t_0, T as specified in [Theorem E.4](#), we have

$$1424 \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(\widehat{P}_{t_0}(\cdot|y), P_0(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{\beta_0}{2(\tilde{\nu}_4+1)(d_0+d_y+\beta_0)}} \right),$$

1425 where $\tilde{\nu}_4 = \frac{144\tilde{d}\beta_0(\tilde{L}+2)(d_0+2\tilde{d}+1)}{d_0+d_y} + 96\tilde{d}C_\sigma(\tilde{L}+2)(d_0+2\tilde{d}+1) - 4\beta_0$.

1426 *Proof.* The proof closely follows [Corollary 3.4.1](#), with differences highlighted in [\(K1\)](#) and [\(K2\)](#). By
1427 replacing the input dimension d, L to \tilde{d} and \tilde{L} , and the input dimension d_x with d_0 in [Corollary 3.4.1](#),
1428 and under the the β_0 -Hölder smoothness assumption on latent data detailed in [\(K2\)](#), the proof is
1429 complete. Please see [Appendix K.6](#) for a detailed proof. \square

F LATENT CONDITIONAL DiT WITH LIPSCHITZ ASSUMPTION

In this section, we apply our techniques to the setting of (Hu et al., 2024) on DiT approximation and estimation theory. Specifically, we extend their work by using the one-layer self-attention transformer universal approximation framework introduced in Appendix H.1.

Compared to (Hu et al., 2024), we consider classifier-free conditional DiTs, providing a holistic view of the theoretical guarantees under various assumptions. In particular, our sample complexity bounds avoid the gigantic double exponential term $2^{(1/\epsilon)^{2L}}$ reported in (Hu et al., 2024). We adopt the following three assumptions considered by Hu et al. (2024):

(A1) Low-Dimensional Linear Latent Space Data Assumption.

Assumption F.1 (Low-Dimensional Linear Latent Space (Assumption 4.1 Restated)). Data point $x = Uh$, where $U \in \mathbb{R}^{d_x \times d_0}$ is an unknown matrix with orthonormal columns. The latent variable $h \in \mathbb{R}^{d_0}$ follows a distribution P_h with a density function p_h .

Under this data assumption, Chen et al. (2023a) show that the latent score function endows a neat decomposition into on-support s_+ and orthogonal s_- terms (see Lemma E.1).

Lemma F.1 (Score Decomposition, Lemma 1 of (Chen et al., 2023c) (Lemma E.1 Restated)). Let data $x = Uh$ follow Assumption 4.1. The decomposition of score function $\nabla \log p_t(x)$ is

$$\nabla \log p_t(x) = \underbrace{U \nabla \log p_t^h(\bar{h}|y)}_{s_+(\bar{h}, y, t)} - \underbrace{(I_D - UU^\top) x / \sigma_t^2}_{s_-(x, t)}, \quad \bar{h} = U^\top x, \quad (\text{F.1})$$

where $p_t^h(\bar{h}|y) := \int \psi_t(\bar{h}|h) p_h(h|y) dh$, $\psi_t(\cdot|h)$ is the Gaussian density function of $N(\beta_t h, \sigma_t^2 I_{d_0})$, $\beta_t = e^{-t/2}$ and $\sigma_t^2 = 1 - e^{-t}$.

(A2) Lipschitz Score Assumption. We assume the on-support score function $s_+(\bar{h}, y, t)$ to be L_{s_+} -Lipschitz for any \bar{h} and y .

Assumption F.2 (L_{s_+} -Lipschitz of $s_+(\bar{h}, y, t)$). The on-support score function $s_+(\bar{h}, y, t)$ is L_{s_+} -Lipschitz with respect to any $\bar{h} \in \mathbb{R}^{d_0}$ and $y \in \mathbb{R}^{d_y}$ for any $t \in [0, T]$. i.e., there exist a constant L_{s_+} , such that for any \bar{h}, y and \bar{h}', y' :

$$\|s_+(\bar{h}, y, t) - s_+(\bar{h}', y', t)\|_2 \leq L_{s_+} \|\bar{h} - \bar{h}'\|_2 + L_{s_+} \|y - y'\|_2.$$

(A3) Light Tail Data Assumption.

Assumption F.3 (Tail Behavior of P_h). The density function $p_h > 0$ is twice continuously differentiable. Moreover, there exist positive constants A_0, A_1, A_2 such that when $\|h\|_2 \geq A_0$, the density function $p_h(h|y) \leq (2\pi)^{-d_0/2} A_1 \exp(-A_2 \|h\|_2^2/2)$.

We note that, the assumptions (A1) and (A3) are on data, and (A2) are on the score function. Notably, (A2) on the smoothness of score function is stronger than Hölder data smoothness assumptions considered in Sections 3 and 4.

Organization. We study latent conditional DiTs under low-dimensional data Assumption F.1, Lipschitz smoothness Assumption F.2, and tail behavior of P_h Assumption F.3 and states the results in Appendices F.1 to F.3, respectively. Appendix F.1 discusses score approximation. Appendix F.2 discusses score estimation. Appendix F.3 discusses distribution estimation. The proof in this section provided in Appendices F.4 to F.6. The proof strategy in this section follows (Hu et al., 2024).

Here we summarize the key settings of this section:

S1. Lipschitz Smooth Assumption and Tail Behavior. Following (Hu et al., 2024), we introduce two assumptions on Lipschitz smoothness for on-support score function s_+ and tail behavior of P_h in Assumptions F.2 and F.3, respectively. The on-support score function is defined as $s_+(U^\top x, y, t) = U \nabla \log p_t^h(U^\top x|y)$ (see Lemma E.1 for score decomposition).

- 1512 **S2. Low-Dimensional Space.** We consider the setting of latent representation that is the data lives
 1513 in a “Low-Dimensional Subspace” under [Assumption 4.1](#), following ([Hu et al., 2024](#); [Chen](#)
 1514 [et al., 2023c](#)). The raw data $x \in \mathbb{R}^{d_x}$ is supported by latent $h \in \mathbb{R}^{d_0}$ where $d_0 \leq d_x$.
 1515
- 1516 **S3. Transformer Network.** We follow the standard setting of “in-context” conditional DiTs by
 1517 [Peebles and Xie \(2023\)](#) on latent representation. The network settings refer to [Section 4](#). Here we
 1518 apply transformer-block $g_{\mathcal{T}} \in \mathbb{R}^{d_0}$ for the approximation of on-support score function s_+ . For
 1519 each input $x \in \mathbb{R}^{d_x}$ and corresponding label $y \in \mathbb{R}^{d_y}$, we use an adapted transformer network to
 1520 obtain a score estimator $s_W \in \mathbb{R}^{d_0}$. The adapted transformer network as the score estimator has
 1521 the following components. We utilize reshape layer to convert vector input $h \in \mathbb{R}^{d_0}$ to matrix
 1522 (sequence) input $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$. Specifically, the reshape layer in the network [Figure 3](#) is defined as
 1523 $\tilde{R}(\cdot) : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}$ and its reverse $\tilde{R}^{-1}(\cdot) : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{d_0}$, where $d_0 \leq d_x$, $\tilde{d} \leq d$, and $\tilde{L} \leq L$.
 1524 For raw data $x \in \mathbb{R}^{d_x}$, we utilize linear encoder $W_U^T \in \mathbb{R}^{d_0 \times d_x}$ and decoder $W_U \in \mathbb{R}^{d_x \times d_0}$
 1525 to convert the raw $x \in \mathbb{R}^{d_x}$ to latent $h \in \mathbb{R}^{d_0}$ data representations. Importantly, $x = Uh$ with
 1526 $U \in \mathbb{R}^{d_x \times d_0}$ by [Assumption 4.1](#).

1527
 1528 Under the [Assumptions F.1](#) to [F.3](#) with the network setting following ([S3](#)), the theoretical results in
 1529 [Appendices F.1](#) to [F.3](#) achieve tighter approximation rates and efficient recovery accuracy of latent
 1530 data detailed in ([R1](#)), ([R2](#)), and ([R3](#)).

1531 We summarize the theoretical comparisons from [Appendix E](#) and [Appendix F](#) as follows:

1532
 1533 **R1.** For score approximation (see [Theorems E.1](#), [E.2](#) and [F.1](#)):

- 1534
 1535 – Under Hölder data assumption the approximation rates gives $\tilde{\mathcal{O}}(\epsilon^{1/(d_0+d_y)})$, where $\tilde{\mathcal{O}}$ ignores
 1536 B_0 , $\log \epsilon$, and $\log n$.
 1537
 1538 – Under Lipschitz score assumption the approximation rate gives $\tilde{\mathcal{O}}(\epsilon \cdot \sqrt{d_0 + d_y})$, where $\tilde{\mathcal{O}}$
 1539 ignores B_0 , $\log \epsilon$, and $\log n$.
 1540
 1541 – For any precision $0 < \epsilon < 1$, the Lipschitz score assumption provides a tighter approximate
 1542 rate for high dimension data $d_0 \gg 1$ compared with under Hölder data assumption.

1543 **R2.** For score estimation (see [Theorems E.3](#) and [F.2](#)):

- 1544
 1545 – Under Hölder data assumption the score estimation error gives $\tilde{\mathcal{O}}\left(n^{-\frac{1}{\tilde{\nu}_3} \cdot \frac{\beta_0}{d_0+d_y+2\beta_0}}\right)$, where
 1546 $\tilde{\mathcal{O}}$ ignores B_0 , $\log \epsilon$, and $\log n$.
 1547
 1548 – Under Lipschitz score assumption the score estimation error gives $\tilde{\mathcal{O}}\left(n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}}\right)$, where
 1549 $\tilde{\mathcal{O}}$ ignores B_0 , $\log \epsilon$, and $\log n$.
 1550
 1551 – Under minimax optimal condition (see [Section 3.4](#)) by setting $\tilde{\nu}_3 = 1/2$, Hölder data
 1552 assumption gives $\tilde{\mathcal{O}}\left(n^{-\frac{\beta_0}{2(d_0+d_y+2\beta_0)}}\right)$. On the other hand, Lipschitz assumption gives
 1553 $\tilde{\mathcal{O}}\left(n^{-\frac{\tilde{d}}{(3/4)d_0+(2/3)\tilde{d}+2}}\right)$. Therefore, the Lipschitz assumption gives a better sample complex-
 1554 ity guarantee for high dimensional data $d_0 = \tilde{d}\tilde{L} \gg 1$.
 1555
 1556
 1557
 1558
 1559

1560 **R3.** For distribution estimation (see [Theorems E.4](#) and [F.3](#)):

- 1561
 1562 – Under Hölder data assumption: $\tilde{\mathcal{O}}\left(n^{-\frac{1}{\tilde{\nu}_3} \cdot \frac{\beta_0}{2(d_0+d_y+2\beta_0)}}\right)$.
 1563
 1564
 1565 – Under Lipschitz score assumption: $\tilde{\mathcal{O}}\left(n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}}\right)$.

- Follow the arguments in (R2), Lipschitz assumption gives a better distribution estimation guarantee for high dimensional data.

Note that d_0, d_y is the latent data dimension and conditioning label dimension and $\tilde{\nu}_3 = \frac{4(12\beta_0 d_0 + 31\beta_0 \tilde{d} + 6\beta_0)}{\tilde{d}(d_0 + d_y)} + \frac{12(12C_\alpha d_0 + 25C_\alpha \tilde{d} + 6C_\alpha)}{\tilde{d}} + 72C_\sigma$.

From (R1), (R2), and (R3), we conclude that stronger approximations yield sharper rates.

F.1 SCORE APPROXIMATION

For completeness, we follow the proofs from (Hu et al., 2024) for score approximation of the conditional latent diffusion model.

Here we use stricter assumptions on the latent density function, instead of assuming Hölder smoothness of the initial conditional data distribution as in Section 4. To be specific, we directly approximate the on-support latent score function, instead of approximating the denominator and nominator separately. From the score decomposition in (4.1), we define the on-support score function s_+ as following:

$$\begin{aligned} s_+(U^\top x, y, t) &= U \int \frac{\nabla_{\bar{h}} \psi_t(\bar{h}|h) p_h(h|y)}{\int \psi_t(\bar{h}|h') p_{h'}(h'|y) dh'} dh \\ &= U \nabla \log p_t^h(U^\top x|y). \end{aligned} \quad (\text{F.2})$$

Here we require two assumptions following the proof of (Hu et al., 2024) on tail behavior of density function and Lipschitz continuous for on-support score function. **Assumption F.3** is the analogy of **Assumption 3.1** for assuming the tail behavior of the density function. On the other hand, **Assumption F.2** further assume the on-support score function s_+ to be L_{s_+} -Lipshitz. Note that this assumption is stricter than **Assumption 3.1** since we make the Lipschitz assumption directly on the score function instead of on the latent density function.

Theorem F.1 (Latent Score Approximation of Conditional DiT, modified from Theorem 3.1 in Hu et al. (2024)). For any approximation error $\epsilon > 0$ and any data distribution P_0 under **Assumptions 4.1, F.2** and **F.3**, there exists a DiT score network $\mathcal{T}_{\text{score}}(\bar{h}, y, t) \in \mathcal{T}_R^{h,s,r}$ where $W = \{W_U, \mathcal{T}_{\text{score}}\}$, such that for any $t \in [t_0, T]$, we have:

$$\|\mathcal{T}_{\text{score}}(\cdot, t) - \nabla \log p_t(\cdot)\|_{L^2(P_t)} \leq \epsilon \cdot \sqrt{d_0 + d_y} / \sigma_t^2,$$

where $\sigma_t^2 = 1 - e^{-t}$ and the parameter bounds in the transformer network class satisfy

$$\begin{aligned} \|W_Q\|_2 &= \|W_K\|_2 = \mathcal{O}\left(\tilde{d} \cdot \epsilon^{-\left(\frac{1}{\tilde{d}} + 2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\ \|W_Q\|_{2,\infty} &= \|W_K\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{3}{2}} \cdot \epsilon^{-\left(\frac{1}{\tilde{d}} + 2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\ \|W_O\|_2 &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{\frac{1}{\tilde{d}}}\right); \|W_O\|_{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{\tilde{d}}}\right); \\ \|W_V\|_2 &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}}\right); \|W_V\|_{2,\infty} = \mathcal{O}(\tilde{d}); \\ \|W_1\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{\tilde{d}}}\right), \|W_1\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{-\frac{1}{\tilde{d}}}\right); \\ \|W_2\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{\tilde{d}}}\right); \|W_2\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{-\frac{1}{\tilde{d}}}\right); \\ \|E^\top\|_{2,\infty} &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \tilde{L}^{\frac{3}{2}}\right). \end{aligned}$$

Proof. Please see **Appendix F.4** for a detailed proof. \square

Remark F.1 (Comparing with Hölder Assumption Results in Low-Dimensional Data). Under **Assumptions 3.1** and **3.2**, the score approximation give us $\tilde{\mathcal{O}}\left(\epsilon^{\frac{1}{d_x+d_y}}/\sigma_t^4\right)$ and $\tilde{\mathcal{O}}\left(\epsilon^{\frac{1}{d_x+d_y}}/\sigma_t^2\right)$ in **Theorems E.1** and **E.2**, respectively. On the other hand, the direct approximation of the Lipschitz smooth on-support score function gives us the approximation error of $\mathcal{O}\left(\epsilon \cdot \sqrt{d_0+d_y}/\sigma_t^2\right)$. For $(d_0+d_y) \gg 1$, **Theorem F.1** delivers superior approximation error compare with **Theorems E.1** and **E.2**.

F.2 SCORE ESTIMATION

Theorem F.2 (Score Estimation of Latent DiT). Under the **Assumptions F.1** to **F.3**, we choose the score network $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_{\tilde{R}}^{h, s, r}$ from **Theorem F.1** using $\epsilon \in (0, 1)$ and $\tilde{L} > 1$. With probability $1 - 1/\text{poly}(n)$, we have

$$\frac{1}{T-t_0} \int_{t_0}^T \|\mathcal{T}_{\text{score}}(\cdot, t) - \nabla \log p_t(\cdot)\|_{L^2(P_t)} dt = \tilde{\mathcal{O}}\left(\frac{1}{t_0^2} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \log^3 \tilde{L} \log^3 n\right),$$

where $\tilde{\mathcal{O}}$ hides the factor about $d_x, d_y, d_0, \tilde{d}, L_{s+}$ and $\delta(n)$ is negligible for sufficiently large n .

Proof. Please see **Appendix F.5** for a detailed proof. \square

Remark F.2 (Comparing Score Estimation in **Theorems E.3** and **F.2**). Under Hölder data assumption, the sample complexity of L_2 estimator for achieving ϵ -error are bound by $\tilde{\mathcal{O}}\left(\epsilon^{-\tilde{\nu}_1(d_0+d_y+\beta_0)}/\beta_0\right)$ and $\tilde{\mathcal{O}}\left(\epsilon^{-\tilde{\nu}_3(d_0+d_y+2\beta_0)}/\beta_0\right)$. In contrast, **Theorem F.2** has the sample complexity bound of $\tilde{\mathcal{O}}\left(\epsilon^{-2(1+3/\tilde{d}+4\tilde{L})/3}\right)$. Therefore, a direct approximation of the Lipschitz smooth score function offers a better sample complexity bound than Hölder data assumption.

F.3 DISTRIBUTION ESTIMATION

In practice, DiTs generate data using the discretized version with step size μ . Let \hat{P}_{t_0} be the distribution generated by $\mathcal{T}_{\text{score}}(x, y, t)$ in **Theorem F.2**. Let $P_{t_0}^h$ and $p_{t_0}^h$ be the distribution and density function of on-support latent variable \bar{h} at t_0 . We have the following results for distribution estimation.

Theorem F.3 (Distribution Estimation of DiT, Modified From Theorem 3 of (Chen et al., 2023c)). Let $T = \mathcal{O}(\log n)$, $t_0 = \mathcal{O}(\min\{c_0, 1/L_{s+}\})$, where c_0 is the minimum eigenvalue of $\mathbb{E}_{P_h}[hh^\top]$. With the estimated DiT score network $\mathcal{T}_{\text{score}}(x, y, t)$ in **Theorem F.2**, we have the following with probability $1 - 1/\text{poly}(n)$.

(i) The accuracy to recover the subspace U is

$$\|W_U W_U^\top - U U^\top\|_F^2 = \tilde{\mathcal{O}}\left(\frac{1}{c_0} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \cdot \log^3 n\right). \quad (\text{F.3})$$

(ii) $(W_B U)_{\#}^\top \hat{P}_{t_0}$ denotes the pushforward distribution. With the conditions $\text{KL}(P_h \| N(0, I_{d_0})) < \infty$, and step size $\mu \leq \xi(n, t_0, L) \cdot t_0^2 / (d_0 \sqrt{\log d_0})$. There exists an orthogonal matrix $U \in \mathbb{R}^{d \times d}$ such that we have the following upper bound for the total variation distance

$$\text{TV}(P_{t_0}^h, (W_B U)_{\#}^\top \hat{P}_{t_0}) = \tilde{\mathcal{O}}\left(\frac{1}{t_0 \sqrt{c_0}} n^{\frac{-3}{4(1+3/\tilde{d}+4\tilde{L})}} \cdot \log^4 n\right), \quad (\text{F.4})$$

where $\tilde{\mathcal{O}}$ hides the factor about d_x, d_0, d , and L_{s+} .

(iii) For the generated data distribution \hat{P}_{t_0} , the orthogonal pushforward $(I - W_B W_B^\top)_{\#} \hat{P}_{t_0}$ is $N(0, \Sigma)$, where $\Sigma \preceq a t_0 I$ for a constant $a > 0$.

1674 *Proof.* Please see [Appendix F.6](#) for a detailed proof. □

1675
1676 **Remark F.3** (Compare with Existing Work). In ([Chen et al., 2023c](#), Theorem 3), the upper bound
1677 for total variation distance with ReLU network is $\tilde{\mathcal{O}}\left(\sqrt{1/(c_0 t_0)} n^{-1/(d+5)} \log^2 n\right)$. Therefore, for
1678 $n \gg 1$, [Theorem F.3](#) gives tighter accuracy if $3d + 11 > 12/\tilde{d} + 16\tilde{L}$ where $\tilde{d} \leq d$ and $\tilde{L} \leq L$. On
1679 the other hand, under similar conditions for d and L , [Theorem F.3](#) suggest to achieve similar total
1680 variation distance we only require $\sqrt{t_0}$ early stopping time which is beneficial for empirical setting.
1681
1682

1683 F.4 PROOF OF SCORE APPROXIMATION ([THEOREM F.1](#))

1684
1685 To begin the proof of the approximate theorem, we first restate some auxiliary lemmas and their
1686 proofs here from ([Chen et al., 2023c](#)) for later convenience. Note that some of the proofs extend to
1687 the latent density function.

1688 **Lemma F.2** (Modified from Lemma 16 in ([Chen et al., 2023c](#))). Consider a probability density
1689 function $p_h(h|y) = \exp(-C\|h\|_2^2/2)$ for $h \in \mathbb{R}^{d_0}$ and constant $C > 0$. Let $r_h > 0$ be a fixed
1690 radius. Then it holds

$$1691 \int_{\|h\|_2 > r_h} p_h(h|y) dh \leq \frac{2d_0\pi^{d_0/2}}{C\Gamma(d_0/2 + 1)} r_h^{d_0-2} \exp(-Cr_h^2/2),$$

$$1692 \int_{\|h\|_2 > r_h} \|h\|_2^2 p_h(h|y) dh \leq \frac{2d_0\pi^{d_0/2}}{C\Gamma(d_0/2 + 1)} r_h^{d_0} \exp(-Cr_h^2/2).$$

1693
1694 **Lemma F.3** (Modified from Lemma 2 in ([Chen et al., 2023c](#))). Suppose Assumption [Assumption F.3](#)
1695 holds and q is defined as:

$$1696 q(\bar{h}, y, t) = \int \frac{h\psi_t(\bar{h}|h) p_h(h|y)}{\int \psi_t(\bar{h}|h) p_h(h|y) dh} dh, \quad \bar{h} = B^\top x.$$

1700
1701 Given $\epsilon > 0$, with $r_h = c\left(\sqrt{d_0 \log(d_0/t_0)} + \log(1/\epsilon)\right)$ for an absolute constant c , it holds

$$1702 \|q(\bar{h}, y, t) \mathbb{1}\{\|\bar{h}\|_2 \geq r_h\}\|_{L^2(P_t)} \leq \epsilon, \text{ for } t \in [t_0, T].$$

1703
1704 **Lemma F.4** (Modified from Theorem 1 in ([Chen et al., 2023c](#))). We denote

$$1705 \tau(r_h) = \sup_{t \in [t_0, T]} \sup_{\bar{h} \in [0, r_h]^{d_0}} \sup_{y \in [0, 1]^{d_y}} \left\| \frac{\partial}{\partial t} q(\bar{h}, y, t) \right\|_2.$$

1706
1707 With $q(\bar{h}, y, t) = \int h\psi_t(\bar{h}|h)p_h(h|y)/(\int \psi_t(\bar{h}|h)p_h(h|y)dh)$ and p_h satisfies [Assumption F.3](#),
1708 we have a coarse upper bound for $\tau(r_h)$

$$1709 \tau(r_h) = \mathcal{O}\left(\frac{1 + \beta_t^2}{\beta_t} \left(L_{s_+} + \frac{1}{\sigma_t^2}\right) \sqrt{d_0} r_h\right) = \mathcal{O}\left(e^{T/2} L_{s_+} r_h \sqrt{d_0}\right).$$

1710
1711 *Proof of [Lemma F.4](#).*

$$1712 \frac{\partial}{\partial t} q(\bar{h}, y, t) = U \int \frac{h \frac{\partial}{\partial t} \psi_t(\bar{h}|h) p_h(h|y)}{\int \psi_t(\bar{h}|h) p_h(h|y) dh} dh - U \int \frac{h \psi_t(\bar{h}|h) p_h(h|y) \int \frac{\partial}{\partial t} \psi_t(\bar{h}|h) p_h(h|y) dh}{\left(\int \psi_t(\bar{h}|h) p_h(h|y) dh\right)^2} dh$$

$$1713 = U \int \frac{h \frac{\beta_t}{\sigma_t^2} \left(\|h\|_2^2 - (1 + \beta_t^2) h^\top \bar{h} + \beta_t \|\bar{h}\|_2^2\right) \psi_t(\bar{h}|h) p_h(h|y)}{\int \psi_t(\bar{h}|h) p_h(h|y) dh} dh$$

$$\begin{aligned}
& -U \int \frac{h\psi_t(\bar{h}|h)p_h(h|y) \int \frac{\beta_t}{\sigma_t^2} \left(\|h\|_2^2 - (1 + \beta_t^2)h^\top \bar{h} + \beta_t \|\bar{h}\|_2^2 \right) \psi_t(\bar{h}|h)p_h(h|y)dh}{\left(\int \psi_t(\bar{h}|h)p_h(h|y)dh \right)^2} dh \\
& \stackrel{(i)}{=} \frac{\beta_t}{\sigma_t^2} U \left[\mathbb{E}_{P_h} \left[h\|h\|_2^2 \right] - (1 + \beta_t^2) \text{Cov} \left[h|\bar{h} \right] \bar{h} \right],
\end{aligned}$$

where we plug in $\partial\psi_t(\bar{h}|h)/\partial t = \beta_t \left(\|h\|_2^2 - (1 + \beta_t^2)h^\top \bar{h} + \beta_t \|\bar{h}\|_2^2 \right) \psi_t(\bar{h}|h)/\sigma_t^2$ and collect terms in (i). Since P_h has a Gaussian tail, its third moment is bounded.

Then we bound $\|\text{Cov}[h|\bar{h}]\|_{\text{op}}$ by taking derivative of $s_+(\bar{h}, y, t)$ with respect to \bar{h} , here

$$s_+(\bar{h}, y, t) = U \frac{\beta_t}{\sigma_t^2} \int \frac{h \cdot \psi_t(\bar{h}|h)p_h(h|y)}{\int \psi_t(\bar{h}|h)p_h(h|y)dh} dh - U \frac{\bar{h}}{\sigma_t^2}.$$

Then we have

$$\begin{aligned}
\frac{\partial}{\partial \bar{h}} s_+(\bar{h}, y, t) &= \left(\frac{\beta_t}{\sigma_t^2} \right)^2 U \left[\int h h^\top \varphi(\bar{h}, y) dh - \int h \varphi(\bar{h}, y) dh \int h^\top \varphi(\bar{h}, y) dh \right] - \frac{1}{\sigma_t^2} U \\
&= \left(\frac{\beta_t}{\sigma_t^2} \right)^2 U \left[\text{Cov}(h|\bar{h}) - \frac{1}{\sigma_t^2} I_{d_0} \right],
\end{aligned}$$

where

$$\varphi(\bar{h}, y) = \frac{\psi_t(\bar{h}|h)p_h(h|y)}{\int \psi_t(\bar{h}|h)p_h(h|y)dh}.$$

Along with the L_{s_+} -Lipschitz property of s_+ , we obtain

$$\|\text{Cov}(h|\bar{h})\|_{\text{op}} \leq \frac{\sigma_t^4}{\beta_t^2} \left(L_{s_+} + \frac{1}{\sigma_t^2} \right).$$

Therefore, we deduce

$$\tau(r_h) = \mathcal{O} \left(\frac{1 + \beta_t^2}{\beta_t} \left(L_{s_+} + \frac{1}{\sigma_t^2} \right) \sqrt{d_0} r_h \right) = \mathcal{O} \left(e^{T/2} L_{s_+} r_h \sqrt{d_0} \right),$$

as P_h having sub-Gaussian tail implies $\mathbb{E}_{P_h} \left[h\|h\|_2^2 \right]$ is bounded. \square

Lemma F.5 (Modified from Lemma 10 in (Chen et al., 2023c)). For any given $\epsilon > 0$, and L -Lipschitz function g defined on $[0, 1]^{d_0} \times [0, 1]^{d_y}$, there exists a continuous function \bar{f} constructed by trapezoid function that

$$\|g - \bar{f}\|_{\infty} \leq \epsilon.$$

Moreover, the Lipschitz continuity of \bar{f} is bounded by

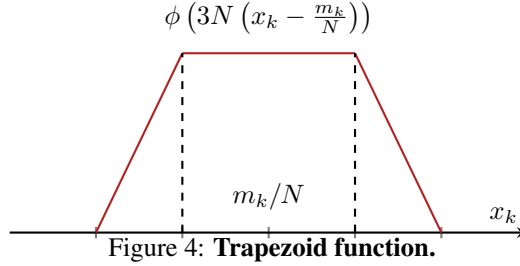
$$|\bar{f}(x, y) - \bar{f}(x', y')| \leq 10d_0L\|x - x'\|_2 + 10d_yL\|y - y'\|_2,$$

for any $x, x' \in [0, 1]^{d_0}$ and $y, y' \in [0, 1]^{d_y}$

Proof of Lemma F.5. This proof closely follows Lemma 10 in (Chen et al., 2023c). We divide the proof into two parts: First, we use a collection of Trapezoid function \bar{f} to approximate the function g defined on $[0, 1]^{d_0} \times [0, 1]^{d_y}$. Then we establish the Lipschitz continuity of the function \bar{f} to facilitate the approximation with a transformer.

1822 1. **Approximation by Trapezoid Function.** Given an integer $N > 0$, we choose $(N + 1)^{d_0}$ points
 1823 in the hypercube $[0, 1]^{d_0}$ and $(N + 1)^{d_y}$ points in the hypercube $[0, 1]^{d_y}$. We denote the index
 1824 of the hypercubes as $m = [m_1, m_2, \dots, m_{d_0}]^\top \in \{0, \dots, N\}$ and $n = [n_1, n_2, \dots, n_{d_y}]^\top \in$
 1825 $\{0, \dots, N\}$. Next, we define a univariate trapezoid function (see Figure 4) as follow
 1826

$$1827 \phi(a) = \begin{cases} 1, & |a| < 1 \\ 2 - |a|, & |a| \in [1, 2] \\ 0, & |a| > 2 \end{cases}. \quad (\text{F.5})$$



1800 Figure 4: Trapezoid function.

1803 For any $x \in [0, 1]^{d_0}$ and $y \in [0, 1]^{d_y}$, we define a partition of unity based on a product of trapezoid
 1804 functions indexed by m and n ,

$$1806 \xi_{m,n}(x, y) = \mathbf{1} \left\{ y \in \left(\frac{n-1}{N}, \frac{n}{N} \right] \right\} \prod_{k=1}^{d_0} \phi \left(3N \left(x_k - \frac{m_k}{N} \right) \right). \quad (\text{F.6})$$

1809 For example, the product of trapezoid function $\xi_{m,n}(x, y) \neq 0$ only if $y \in \left(\frac{n-1}{N}, \frac{n}{N} \right]$ and
 1810 $x \in \left[\frac{m-2 \cdot 1 \cdot 3}{N}, \frac{m+2 \cdot 1 \cdot 3}{N} \right]$. For any target L -Lipschitz function g with respect to x and y , it is more
 1811 convenient to write its Lipschitz continuity with respect to the ℓ_∞ norm, i.e.,
 1812

$$1813 |g(x, y) - g(x', y')| \leq L \|x - x'\|_2 + L \|y - y'\|_2 \\ 1814 \leq L \sqrt{d_0} \|x - x'\|_\infty + L \sqrt{d_y} \|y - y'\|_\infty. \quad (\text{F.7})$$

1816 We now define a collection of piecewise-constant functions as

$$1818 P_{m,n}(x, y) = g(m, n) \quad \text{for } m \in \{0, \dots, N\}^{d_0} \text{ and } n \in \{0, \dots, N\}^{d_y}.$$

1820 We claim that $\bar{f}(x, y) = \sum_{m,n} \xi_{m,n}(x, y) P_{m,n}(x, y)$ is an approximation of g , with an approxi-
 1821 mation error evaluated as
 1822

$$1823 \sup_{x \in [0, 1]^{d_0}} \sup_{y \in [0, 1]^{d_y}} |\bar{f}(x, y) - g(x, y)| \\ 1824 = \sup_{x \in [0, 1]^{d_0}} \sup_{y \in [0, 1]^{d_y}} \left| \sum_{m,n} \xi_{m,n}(x, y) (P_{m,n}(x, y) - g(x, y)) \right| \\ 1825 \leq \sup_{x \in [0, 1]^{d_0}} \sup_{y \in [0, 1]^{d_y}} \sum_{\substack{m: |x_k - m_k/N| \leq \frac{2}{3N} \\ n: |y_j - n_j/N| \in (-\frac{1}{2N}, \frac{1}{2N}]} |P_{m,n}(x, y) - g(x, y)| \\ 1826 = \sup_{x \in [0, 1]^{d_0}} \sup_{y \in [0, 1]^{d_y}} \sum_{\substack{m: |x_k - m_k/N| \leq \frac{2}{3N} \\ n: |y_j - n_j/N| \in (-\frac{1}{2N}, \frac{1}{2N}]} |g(m, n) - g(x, y)| \\ 1827 \leq L \sqrt{d_0} 2^{d_0+1} \frac{1}{3N} + L \sqrt{d_y} 1^{d_y} \frac{1}{2N} \quad (\text{By Lipschitz continuity in (F.7)})$$

$$= \frac{L}{N} \left(\frac{\sqrt{d_0} 2^{d_0+1}}{3} + \frac{\sqrt{d_y}}{2} \right),$$

where the last inequality follows the Lipschitz continuity in (F.7) and using the fact that there are at most 2^{d_0} terms in the summation of m and at most 1^{d_y} terms in the summation of n . By choosing $N = \lceil L (\sqrt{d_0} 2^{d_0+1}/3 + \sqrt{d_y}/2) / \epsilon \rceil$, we have $\|g - \bar{f}\|_\infty \leq \epsilon$.

2. **Lipschitz Continuity.** Next we compute the Lipschitz of the function \bar{f} with respect to x and y . Suppose the approximation error $\epsilon > 0$ is small enough, then we have

$$\begin{aligned} & |\bar{f}(x, y) - \bar{f}(x', y')| \\ & \leq |\bar{f}(x, y) - g(x, y)| + |g(x, y) - g(x', y')| + |g(x', y') - \bar{f}(x', y')| \\ & \leq 2\epsilon + L\sqrt{d_0}\|x - x'\|_\infty + L\sqrt{d_y}\|y - y'\|_\infty \\ & \leq 10L\sqrt{d_0}\|x - x'\|_\infty + 10L\sqrt{d_y}\|y - y'\|_\infty \\ & \leq 10Ld_0\|x - x'\|_2 + 10Ld_y\|y - y'\|_2. \end{aligned}$$

This completes the proof. \square

Main Proof of Theorem F.1. Now we are ready to state the main proof.

Proof of Theorem F.1. From low-dimensional data assumption, the score function $\log p_t(x|y)$ decomposes as the on-support and orthogonal component (see Lemma E.1). Recall the on-support score function is given by $\nabla \log p_t^h(\bar{h}|y) = U^\top s_+(\bar{h}, y, t)$ from (F.7). We use a latent score network to approximate the score function (see (K3)). Specifically, the latent score network includes a latent encoder and a latent decoder. The encoder approximates $U^\top \in \mathbb{R}^{d_0 \times d_x}$, and decoder approximates $U \in \mathbb{R}^{d_x \times d_0}$. At its core, we use the transformer $g_T(W_U^\top x, y, t) \in \mathcal{T}^{h,s,r}$ to approximate $q(\bar{h}, y, t)$ as defined in (E.1). The expression for $q(\bar{h}, y, t)$ is given by:

$$q(\bar{h}, y, t) = \sigma_t^2 \nabla \log p_t^h(U^\top x|y) + U^\top x = \sigma_t^2 U^\top (s_+(\bar{h}, y, t) + x/\sigma_t^2). \quad (\text{F.8})$$

We proceed as follows:

- **Step 1.** Approximate $q(\bar{h}, y, t)$ with a compact-supported continuous function $\bar{f}(\bar{h}, y, t)$.
- **Step 2.** Approximate $\bar{f}(\bar{h}, y, t)$ with a one-layer single-head transformer network.

Step 1. Approximate $q(\bar{h}, y, t)$ with a Compact-Supported Continuous Function $\bar{f}(\bar{h}, y, t)$. First, we partition \mathbb{R}^{d_0} into a compact subset $H_1 := \{\bar{h} \mid \|\bar{h}\|_2 \leq r_h\}$ and its complement H_2 , where the choice of r_h comes from Lemma F.3. Next, we approximate $q(\bar{h}, y, t)$ on the two subsets by using the compact-supported continuous function $\bar{f}(\bar{h}, y, t)$. Finally, calculating the continuity of \bar{f} gives an estimation error of $\sqrt{d_0 + d_y}\epsilon$ between $q(\bar{h}, y, t)$ and $\bar{f}(\bar{h}, y, t)$. We present the main proof as follows.

- **Approximation on $H_2 \times [0, 1] \times [t_0, T]$.** For any $\epsilon > 0$, by taking $r_h = c(\sqrt{d_0 \log(d_0/t_0) - \log \epsilon})$, we obtain from Lemma F.3 that

$$\|q(\bar{h}, y, t) \mathbb{1}\{\|\bar{h}\|_2 \geq r_h\}\|_{L^2(P_t)} \leq \epsilon \quad \text{for } t \in [t_0, T] \quad \text{and } y \in [0, 1].$$

So we set $\bar{f}(\bar{h}, y, t) = 0$ on $H_2 \times [0, 1] \times [t_0, T]$.

- 1890 • **Approximation on $H_1 \times [0, 1] \times [t_0, T]$.** On $H_1 \times [0, 1] \times [t_0, T]$, we approximate

1891
1892
$$q(\bar{h}, y, t) = [q_1(\bar{h}, y, t), q_2(\bar{h}, y, t), \dots, q_{d_0}(\bar{h}, y, t)],$$

1893
1894 by approximating each coordinate $q_k(\bar{h}, y, t)$ separately.

1895 We firstly rescale the input by $h' = (\bar{h} + r_h \mathbf{1})/2r_h$ and $t' = t/T$, so that the transformed input
1896 space is $[0, 1]^{d_0} \times [0, 1]^{d_y} \times [t_0/T, 1]$. Here we do not need to rescale y , since it is already in $[0, 1]$
1897 by definition. We implement such transformation by a single feed-forward layer.

1898 By [Assumption F.2](#), the on-support score $s_+(\bar{h}, y, t)$ is L_{s_+} -Lipschitz with respect to any $\bar{h} \in \mathbb{R}^{d_0}$
1899 and $y \in \mathbb{R}^{d_y}$. This implies $q(\bar{h}, y, t)$ is $(1 + L_{s_+})$ -Lipschitz in \bar{h} and y . When taking the
1900 transformed inputs, $g(h', y, t') = q(2r_h h' - r_h \mathbf{1}, T t')$ becomes $2r_h(1 + L_{s_+})$ -Lipschitz in h' ;
1901 each coordinate $g_k(h', y, t)$ is also $2r_h(1 + L_{s_+})$ -Lipschitz in h' . Here we denote $L_* = 1 + L_{s_+}$.
1902 Besides, $g(h', y, t')$ is $T\tau(r_h)$ -Lipschitz with respect to t , where

1903
1904
1905
$$\tau(r_h) = \sup_{t \in [t_0, T]} \sup_{\bar{h} \in [0, r_h]^{d_0}} \sup_{y \in [0, 1]^{d_y}} \left\| \frac{\partial}{\partial t} q(\bar{h}, y, t) \right\|_2.$$

1906 We have a coarse upper bound for $\tau(r_h)$ in [Lemma F.4](#). We restate it as follows:

1907
1908
$$\tau(r_h) = \mathcal{O} \left(\frac{1 + \beta_t^2}{\beta_t} \left(L_{s_+} + \frac{1}{\sigma_t^2} \right) \sqrt{d_0} r_h \right) = \mathcal{O} \left(e^{T/2} L_{s_+} r_h \sqrt{d_0} \right).$$

1909 Since each $g_k(h', y, t)$ is Lipschitz continuous, we apply [Lemma F.5](#) to construct a collection of
1910 coordinate-wise functions, denoted as $\bar{f}_k(h', y, t)$. We concatenate \bar{f}_k 's together and construct
1911 $\bar{f} = [\bar{f}_1, \dots, \bar{f}_{d_0}]^\top$. According to the construction of trapezoid function in [Lemma F.5](#), for any
1912 given ϵ , we have the following relations:

1913
1914
$$\sup_{h', y, t' \in [0, 1]^{d_0} \times [0, 1]^{d_y} \times [t_0/T, 1]} \left\| \bar{f}(h', y, t') - g(h', y, t') \right\|_\infty \leq \epsilon.$$

1915 Considering the input rescaling (i.e., $\bar{h} \rightarrow h'$, $y \rightarrow y$ and $t \rightarrow t'$), we obtain:

- 1916 – The constructed function is Lipschitz continuous in \bar{h} and y , i.e., for any $\bar{h}_1, \bar{h}_2 \in H_1$, $y_1, y_2 \in$
1917 $[0, 1]$ and $t \in [t_0, T]$, it holds

1918
1919
$$\left\| \bar{f}(\bar{h}_1, y_1, t) - \bar{f}(\bar{h}_2, y_2, t) \right\|_\infty \leq 10d_0 L_* \|\bar{h}_1 - \bar{h}_2\|_2 + 10d_y L_* \|y_1 - y_2\|_2. \quad (\text{F.9})$$

- 1920 – The function is also Lipschitz in t , i.e., for any $t_1, t_2 \in [t_0, T]$ and $\|\bar{h}\|_2 \leq r_h$, it holds

1921
1922
$$\left\| \bar{f}(\bar{h}, y, t_1) - \bar{f}(\bar{h}, y, t_2) \right\|_\infty \leq 10\tau(r_h) \|t_1 - t_2\|_2.$$

1923
1924 To conclude, the construction of $\bar{f}(\bar{h}, y, t)$ uses a collection of trapezoid functions, as described
1925 in [Lemma F.5](#). This ensures that $\bar{f}(\bar{h}, y, t) = 0$ for $\|\bar{h}\|_2 > r_h$, for all $t \in [t_0, T]$ and $y \in [0, 1]$.
1926 Consequently, the Lipschitz continuity of $\bar{f}(\bar{h}, y, t)$ with respect to \bar{h} extends over the entire space
1927 \mathbb{R}^{d_0} .
1928

- 1929 • **Approximation Error Analysis under L^2 Norm.** We first decompose the L^2 approximation
1930 error of \bar{f} into two terms ($\|\bar{h}\|_2 < r_h$ and $\|\bar{h}\|_2 > r_h$):

1931
1932
$$\left\| q(\bar{h}, y, t) - \bar{f}(\bar{h}, y, t) \right\|_{L^2(P_t^h)}$$

$$= \|(q(\bar{h}, y, t) - \bar{f}(\bar{h}, y, t))\mathbb{1}\{\|\bar{h}\|_2 < r_h\}\|_{L^2(P_t^h)} + \|q(\bar{h}, y, t)\mathbb{1}\{\|\bar{h}\|_2 > r_h\}\|_{L^2(P_t^h)}.$$

By selecting $r_h = \mathcal{O}\left(\sqrt{d_0 \log(d_0/t_0) + \log(1/\epsilon)}\right)$ (see Lemma F.3), we bound the second term on the RHS of above expression as:

$$\|g(\bar{h}, y, t)\mathbb{1}\{\|\bar{h}\|_2 > r_h\}\|_{L^2(P_t^h)} \leq \epsilon.$$

For the first term, we bound

$$\begin{aligned} & \|(q(\bar{h}, y, t) - \bar{f}(\bar{h}, y, t))\mathbb{1}\{\|\bar{h}\|_2 < r_h\}\|_{L^2(P_t^h)} \\ & \leq \sqrt{d_0 + d_y} \sup_{h', y, t' \in [0,1]^{d_0} \times [0,1]^{d_y} \times [t_0/T, 1]} \|\bar{f}(h', y, t') - g(h', y, t')\|_\infty \\ & \leq \sqrt{d_0 + d_y} \epsilon. \end{aligned}$$

So we obtain

$$\|q(\bar{h}, y, t) - \bar{f}(\bar{h}, y, t)\|_{L^2(P_t^h)} \leq \left(\sqrt{d_0 + d_y} + 1\right)\epsilon.$$

Substituting ϵ with $\epsilon/2$ gives an approximation error for $\bar{f}(\bar{h}, y, t)$ of $\sqrt{d_0 + d_y}\epsilon$.

Step 2. Approximate $\bar{f}(\bar{h}, y, t)$ with One-Layer Self-Attention. This step is based on the universal approximation of single-layer single-head transformers for compact-supported continuous function in Theorem H.2.

Recall the reshape layer $\tilde{R}(\cdot)$ from Definition 2.3. We use $f(\cdot) := \tilde{R}^{-1} \circ \hat{g}_T \circ \tilde{R}(\cdot)$ to approximate $\bar{f}_t(\cdot) := \bar{f}(\cdot, t)$, where $\hat{g}_T(\cdot) \in \mathcal{T}^{h,s,r} = \{f_2^{(\text{FF})} \circ f^{(\text{SA})} \circ f_1^{(\text{FF})} : \mathbb{R}^{\tilde{d} \times \tilde{L}} \rightarrow \mathbb{R}^{\tilde{d} \times \tilde{L}}\}$.

We first use $\hat{f}_t(\cdot) := \tilde{R}^{-1} \circ \hat{g}_T \circ \tilde{R}(\cdot)$ to approximate the function $\bar{f}_t(\cdot)$ constructed at Step 1 and denote $H = R(\bar{h})$. Using Theorem H.2, we have:

$$\|\bar{f}_t(\bar{h}, y) - \hat{f}_t(\bar{h}, y)\|_{L^2(P_t^h)} = \left(\int_{P_t^h} \|\bar{f}_t(\bar{h}, y) - \hat{f}_t(\bar{h}, y)\|_2^2 dh \right)^{1/2} \quad (\text{F.10})$$

$$= \left(\int_{P_t^h} \|\tilde{R} \circ \bar{f}_t \circ \tilde{R}^{-1}(H) - \tilde{R} \circ \hat{g}_T \circ \tilde{R}^{-1}(H)\|_F^2 dh \right)^{1/2}$$

$$= \left(\int_{P_t^h} \|\tilde{R} \circ \bar{f}_t \circ \tilde{R}^{-1}(H) - \hat{g}_T(H)\|_F^2 dh \right)^{1/2}$$

$$\leq \epsilon. \quad (\text{F.11})$$

Along with Step 1, we obtain

$$\begin{aligned} \|(q(\bar{h}, y, t) - \hat{f}(\bar{h}, y))\|_{L^2(P_t^h)} & \leq \|q(\bar{h}, y, t) - \bar{f}(\bar{h}, y, t)\|_{L^2(P_t^h)} + \|\bar{f}(\bar{h}, y, t) - \hat{g}_T(\bar{h}, y)\|_{L^2(P_t^h)} \\ & \leq \left(1 + \sqrt{d_0 + d_y}\right)\epsilon. \end{aligned}$$

The approximator $s_{\widehat{W}}$ for the score function $\nabla \log p_t(\bar{h}|y)$ is define in (E.2) where $s_{\widehat{W}} = (W_U \hat{f}(U^\top x, y, t) - x)/\sigma_t^2$. The approximation error for such an approximator is

$$\|\nabla \log p_t(\cdot) - s_{\widehat{W}}(\cdot, t)\|_{L^2(P_t)} \leq \frac{1 + \sqrt{d_0 + d_y}}{\sigma_t^2} \epsilon, \quad \text{for all } t \in [t_0, T].$$

Finally, the parameter bounds in the transformer network class satisfy

$$\begin{aligned}
\|W_Q\|_2 &= \|W_K\|_2 = \mathcal{O}\left(\tilde{d} \cdot \epsilon^{-\left(\frac{1}{a}+2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\
\|W_Q\|_{2,\infty} &= \|W_K\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{3}{2}} \cdot \epsilon^{-\left(\frac{1}{a}+2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\
\|W_O\|_2 &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{\frac{1}{a}}\right); \|W_O\|_{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{a}}\right); \\
\|W_V\|_2 &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}}\right); \|W_V\|_{2,\infty} = \mathcal{O}(\tilde{d}); \\
\|W_1\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{a}}\right), \|W_1\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{-\frac{1}{a}}\right); \\
\|W_2\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{a}}\right); \|W_2\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{-\frac{1}{a}}\right); \\
\|E^\top\|_{2,\infty} &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \tilde{L}^{\frac{3}{2}}\right).
\end{aligned}$$

We refer to [Appendix H.2](#) for the calculation of the hyperparameters configuration of this network. This completes the proof. \square

F.5 PROOF OF SCORE ESTIMATION ([THEOREM F.2](#))

Lemma F.6 (Lemma 15 of [\(Chen et al., 2023c\)](#)). Let \mathcal{G} be a bounded function class, i.e., there exists a constant b such that any function $g \in \mathcal{G} : \mathbb{R}^{d_0} \mapsto [0, b]$. Let $z_1, z_2, \dots, z_n \in \mathbb{R}^{d_0}$ be i.i.d. random variables. For any $\delta \in (0, 1)$, $a \leq 1$, and $c > 0$, we have

$$\begin{aligned}
P\left(\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n g(z_i) - (1+a)\mathbb{E}[g(z)] > \frac{(1+3/a)B}{3n} \log \frac{\mathcal{N}(c, \mathcal{G}, \|\cdot\|_\infty)}{\delta} + (2+a)c\right) &\leq \delta, \\
P\left(\sup_{g \in \mathcal{G}} \mathbb{E}[g(z)] - \frac{1+a}{n} \sum_{i=1}^n g(z_i) > \frac{(1+6/a)B}{3n} \log \frac{\mathcal{N}(c, \mathcal{G}, \|\cdot\|_\infty)}{\delta} + (2+a)c\right) &\leq \delta.
\end{aligned}$$

Main Proof of [Theorem F.2](#). Now we are ready to state the main proof.

Proof of [Theorem F.2](#). Our proof is built on [\(Chen et al., 2023c, Appendix B.2\)](#).

Recall that the empirical score-matching loss is

$$\mathcal{L}(s_{\widehat{W}}) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; s_{\widehat{W}}), \quad (\text{F.12})$$

with the loss function ℓ for a data sample (x, y) is defined as

$$\ell(x, y, s_{\widehat{W}}) = \int_{t_0}^T \frac{1}{T-t_0} \mathbb{E}_{(x_t|x_0=x, \tau)} \left[\|s(x_t, \tau y, t) - \nabla \log \phi_t(x_t|x_0)\|_2^2 \right] dt.$$

We organize the proof into the following three steps:

- **Step 1. Decomposing $\mathcal{L}(s_{\widehat{W}})$:** We first decompose \mathcal{L} into three terms (A) , (B) , and (C) .
- **Step 2. Bounding Each Term:** We then bound three terms separately using some helper from [Lemma F.2](#) and [Lemma F.6](#).
- **Step 3. Putting All Together:** Finally, we combine the above bounds and substitute the covering number of $\mathcal{S}(C_x)$ from [Lemma K.3](#).

2052 • **Step 1. Decomposing $\mathcal{L}(s_{\widehat{W}})$:**

2053 Following (Chen et al., 2023c, Appendix B.2), for any $a \in (0, 1)$, we have:

$$2054 \mathcal{L}(s_{\widehat{W}}) \\ 2055 \leq \underbrace{\mathcal{L}^{\text{trunc}}(s_{\widehat{W}}) - (1+a)\widehat{\mathcal{L}}^{\text{trunc}}(s_{\widehat{W}})}_{(A)} + \underbrace{\mathcal{L}(s_{\widehat{W}}) - \mathcal{L}^{\text{trunc}}(s_{\widehat{W}})}_{(B)} + (1+a) \underbrace{\inf_{s_W \in \mathcal{T}_{\bar{R}}^{h,s,r}} \widehat{\mathcal{L}}(s_W)}_{(C)}.$$

2060 where

$$2061 \mathcal{L}^{\text{trunc}}(s_{\widehat{W}}) := \mathbb{E}_{x \sim P_0} [\ell(x, \tau y, s_{\widehat{W}}) \mathbb{1}\{\|x\|_2 \leq r_x\}], \quad r_x > B,$$

2062 We denote

$$2063 \eta := 4C_{\mathcal{T}}(C_{\mathcal{T}} + r_x)(r_x/d_x)^{d_x-2} \cdot \exp(-r_x^2/\sigma_t^2)/t_0(T-t_0), \\ 2064 r_x := \mathcal{O}\left(\sqrt{d_0 \log d_0 + \log C_{\mathcal{T}} + \log(n/\bar{\delta})}\right).$$

2070 • **Step 2. Bounding Each Term:** We bound (A), (B), and (C) term separately using some helper from Lemma F.2 and Lemma F.6.

2071 **Bounding term (A).** For any $\bar{\delta} > 0$, following (Chen et al., 2023c, Appendix B.2) and applying Lemma F.6, we have the following for term (A) with probability $1 - \bar{\delta}$,

$$2072 (A) = \mathcal{O}\left(\frac{(1+3/a)(C_{\mathcal{T}}^2 + r_x^2)}{nt_0(T-t_0)} \log \frac{\mathcal{N}\left(\frac{(T-t_0)(\epsilon_c - \eta)}{(C_{\mathcal{T}} + r_x)\log(T/t_0)}, \mathcal{T}^{h,s,r}, \|\cdot\|_2\right)}{\bar{\delta}} + (2+a)c\right),$$

2073 where $c \leq 0$ is a constant, and $\epsilon_c > 0$ is another constant to be determined later.

2074 By setting $\epsilon_c = \log(2/(nt_0(T-t_0)))$, then we have

$$2075 (A) = \mathcal{O}\left(\frac{(1+3/a)(C_{\mathcal{T}}^2 + r_x^2)}{nt_0(T-t_0)} \log \frac{\mathcal{N}\left((n(C_{\mathcal{T}} + r_x)t_0 \log(T/t_0))^{-1}, \mathcal{T}^{h,s,r}, \|\cdot\|_2\right)}{\bar{\delta}} + \frac{1}{n}\right), \\ 2076 (F.13)$$

2077 with probability $1 - \bar{\delta}$.

2078 **Bounding term (B).** Following (Chen et al., 2023c, Appendix B.2) and applying Lemma F.2, we has the following bound for term (B):

$$2079 (B) = \mathcal{O}\left(\frac{1}{t_0(T-t_0)} C_{\mathcal{T}}^2 r_x^{d_0} \frac{2^{-2/d_0+2} d_0}{\Gamma(d_0/2+1)} \exp(-C_2 r_x^2/2)\right). \quad (F.14)$$

2080 **Bounding term (C).** In Theorem F.1, we approximate the score function with the network \widehat{s}_W for any $\epsilon > 0$. We decompose the term (C) into statistical error (C_1) and approximation error (C_2):

$$2081 (C) \leq \underbrace{\widehat{\mathcal{L}}(\widehat{s}_W) - (1+a)\mathcal{L}^{\text{trunc}}(\widehat{s}_W)}_{(C_1)} + (1+a) \underbrace{\mathcal{L}^{\text{trunc}}(\widehat{s}_W)}_{(C_2)}.$$

Following (Chen et al., 2023c, Appendix B.2) and applying Lemma F.2 and Lemma F.6, we have the following bound for term (C_1) :

$$(C_1) = \widehat{\mathcal{L}}^{\text{trunc}}(\widehat{s}_W) - (1+a)\mathcal{L}^{\text{trunc}}(\widehat{s}_W) = \mathcal{O}\left(\frac{(1+6/a)(C_{\mathcal{T}}^2 + r_x^2)}{nt_0(T-t_0)} \log \frac{1}{\delta}\right),$$

with probability $1 - \delta$.

Finally, for the term (C_2) we use Theorem F.1 for score function approximation of $\mathcal{L}(\widehat{s}_W)$:

$$(C_2) = \mathcal{O}\left(\frac{d_0 + d_y}{t_0(T-t_0)} \epsilon^2\right) + (\text{const.}).$$

This give us the bound for term $(C) \leq (C_1) + (1+a)(C_2)$ as

$$(C) \leq \mathcal{O}\left(\frac{(1+6/a)(C_{\mathcal{T}}^2 + r_x^2)}{nt_0(T-t_0)} \log \frac{1}{\delta} + \frac{d_0 + d_y}{t_0(T-t_0)} \epsilon^2\right) + (\text{const.}). \quad (\text{F.15})$$

- **Step 3. Putting All Together:** In the final steps, we combine three terms and substitute the covering number to get the score estimation bound for latent DiT.

Combining (A), (B) and (C). Following (Chen et al., 2023c, Appendix B.2), we set $a = \epsilon^2$ and get the overall bound:

$$\begin{aligned} & \frac{1}{T-t_0} \int_{t_0}^T \|s_{\widehat{W}}(\cdot, t) - \nabla \log p_t(\cdot)\|_{L^2(P_t)}^2 dt \\ &= \mathcal{O}\left(\frac{(C_{\mathcal{T}}^2 + r_x^2)}{\epsilon^2 nt_0(T-t_0)} \log \frac{\mathcal{N}((n(C_{\mathcal{T}} + r_x)t_0 \log(T/t_0))^{-1}, \mathcal{S}_{\mathcal{T}^{h,s,r}}, \|\cdot\|_2)}{\delta} + \frac{1}{n} + \frac{d_0 + d_y}{t_0(T-t_0)} \epsilon^2\right), \end{aligned} \quad (\text{F.16})$$

with probability $1 - 3\delta$.

Before we move on to the covering number of $\mathcal{T}_{\widehat{R}}^{h,s,r}$, we first compute the Lipschitz upper bound $L_{\mathcal{T}}$ and model output bound $C_{\mathcal{T}}$.

Lipschitz Upper Bound $L_{\mathcal{T}}$ and Model Output Bound $C_{\mathcal{T}}$. We then compute the Lipschitz upper bound $L_{\mathcal{T}}$ for the transformer. We denote $\bar{f}_{t,R}(\cdot) = \widetilde{R} \circ \widehat{g}_t \circ \widetilde{R}^{-1}(\cdot)$ and $H = (\widetilde{R}(\bar{h}), y)$.

We get the Lipschitz upper bound for $\widehat{f}_{\mathcal{T}} \in \mathcal{T}_{\widehat{R}}^{h,s,r}$:

$$\begin{aligned} \left\| \widehat{f}_{\mathcal{T}}(H_1) - \widehat{f}_{\mathcal{T}}(H_2) \right\|_F &\leq \left\| \widehat{f}_{\mathcal{T}}(H_1) - \bar{f}_{t,\widetilde{R}}(H_1) \right\|_F + \left\| \bar{f}_{t,\widetilde{R}}(H_1) - \bar{f}_{t,\widetilde{R}}(H_2) \right\|_F \\ &\quad + \left\| \bar{f}_{t,\widetilde{R}}(H_2) - \widehat{f}_{\mathcal{T}}(H_2) \right\|_F \\ &\leq 2\epsilon + \left\| \bar{f}_{t,\widetilde{R}}(H_1) - \bar{f}_{t,\widetilde{R}}(H_2) \right\|_F \quad (\text{By (F.10)}) \\ &\leq 2\epsilon + 10(d_0 + d_y)L_{s_+} \|H_1 - H_2\|_F. \quad (\text{By (F.9)}) \end{aligned}$$

Then we get the upper bound of Lipschitzness of $\mathcal{T}_{\widehat{R}}^{h,s,r}$:

$$L_{\mathcal{T}} = \mathcal{O}((d_0 + d_y)L_{s_+}). \quad (\text{F.17})$$

Next, we compute the model output bound for $\mathcal{T}_{\widehat{R}}^{h,s,r}$. For the output of the constructed transformer $\widehat{f}_{\mathcal{T}} \in \mathcal{T}^{h,s,r}$, according to (H.20), the output of the network is lower bounded by $\mathcal{O}(1)$. Thus with the Lipschitz upper bound $L_{\mathcal{T}} = \mathcal{O}((d_0 + d_y)L_{s_+})$, we have $\|\widehat{f}_{\mathcal{T}}(H)\|_F = \mathcal{O}((d_0 + d_y)L_{s_+}r_h)$,

where $\|H\|_F \leq r_h$. With $r_h = c(\sqrt{d_0 \log(d_0/t_0)} + \log(1/\epsilon))$, we obtain

$$C_{\mathcal{T}} = \mathcal{O}\left((d_0 + d_y)L_{s_+} \cdot \sqrt{d_0 \log(d_0/t_0)} + \log(1/\epsilon)\right). \quad (\text{F.18})$$

Covering Number of $\mathcal{T}_{\tilde{R}}^{h,s,r}$. The next step is to calculate the covering number of $\mathcal{T}_{\tilde{R}}^{h,s,r}$. In particular, $\mathcal{T}_{\tilde{R}}^{h,s,r}$ consists of two components: (i) Matrix W_U with orthonormal columns; (ii) Network function $g_{\mathcal{T}}$. Suppose we have W_{U1}, W_{U2} and g_1, g_2 such that $\|W_{U1} - W_{U2}\|_F \leq \delta_1$ and $\sup_{\|x\|_2 \leq 3r_x + \sqrt{d_x \log d_x}, y \in [0,1], t \in [t_0, T]} \|g_1(x, y, t) - g_2(x, y, t)\|_2 \leq \delta_2$, where $g_1 = \tilde{R}^{-1} \circ g_{\mathcal{T}1} \circ \tilde{R}$ and $g_2 = \tilde{R}^{-1} \circ g_{\mathcal{T}2} \circ \tilde{R}$. Then we evaluate

$$\begin{aligned} & \sup_{\|x\|_2 \leq 3r_x + \sqrt{d_x \log d_x}, y \in [0,1], t \in [t_0, T]} \|s_{W_{U1}, g_{\mathcal{T}1}}(x, y, t) - s_{W_{U2}, g_{\mathcal{T}2}}(x, y, t)\|_2 \\ &= \frac{1}{\sigma_t^2} \sup_{\|x\|_2 \leq 3r_x + \sqrt{d_x \log d_x}, y \in [0,1], t \in [t_0, T]} \|W_{U1}g_1(W_{U1}^\top x, y, t) - W_{U2}g_2(W_{U2}^\top x, y, t)\|_2 \\ &\leq \frac{1}{\sigma_t^2} \sup_{\|x\|_2 \leq 3r_x + \sqrt{d_x \log d_x}, y \in [0,1], t \in [t_0, T]} \left(\underbrace{\|W_{U1}g_1(W_{U1}^\top x, y, t) - W_{U1}g_1(W_{U2}^\top x, y, t)\|_2}_{1^{\text{st}} \text{ term}} \right. \\ & \quad \left. + \underbrace{\|W_{U1}g_1(W_{U2}^\top x, y, t) - W_{U1}g_2(W_{U2}^\top x, y, t)\|_2}_{2^{\text{nd}} \text{ term}} + \underbrace{\|W_{U1}g_2(W_{U2}^\top x, y, t) - W_{U2}g_2(W_{U2}^\top x, y, t)\|_2}_{3^{\text{rd}} \text{ term}} \right) \\ &\leq \frac{1}{\sigma_t^2} \left(\underbrace{L_{\mathcal{T}}\delta_1\sqrt{d_0}(3r_x + \sqrt{d_x \log d_x})}_{1^{\text{st}} \text{ term}} + \underbrace{\delta_2}_{2^{\text{nd}} \text{ term}} + \underbrace{\delta_1}_{3^{\text{rd}} \text{ term}} \right), \end{aligned} \quad (\text{F.19})$$

where $L_{\mathcal{T}}$ upper bounds the Lipschitz constant of $g_{\mathcal{T}}$ (see (F.17)).

For the set $\{W_B \in \mathbb{R}^{d_x \times d_0} : \|W_B\|_2 \leq 1\}$, its δ_1 -covering number is $(1 + 2\sqrt{d_0}/\delta_1)^{d_x d_0}$ (Chen et al., 2023c, Lemma 8). The δ_2 -covering number of f needs further discussion as there is a reshaping process in our network. For the input reshaped from $\bar{h} \in \mathbb{R}^{d_0}$ to $H \in \mathbb{R}^{\tilde{d} \times \tilde{L}}$, we have

$$\|\bar{h}\|_2 \leq r_x \iff \|H\|_F \leq r_x,$$

Thus we have

$$\begin{aligned} & \sup_{\|\bar{h}\|_2 \leq 3r_x + \sqrt{D \log D}, y \in [0,1], t \in [t_0, T]} \|g_1(\bar{h}, y, t) - g_2(\bar{h}, y, t)\|_2 \leq \delta_2, \\ & \iff \sup_{\|H\|_F \leq 3r_x + \sqrt{D \log D}, y \in [0,1], t \in [t_0, T]} \|g_{\mathcal{T}1}(H) - g_{\mathcal{T}2}(H)\|_2 \leq \delta_2. \end{aligned}$$

Next we follow the covering number property for sequence-to-sequence transformer $\mathcal{T}_{\tilde{R}}^{h,s,r}$, i.e., Lemma K.2 and get the following δ_2 -covering number

$$\log \mathcal{N}\left(\epsilon_c, \mathcal{T}_{\tilde{R}}^{h,s,r}, \|\cdot\|_2\right) \quad (\text{F.20})$$

$$\leq \frac{\log(nL)}{\epsilon_c^2} \cdot \alpha^2 \left(\underbrace{\left((C_F)^2 C_{OV}^{2,\infty}\right)^{\frac{2}{3}}}_{1^{\text{st}} \text{ term}} + \underbrace{(d + d_y)^{\frac{2}{3}} (C_F^{2,\infty})^{\frac{4}{3}}}_{2^{\text{nd}} \text{ term}} + \underbrace{(d + d_y)^{\frac{2}{3}} (2(C_F)^2 C_{OV} C_{KQ}^{2,\infty})^{\frac{2}{3}}}_{3^{\text{rd}} \text{ term}} \right)^3, \quad (\text{F.21})$$

where

$$\alpha := \prod_{j < i} (C_F)^2 C_{OV} (1 + 4C_{KQ}) (C_X + C_E).$$

Recall that from the network configuration in [Theorem F.1](#), we have the following bound:

$$\begin{aligned} \|W_Q\|_2 &= \|W_K\|_2 = \mathcal{O}\left(\tilde{d} \cdot \epsilon^{-\left(\frac{1}{\tilde{d}} + 2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\ \|W_Q\|_{2,\infty} &= \|W_K\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{3}{2}} \cdot \epsilon^{-\left(\frac{1}{\tilde{d}} + 2\tilde{L}\right)} (\log \tilde{L})^{\frac{1}{2}}\right); \\ \|W_O\|_2 &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \epsilon^{\frac{1}{\tilde{d}}}\right); \|W_O\|_{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{\tilde{d}}}\right); \\ \|W_V\|_2 &= \mathcal{O}(\tilde{d}^{\frac{1}{2}}); \|W_V\|_{2,\infty} = \mathcal{O}(\tilde{d}); \\ \|W_1\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{\tilde{d}}}\right), \|W_1\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{3}{2}} \epsilon^{-\frac{1}{\tilde{d}}}\right); \\ \|W_2\|_2 &= \mathcal{O}\left(\tilde{d} \epsilon^{-\frac{1}{\tilde{d}}}\right); \|W_2\|_{2,\infty} = \mathcal{O}\left(\tilde{d}^{\frac{3}{2}} \epsilon^{-\frac{1}{\tilde{d}}}\right); \\ \|E^\top\|_{2,\infty} &= \mathcal{O}\left(\tilde{d}^{\frac{1}{2}} \tilde{L}^{\frac{3}{2}}\right). \end{aligned}$$

Note that $W_{K,Q} = W_Q W_K^\top$ and $W_{O,V} = W_O W_V^\top$. Combining every component and substitute into [\(F.20\)](#), we have three respective terms bounded as

$$\begin{aligned} \text{1st term} &= \mathcal{O}\left(\tilde{d}^2 \epsilon^{-2/(3\tilde{d})}\right), \\ \text{2nd term} &= \mathcal{O}\left((d_0 + d_y)^{2/3} \tilde{d}^{2/3} \epsilon^{-4/(3\tilde{d})}\right), \\ \text{3rd term} &= \mathcal{O}\left((d_0 + d_y)^{2/3} \cdot (\log \tilde{L})^{2/3} \cdot \tilde{d}^4 \cdot \epsilon^{(-2/3)(3\tilde{d}+4\tilde{L})}\right). \end{aligned}$$

Apparently the 3rd term dominates the other two. For the α^2 term, we write

$$\alpha^2 = \mathcal{O}\left(\tilde{d}^{10} \epsilon^{-2(3\tilde{d}+4\tilde{L})} (\log \tilde{L}) C'_x\right),$$

where $C'_x = (C_x + (d_0 + d_y)^{3/2})^2$.

Combining the above bound we get the log-covering number of \mathcal{T}_2 as

$$\log \mathcal{N}\left(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2\right) \lesssim \mathcal{O}\left(\frac{\log(n\tilde{L}) \log^3(\tilde{L})}{\epsilon_c^2} \tilde{d}^{22} (d_0 + d_y)^2 \epsilon^{-4(3\tilde{d}+4\tilde{L})} C_x^2\right). \quad (\text{F.22})$$

Substituting the log-covering number of $\mathcal{T}_R^{h,s,r}$ into [\(F.16\)](#), we have

$$\begin{aligned} &\frac{1}{T-t_0} \int_{t_0}^T \|s_{\widehat{W}}(\cdot, t) - \nabla \log p_t(\cdot)\|_{L^2(P_t)}^2 dt \\ &= \mathcal{O}\left(\frac{(C_T^2 + \log(n/\delta))}{\epsilon^2 n t_0 (T-t_0)} \left(\frac{\log(n\tilde{L}) \log^3(\tilde{L})}{(T-t_0)n^2} \tilde{d}^{22} (d_0 + d_y)^2 \epsilon^{-4(3\tilde{d}+4\tilde{L})} C_x^2\right) + \frac{1}{n} + \frac{d_0 + d_y}{t_0(T-t_0)} \epsilon^2\right) \\ &\hspace{15em} (\text{By (F.16)}) \end{aligned}$$

$$\begin{aligned}
&= \mathcal{O} \left(\frac{((\tilde{d} + d_0)^2 L_{s^+}^2 (d_0 \log(d_0/t_0) + \log(1/\epsilon)) + \log(n/\bar{\delta}))}{\epsilon^2 n t_0 (T - t_0)} \left(\frac{\log(n\tilde{L}) \log^3(\tilde{L})}{(T - t_0) n^2} \tilde{d}^{22} (\tilde{d} + d_y)^2 \epsilon^{-4(3/\tilde{d} + 4\tilde{L})} C_x^2 \right) \right. \\
&\quad \left. + \frac{d_0 + d_y}{t_0 (T - t_0)} \epsilon^2 \right). \tag{By (F.17) and (F.18)}
\end{aligned}$$

Balancing Error Terms. To balance the error term, we set $\epsilon = n^{-3/4(1+3/\tilde{d}+4\tilde{L})}$. Also setting $\bar{\delta} = 1/3n$ then we have

$$\frac{1}{T - t_0} \int_{t_0}^T \|\widehat{s_{\tilde{W}}}(\cdot, t) - \nabla \log p_t(\cdot)\|_{L^2(P_t)}^2 dt = \mathcal{O} \left(\frac{\tilde{d}^{22} (\tilde{d} + d_0)^2 (\tilde{d} + d_y)^2}{t_0^2} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \log^3 \tilde{L} \log^3 n \right) \tag{F.23}$$

with probability of $1 - \frac{1}{n}$.

This completes the proof. \square

F.6 PROOF OF DISTRIBUTION ESTIMATION (THEOREM F.3)

Our proof is built on [Chen et al. \(2023c, Appendix C\)](#). The main difference between our work and [Chen et al. \(2023c\)](#) is our score estimation error from [Theorem F.2](#). This is based on our universal approximation of transformers in [Corollary H.2.1](#). Consequently, only the subspace error and the total variation distance differ from [Chen et al. \(2023c, Theorem 3\)](#).

Proof Sketch of (i). We show that if the orthogonal score increases significantly, the mismatch between the column span of U and W_U will be greatly amplified. Therefore, an accurate score network estimator forces U and W_U to align with each other.

Proof Sketch of (ii). We conduct the proof via 2 steps:

- **Step 1: Total Variation Distance Bound.** We obtain the discrete result from the continuous-time generated distribution \widehat{P}_{t_0} by adding discretization error ([Chen et al., 2023c, Lemma 4](#)). It suffices to bound the divergence between the following two stochastic processes:

- For the ground-truth backward process, consider $h_t^{\leftarrow} = B^\top y_t$ and the following SDE:

$$dh_t^{\leftarrow} = \left[\frac{1}{2} h_t^{\leftarrow} + \nabla \log p^h T - t(h_t^{\leftarrow}) \right] dt + d\bar{U}_t^h.$$

Denote the marginal distribution of the ground-truth process as $P_{t_0}^h$.

- For the learned process, consider $\tilde{h}_t^{\leftarrow, r}$ and the following SDE:

$$d\tilde{h}_t^{\leftarrow, r} = \left[\frac{1}{2} \tilde{h}_t^{\leftarrow, r} + \tilde{s}_{f, M}^h(\tilde{h}_t^{\leftarrow, r}, T - t) \right] dt + d\bar{U}_t^h,$$

where $\tilde{s}_{f, M}^h(z, t) := [M^\top f(Mz, t) - z]/\sigma_t^2$ and M is an orthogonal matrix. Following the notation in ([Chen et al., 2023c](#)), we use $(W_U M)^\top \widehat{P}_{t_0}$ to denote the marginal distribution of \widehat{P}_{t_0} . We first calculate the latent score matching error, i.e., the error between $\nabla \log p_t^h(h, y)$ and $\tilde{s}_{M, f}^h(h, y, t)$. Then, we adopt Girsanov's Theorem ([Chen et al., 2022](#)) and bound the difference in the KL divergence of the above two processes to derive the score-matching error bound.

Proof Sketch of (iii). We derive item (iii) by solving the orthogonal backward process of the diffusion model.

Definition F.1. For later convenience, let us define $\xi(n, t_0, \tilde{d}, \tilde{L}) := \frac{1}{t_0^2} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \log^3 n$.

Here we include a few auxiliary lemmas from [Chen et al. \(2023c\)](#) without proofs. Recall the definition of Lipschitz norm: for a given function f , $\|f(\cdot)\|_{Lip} = \sup_{x \neq y} (\|f(x) - f(y)\|_2 / \|x - y\|_2)$.

Lemma F.7 (Lemma 3 of [Chen et al. \(2023c\)](#)). Assume that the following holds

$$\mathbb{E}_{h \sim P_h} \|\nabla \log p_h(h|y)\|_2^2 \leq C_{sh}, \quad \lambda_{\min} \mathbb{E}_{h \sim P_h} [hh^\top] \geq c_0, \quad \mathbb{E}_{h \sim P_h} \|h\|_2^2 \leq C_h,$$

where λ_{\min} denotes the smallest eigenvalue. We denote

$$\bar{\mathbb{E}}[\phi(\cdot, t)] = \int_{t_0}^T \frac{1}{\sigma_t^4} \mathbb{E}_{x \sim P_t} [\phi(\cdot, t)] dt.$$

We set $t_0 \leq \min\{2 \log(d_0/C_{sh}), 1, 2 \log(c_0), c_0\}$ and $T \geq \max\{2 \log(C_h/d_0), 1\}$. Suppose we have

$$\bar{\mathbb{E}} \|W_B f(W_B^\top x, y, t) - Uq(B^\top x, y, t)\|_2^2 \leq \epsilon.$$

Then we have

$$\|W_U W_U^\top - U U^\top\|_F^2 = \mathcal{O}(\epsilon t_0 / c_0),$$

and there exists an orthogonal matrix $M \in \mathbb{R}^{d_0 \times d_0}$, such that:

$$\begin{aligned} & \bar{\mathbb{E}} \|M^\top f(Mh, y, t) - q(h, y, t)\|_2^2 \\ &= \epsilon \cdot \mathcal{O} \left(1 + \frac{t_0}{c_0} \left[(T - \log t_0) d_0 \cdot \max_t \|f(\cdot, t)\|_{Lip}^2 + C_{sh} \right] + \frac{\max_t \|f(\cdot, t)\|_{Lip}^2 \cdot C_h}{c_0} \right). \end{aligned}$$

Lemma F.8 (Lemma 4 of [Chen et al. \(2023c\)](#)). Assume that P_h is sub-Gaussian, $f(h, y, t)$ and $\nabla \log p_t^h(h|y)$ are Lipschitz in both h, y and t . Assume we have the latent score matching error-bound

$$\int_{t_0}^T \mathbb{E}_{h \sim P_t^h} \|\tilde{s}_{M,f}^h(h_t, y, t) - \nabla \log p_t^h(h_t|y)\|_2^2 dt \leq \epsilon_{\text{latent}} (T - t_0).$$

Then we have the following latent distribution estimation error for the undiscretized backward SDE

$$\text{TV} \left(P_{t_0}^h, \hat{P}_{t_0}^h \right) \lesssim \sqrt{\epsilon_{\text{latent}} (T - t_0)} + \sqrt{\text{KL}(P_h \| N(0, I_{d_0}))} \cdot \exp(-T).$$

Furthermore, we have the following latent distribution estimation error for the discretized backward SDE

$$\text{TV} \left(P_{t_0}^h, \hat{P}_{t_0}^{h, \text{dis}} \right) \lesssim \sqrt{\epsilon_{\text{latent}} (T - t_0)} + \sqrt{\text{KL}(P_h \| N(0, I_{d_0}))} \cdot \exp(-T) + \sqrt{\epsilon_{\text{dis}} (T - t_0)},$$

where

$$\begin{aligned} \epsilon_{\text{dis}} &= \left(\frac{\max_h \|f(h, y, \cdot)\|_{Lip}}{\sigma(t_0)} + \frac{\max_{h,t} \|f(h, y, t)\|_2}{t_0^2} \right)^2 \eta^2 \\ &+ \left(\frac{\max_t \|f(\cdot, y, t)\|_{Lip}}{\sigma(t_0)} \right)^2 \eta^2 \max \left\{ \mathbb{E} \|h_0\|^2, d_0 \right\} + \eta d_0, \end{aligned}$$

and η is the step size in the backward process.

Lemma F.9 (Lemma 6 of [Chen et al. \(2023c\)](#)). Consider the following discretized SDE with step size μ satisfying $T - t_0 = K_T \mu$

$$dy_t = \left[\frac{1}{2} - \frac{1}{\sigma(T - k\mu)} \right] y_{k\mu} dt + dU_t, \text{ for } t \in [k\mu, (k+1)\mu),$$

where $Y_0 \sim N(0, I)$. Then when $T > 1$ and $t_0 + \mu \leq 1$, we have $Y_{T-t_0} \sim N(0, \sigma^2 I)$ with $\sigma^2 \leq e(t_0 + \mu)$.

Lemma F.10 (Lemma 10 in [Chen et al. \(2023c\)](#)). Assume that $\nabla \log p_h(h|y)$ is L_h -Lipschitz. Then we have $\mathbb{E}_{h \sim P_h} \|\nabla \log p_h(h|y)\|_2^2 \leq d_0 L_h$.

Main Proof of Theorem F.3. Now we are ready to state the main proof.

Proof of Theorem F.3. Recall that in (F.23), we have

$$\xi(n, t_0, \tilde{d}, \tilde{L}) := \frac{1}{t_0^2} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \log^3 L \log^3 n.$$

• **Proof of (i).** With [Lemma F.7](#), we replace ϵ to be $\epsilon(T - t_0)^2$ and we set $C_{sh} = L_h d_0$ by [Lemma F.10](#), we have

$$\|W_U W_U^\top - U U^\top\|_F^2 = \mathcal{O}\left(\frac{t_0^2 \xi(n, t_0, \tilde{d}, \tilde{L})}{c_0}\right).$$

We substitute the score estimation error in [Theorem F.2](#) and $T = \mathcal{O}(\log n)$ into the bound above, we deduce

$$\|W_U W_U^\top - U U^\top\|_F^2 = \tilde{\mathcal{O}}\left(\frac{1}{c_0} n^{\frac{-3}{2(1+3/\tilde{d}+4\tilde{L})}} \cdot \log^3 n\right).$$

We note that $\log n$ is great enough to make T satisfies $T \geq \max\{\log(C_h/d_0 + 1), 1\}$ where $C_h \geq \mathbb{E}_{h \sim P_h} \|h\|_2^2$.

• **Proof of (ii).** [Lemma F.7](#) and [Lemma F.10](#) imply that

$$\mathbb{E} \|M^\top f(Mh, y, t) - q(h, y, t)\|_2^2 = \mathcal{O}(\epsilon_{\text{latent}}(T - t_0)),$$

where

$$\epsilon_{\text{latent}} = \epsilon \cdot \mathcal{O}\left(\frac{t_0}{c_0} \left[(T - \log t_0) d_0 \cdot L_{s_+}^2 + d_0 L_h \right] + \frac{L_{s_+}^2 \cdot C_h}{c_0}\right).$$

Through the algebra calculation, we get

$$\begin{aligned} \mathbb{E} \|M^\top f(Mh, y, t) - q(h, y, t)\|_2^2 &= \int_{t_0}^T \mathbb{E}_{h \sim P_t^h} \left\| \frac{U^\top f(Uh, y, t) - h}{\sigma_t^2} - \nabla \log p_t^h(h|y) \right\|_2^2 dt \\ &\leq \epsilon_{\text{latent}}(T - t_0). \end{aligned}$$

With ϵ_{latent} and [Lemma F.8](#), we obtain

$$\text{TV}(P_{t_0}^h, (W_U M)^\top \hat{P}_{t_0}^{\text{dis}})$$

$$\begin{aligned}
&\lesssim \sqrt{\epsilon_{\text{latent}}(T - t_0)} + \sqrt{\text{KL}(P_h \| N(0, I_{d_0}))} \exp(-T) + \sqrt{\epsilon_{\text{dis}}(T - t_0)} \\
&= \tilde{\mathcal{O}} \left(\frac{1}{t_0 \sqrt{c_0}} n^{\frac{-3}{2(1+3/d+4L)}} \cdot \log^3 n + \frac{1}{n} + \mu \frac{\sqrt{d_0^2 \log d_0}}{t_0^2} + \sqrt{\mu} \sqrt{d_0} \right).
\end{aligned}$$

2434

2435

2436

2437

As we choose time step $\mu = \mathcal{O} \left(t_0^2 / d_0 \sqrt{\log d_0} n^{\frac{-3}{4(1+3/d+4L)}} \right)$, we obtain

2438

2439

2440

$$\text{TV}(P_{t_0}^h, (W_U M)_\#^\top \hat{P}_{t_0}^{\text{dis}}) = \tilde{\mathcal{O}} \left(\frac{1}{t_0 \sqrt{c_0}} n^{\frac{-3}{2(1+3/d+4L)}} \cdot \log^3 n \right).$$

2441

2442

By definition, $\hat{P}_{t_0}^{h, \text{dis}} = (UW_B)_\#^\top \hat{P}_{t_0}^{\text{dis}}$. This completes the proof of the total variation distance.

2443

2444

2445

- **Proof of (iii).** We apply [Lemma F.9](#) due to our score decomposition. With the marginal distribution at time $T - t_0$ and observing $\mu \ll t_0$, we obtain the last property.

2446

This completes the proof. □

2447

2448

2449

2450

2451

2452

2453

2454

2455

2456

2457

2458

2459

2460

2461

2462

2463

2464

2465

2466

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

G SUPPLEMENTARY THEORETICAL BACKGROUND

In this section, we provide an overview of the conditional diffusion model and classifier guidance in Appendix G.1 and classifier-free guidance in Appendix G.2.

G.1 CONDITIONAL DIFFUSION PROCESS

Conditional diffusion models use the conditional information (guidance) y to generate samples from conditional data distribution $P(\cdot|y = \text{guidance})$. Depending on the model’s objective, the guidance is either a label for generating categorical images, a text prompt for generating images from input sentences, or an image region for tasks like image editing and restoration. Throughout this paper, we coin diffusion models with label guidance y as conditional diffusion models (CDMs). Practically, implement a conditional diffusion model characterized as classifier and classifier-free guidance. The classifier guidance diffusion model combines the unconditional score function with the gradient of an external classifier trained on corrupted data. On the other hand, classifier-free guidance integrates the conditional and unconditional score function by randomly ignoring y with mask signal (see (G.6)). In this paper, we focus on the latter approach.

Specifically, we consider data $x \in \mathbb{R}^{d_x}$ and label $y \in \mathbb{R}^{d_y}$ with initial conditional distribution $P(x|y)$. The diffusion process (forward Ornstein–Uhlenbeck process) is characterized by:

$$dX_t = -\frac{1}{2}X_t dt + dW_t \quad \text{with} \quad X_0 \sim P(x|y), \quad (\text{G.1})$$

where W_t is a Wiener process. The distribution at any finite time t is denoted by $P_t(x|y)$, and X_∞ follows standard Gaussian distribution. Up to a sufficiently large terminating time T , we generate samples by the reverse process:

$$dX_t^\leftarrow = \left[\frac{1}{2}X_t^\leftarrow + \nabla \log p_{T-t}(X_t^\leftarrow | y) \right] dt + d\bar{W}_t \quad \text{with} \quad X_0^\leftarrow \sim P_T(x|y), \quad (\text{G.2})$$

where the term $\nabla \log p_{T-t}(X_t^\leftarrow | y)$ represents the conditional score function. We have $X_t | X_0 \sim N(\alpha_t X_0, \sigma_t^2 I)$ with $\alpha_t = e^{-t/2}$ and $\sigma_t^2 = 1 - e^{-t}$.

We use a score network \hat{s} to estimate the conditional score function $\nabla \log p_t(x|y)$, and the quadratic loss of the conditional diffusion model is given by

$$\hat{s} := \underset{s \in \mathcal{T}_R^{h,s,r}}{\operatorname{argmin}} \mathbb{E}_t \left[\mathbb{E}_{(x_0, y)} \left[\mathbb{E}_{(x' \sim x' | x_0)} \left[\|s(x', y, t) - \nabla_{x'} \log p_t(x' | x_0)\|_2^2 \right] \right] \right], \quad (\text{G.3})$$

where $t \sim \operatorname{Unif}(t_0, T)$.

With the estimate score network \hat{s} in (G.3), we generate the conditional sample in the backward process as follows:

$$d\tilde{X}_t^\leftarrow = \left[\frac{1}{2}\tilde{X}_t^\leftarrow + \hat{s}(\tilde{X}_t^\leftarrow, y, T - t) \right] dt + d\bar{W}_t \quad \text{with} \quad \tilde{X}_0^\leftarrow \sim N(0, I_d). \quad (\text{G.4})$$

Classifier guidance (Song et al., 2021; Dhariwal and Nichol, 2021) and classifier-free guidance (Ho and Salimans, 2022) are piratical implementations for conditional score estimation. For classifier guidance (Song et al., 2021; Dhariwal and Nichol, 2021), it use the gradient of the classifier to improve the conditional sample quality of the diffusion model. According to Bayes rule, the conditional score function has the relation:

$$\nabla_x \log p_t(x_t | y) = \underbrace{\nabla \log p_t(x_t)}_{\text{Approximate by } \hat{s}} + \underbrace{\nabla_x \log p_t(y | x_t)}_{\text{Guidance from classifier}}. \quad (\text{G.5})$$

It uses the neural network to approximate the unconditional score function $\nabla \log \widehat{p}_t(x_t)$ along with external classifier to approximate $\widehat{p}_t(y|x_t)$ and compute the gradient of the classifier logits as the guidance $\nabla \log \widehat{p}_t(y|x_t)$.

G.2 CLASSIFIER-FREE GUIDANCE

Classifier-free guidance (Ho and Salimans, 2022) provides a widely used approach for training condition diffusion models. It not only simplifies the training pipeline but also improves performance and removes the need for an external classifier. Classifier-free guidance diffusion model approximates both conditional and unconditional score functions by neural networks s_W , where W is the network parameters.

Our primary goal is to establish the theoretical guarantee for selecting conditional score estimator $\widehat{s}(x, y, t)$ chosen from the transformer architecture class and bound the error for such estimation. Based on previous work by Dhariwal and Nichol (2021); Fu et al. (2024b); Sohl-Dickstein et al. (2015); Ho and Salimans (2022), we adopt the unified setting for the conditional diffusion model. First we define the mask signal as $\tau := \{\emptyset, \text{id}\}$, where \emptyset denotes the absence of guidance y and id denotes otherwise. Unites the learning of conditional and unconditional scores by randomly ignoring the guidance y . Therefore we write the function class of the score estimator as

$$s(x, y, t) = \begin{cases} s_1(x, y, t), & \text{if } y \in \mathbb{R}^{d_y} \\ s_2(x, t), & \text{if } y = \emptyset. \end{cases} \quad (\text{G.6})$$

Both $s_1(x, y, t)$ and $s_2(x, t)$ belong to the transformer function class with slight adaption. Following Fu et al. (2024b), we consider $P(\tau = \text{id}) = P(\tau = \emptyset) = \frac{1}{2}$ without loss of generality, and we have the following objective function for score matching:

$$\widehat{s} := \operatorname{argmin}_{s_W \in \mathcal{T}_R^{h,s,r}} \mathbb{E}_t \left[\mathbb{E}_{(x_0, y)} \left[\mathbb{E}_{(\tau, x' \sim x' | x_0)} \left[\left\| s_W(x', \tau y, t) - \nabla_{x'} \log p_t(x' | x_0) \right\|_2^2 \right] \right] \right].$$

In practice, the loss function is given by

$$\ell(x_0, y; s_W) = \int_{T_0}^T \frac{1}{T - T_0} \mathbb{E}_{\tau, x_t | x_0 \sim N(\alpha_t x_0, \sigma_t^2 I_{d_x})} \left[\left\| s_W(x_t, \tau y, t) - \nabla_{x_t} \log p_t(x_t | x_0) \right\|_2^2 \right] dt, \quad (\text{G.7})$$

where T_0 is a small value for stabilize training (Vahdat et al., 2021). To train s_W , we select n i.i.d. training samples $\{x_{0,i}, y_i\}_{i=1}^n$, where $x_{0,i} \sim P_0(\cdot | y_i)$. We utilize the following empirical loss:

$$\widehat{\mathcal{L}}(s_W) = \frac{1}{n} \sum_{i=1}^n \ell(x_{0,i}, y_i; s_W). \quad (\text{G.8})$$

With the estimate score function $s_W(x, y, t)$ from minimizing the empirical loss in (G.8), we use $s_W(x, y, t)$ to generate new samples. In the classifier-free guidance setting, we generate a new conditional sample by replacing the approximation s_W in (G.4) with \widetilde{s}_W , defined as:

$$\widetilde{s}_W(x, y, t) = (1 + \eta) \cdot s_W(x, y, t) - \eta \cdot s_W(x, \emptyset, t), \quad (\text{G.9})$$

where the strength of guidance $\eta > 0$. The proper choice of η is crucial for balancing trade-offs between conditional guidance and unconditional ones. The choice directly impacts the performance of the generation process. Wu et al. (2024b) theoretically study the effect of guidance η on Gaussian mixture model. They demonstrate that strong guidance improves classification confidence but reduces sample diversity. For more detailed related work, refer to Appendix C.1.

2592 H UNIVERSAL APPROXIMATION OF TRANSFORMERS

2593 H.1 TRANSFORMERS AS UNIVERSAL APPROXIMATORS

2594 **Background: Contextual Mapping.** Let $X, Y \in \mathbb{R}^{d \times L}$ be the input and output label sequences,
2595 respectively. Let $X_{:,i} \in \mathbb{R}^d$ be the i -th token (column) of each X sequence.

2596 **Definition H.1** (Vocabulary). We define the i -th vocabulary set for $i \in [N]$ by $\mathcal{V}^{(i)} = \bigcup_{k \in [L]} X_{:,k}^{(i)} \subset$
2597 \mathbb{R}^d , and the whole vocabulary set \mathcal{V} is defined by $\mathcal{V} = \bigcup_{i \in [N]} \mathcal{V}^{(i)} \subset \mathbb{R}^d$.

2600 To facilitate our analysis, we introduce the idea of input token separation following (Kajitsuka and
2601 Sato, 2024; Kim et al., 2022; Yun et al., 2020).

2602 **Definition H.2** (Tokenwise Separateness). Let $Z^{(1)}, \dots, Z^{(N)} \in \mathbb{R}^{d \times L}$ be input sequences. Then,
2603 $Z^{(1)}, \dots, Z^{(N)}$ are called tokenwise $(\gamma_{\min}, \gamma_{\max}, \delta)$ -separated if the following three conditions hold.

- 2604 (i) For any $i \in [N]$ and $k \in [n]$, $\|Z_{:,k}^{(i)}\| > \gamma_{\min}$ holds.
- 2605 (ii) For any $i \in [N]$ and $k \in [n]$, $\|Z_{:,k}^{(i)}\| < \gamma_{\max}$ holds.
- 2606 (iii) For any $i, j \in [N]$ and $k, l \in [n]$ if $Z_{:,k}^{(i)} \neq Z_{:,l}^{(j)}$, then $\|Z_{:,k}^{(i)} - Z_{:,l}^{(j)}\| > \delta$ holds.

2607 Note that when only conditions (ii) and (iii) hold, we denote this as (γ, δ) -separateness. Moreover, if
2608 only condition (iii) holds, we denote it as (δ) -separateness.

2609 To clarify condition (iii), we consider cases where there are repeated tokens between different input
2610 sequences. Next, we define contextual mapping. Contextual mapping describes a function’s ability to
2611 capture the context of each input sequence as a whole and assign a unique ID to each input sequence.

2612 **Definition H.3** (Contextual Mapping). Let $X^{(1)}, \dots, X^{(N)} \in \mathbb{R}^{d \times L}$ be input sequences. Then, a
2613 map $q : \mathbb{R}^{d \times L} \rightarrow \mathbb{R}^{d \times L}$ is called an (γ, δ) -contextual mapping if the following two conditions hold:

- 2614 1. For any $i \in [N]$ and $k \in [L]$, $\|q(X^{(i)})_{:,k}\| < \gamma$ holds.
- 2615 2. For any $i, j \in [N]$ and $k, l \in [L]$ such that $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$ or $X_{:,k}^{(i)} \neq X_{:,l}^{(j)}$, $\|q(X^{(i)})_{:,k} -$
2616 $q(X^{(j)})_{:,l}\| > \delta$ holds.

2617 Note that $q(X^{(i)})$ for $i \in [N]$ is called a *context ID* of $X^{(i)}$.

2618 **Helper Lemmas.** To prove that 1-layer single-head attention is a contextual mapping, we first
2619 introduce some useful lemmas.

2620 **Lemma H.1** (Boltz Preserves Distance, Lemma 1 of (Kajitsuka and Sato, 2024)). Given (γ, δ) -
2621 tokenwise separated vectors $z^{(1)}, \dots, z^{(N)} \in \mathbb{R}^n$ with no duplicate entries in each vector, that
2622 is

$$2623 z_s^{(i)} \neq z_t^{(i)},$$

2624 where $i \in [N]$ and $s, t \in [L]$, $s \neq t$. Also, let

$$2625 \delta \geq 4 \ln n.$$

2626 Then, the outputs of the Boltzmann operator has the following property:

$$2627 \left| \text{Boltz} \left(z^{(i)} \right) \right| \leq \gamma, \tag{H.1}$$

$$2628 \left| \text{Boltz} \left(z^{(i)} \right) - \text{Boltz} \left(z^{(j)} \right) \right| > \delta' = \ln^2(n) \cdot e^{-2\gamma} \tag{H.2}$$

2629 for all $i, j \in [N]$, $i \neq j$.

Lemma H.2 (Lemma 13 of (Park et al., 2021)). For any finite subset $\mathcal{X} \subset \mathbb{R}^d$, there exists at least one unit vector $u \in \mathbb{R}^d$ such that

$$\frac{1}{|\mathcal{X}|^2} \sqrt{\frac{8}{\pi d}} \|x - x'\| \leq |u^\top (x - x')| \leq \|x - x'\|$$

for any $x, x' \in \mathcal{X}$.

With **Lemma H.2**, we present a configuration for weight matrices of a self-attention layer.

Lemma H.3 (Construction of Weight Matrices). Given a dataset with a $(\gamma_{\min}, \gamma_{\max}, \epsilon)$ -separated finite vocabulary $\mathcal{V} \subset \mathbb{R}^d$. There exists rank- ρ weight matrices $W_K, W_Q \in \mathbb{R}^{s \times d}$ such that

$$\left| (W_K v_a)^\top (W_Q v_c) - (W_K v_b)^\top (W_Q v_c) \right| > \delta,$$

for any $\delta > 0$, any $\min(d, s) \geq \rho \geq 1$ and any $v_a, v_b, v_c \in \mathcal{V}$ with $v_a \neq v_b$. In addition, the matrices are constructed as

$$W_K = \sum_{i=1}^{\rho} p_i q_i^\top \in \mathbb{R}^{s \times d}, \quad W_Q = \sum_{j=1}^{\rho} p'_j q'_j{}^\top \in \mathbb{R}^{s \times d},$$

where for at least one $i, q_i, q'_i \in \mathbb{R}^d$ are unit vectors that satisfy **Lemma H.2**, and $p_i, p'_i \in \mathbb{R}^s$ satisfies

$$|p_i^\top p'_i| = 5 (|\mathcal{V}| + 1)^4 d \frac{\delta}{\epsilon \gamma_{\min}}.$$

Proof of Lemma H.3. We build our proof upon (Kajitsuka and Sato, 2024). We start the proof by applying **Lemma H.2** to $\mathcal{V} \cup \{0\}$. We obtain at least one unit vector $q \in \mathbb{R}^d$ such that for any $v_a, v_b \in \mathcal{V} \cup \{0\}$ and $v_a \neq v_b$, we have

$$\frac{1}{(|\mathcal{V}| + 1)^2 d^{0.5}} \|v_a - v_b\| \leq |q^\top (v_a - v_b)| \leq \|v_a - v_b\|.$$

By choosing $v_b = 0$, we have that for any $v_c \in \mathcal{V}$

$$\frac{1}{(|\mathcal{V}| + 1)^2 d^{0.5}} \|v_c\| \leq |q^\top v_c| \leq \|v_c\|. \quad (\text{H.3})$$

For convenience, we denote the set of all unit vector q that satisfies (H.3) as \mathbb{Q} . Next, we choose some arbitrary vector pairs $p_i, p'_i \in \mathbb{R}^s$ that satisfy

$$|p_i^\top p'_i| = (|\mathcal{V}| + 1)^4 d \frac{\delta}{\epsilon \gamma_{\min}}. \quad (\text{H.4})$$

We construct the weight matrices by setting

$$W_K = \sum_{i=1}^{\rho} p_i q_i^\top \in \mathbb{R}^{s \times d},$$

$$W_Q = \sum_{j=1}^{\rho} p'_j q'_j{}^\top \in \mathbb{R}^{s \times d},$$

where for at least one i, p_i, p'_i satisfies (H.4) and $q_i, q'_i \in \mathbb{Q}$. Here, $\mathbb{Q} = \{q \in \mathbb{R}^n : \|q\| = 1\}$ denotes the set of all unit vectors in \mathbb{R}^n . We arrive at

$$\begin{aligned}
& \left| (W_K v_a)^\top (W_Q v_c) - (W_K v_b)^\top (W_Q v_c) \right| \\
&= \left| (v_a - v_b)^\top (W_K)^\top (W_Q v_c) \right| \\
&= \left| (v_a - v_b)^\top \left(\sum_{i=1}^{\rho} q_i p_i^\top \right) \left(\sum_{j=1}^{\rho} p'_j q'_j{}^\top v_c \right) \right| \\
&= \left| \left(\sum_{i=1}^{\rho} (v_a - v_b)^\top q_i p_i^\top \right) \left(\sum_{j=1}^{\rho} p'_j q'_j{}^\top v_c \right) \right| \\
&= \left| \sum_{i=1}^{\rho} \sum_{j=1}^{\rho} (v_a - v_b)^\top q_i p_i^\top p'_j q'_j{}^\top v_c \right| \\
&= \sum_{i=1}^{\rho} \sum_{j=1}^{\rho} \left| (v_a - v_b)^\top q_i \right| \cdot |p_i^\top p'_j| \cdot |q'_j{}^\top v_c| \\
&\geq \frac{1}{(|\mathcal{V}| + 1)^2 d^{0.5}} \|v_a - v_b\| \cdot (|\mathcal{V}| + 1)^4 d \frac{\delta}{\epsilon \gamma_{\min}} \cdot \frac{1}{(|\mathcal{V}| + 1)^2 d^{0.5}} \|v_c\| \quad (\text{By (H.3) and (H.4)}) \\
&> \delta. \quad (\text{By } (\gamma_{\min}, \gamma_{\max}, \epsilon)\text{-separateness of } \mathcal{V})
\end{aligned}$$

This completes the proof. \square

Any-Rank Attention is Contextual Mapping. Now we present the result showing that a softmax-based 1-head, 1-layer attention block with any-rank weight matrices is a contextual mapping.

Theorem H.1 (Any-Rank Attention is (γ, δ) -Contextual Mapping, Modified from Theorem 2 of (Kajitsuka and Sato, 2024)). Given input sequences $X^{(1)}, \dots, X^{(N)} \in \mathbb{R}^{d \times L}$ which are $(\gamma_{\min}, \gamma_{\max}, \epsilon)$ -tokenwise separated and vocabulary set $\mathcal{V} = \bigcup_{i \in [N]} \mathcal{V}^{(i)} \subset \mathbb{R}^d$. Also, let $X^{(1)}, \dots, X^{(N)} \in \mathbb{R}^{d \times L}$ be sequences with no duplicate word token in each sequence, that is, $X_{:,k}^{(i)} \neq X_{:,l}^{(i)}$, for any $i \in [N]$ and $k, l \in [L]$. Then, there exists a 1-layer single head attention with weight matrices $W_O \in \mathbb{R}^{d \times s}$ and $W_V, W_K, W_Q \in \mathbb{R}^{s \times d}$, that is a (γ, δ) -contextual mapping for the input sequences $X^{(1)}, \dots, X^{(N)}$ with $\gamma = \gamma_{\max} + \epsilon/4$, $\delta = \exp(-5\epsilon^{-1}|\mathcal{V}|^4 d \kappa \gamma_{\max} \log L)$ where $\kappa = \gamma_{\max}/\gamma_{\min}$.

Theorem H.1 indicates that any-rank self-attention function distinguishes input tokens $X_{:,k}^{(i)} = X_{:,l}^{(j)}$ such that $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$. In other words, it distinguishes two identical tokens within a different context.

Remark H.1 (Comparing with Existing Works). In comparison with (Kajitsuka and Sato, 2024), they provide a proof for the case where all self-attention weight matrices $W_V, W_K, W_Q \in \mathbb{R}^{s \times d}$ are strictly rank-1. However, this is almost impossible in practice for any pre-trained transformer-based models. Here, by considering self-attention weight matrices of rank ρ where $1 \leq \rho \leq \min(d, s)$, we show that single-head, single-layer self-attention with matrices of any rank is a contextual mapping, pushing the universality of (prompt tuning) transformers towards more practical scenarios.

Remark H.2. In (Kajitsuka and Sato, 2024), γ and δ are chosen as follows:

$$\gamma = \gamma_{\max} + \frac{\epsilon}{4}, \quad \delta = \frac{2(\ln L)^2 \epsilon^2 \gamma_{\min}}{\gamma_{\max}^2 (|\mathcal{V}| + 1)^4 (2 \ln L + 3) \pi d} \exp \left(-(|\mathcal{V}| + 1)^4 \frac{(2 \ln L + 3) \pi d \gamma_{\max}^2}{4 \epsilon \gamma_{\min}} \right).$$

Since the exponential term dominates the polynomial terms, in Lemma H.1, we simplify δ to $\exp(-\Theta(\epsilon^{-1}|\mathcal{V}|^4 d \kappa \gamma_{\max} \ln L))$.

Proof Sketch. We generalize the results of (Kajitsuka and Sato, 2024, Theorem 2) where all weight matrices have to be rank-1. We eliminate the rank-1 requirement, and extend the lemma for weights

of any rank ρ . This is achieved by constructing the weight matrices as a outer product sum $\sum_i^\rho u_i v_i^\top$, where $u_i \in \mathbb{R}^s, v_i \in \mathbb{R}^d$. Specifically, we divide the proof into two parts:

- We first construct a softmax-based self-attention that maps different input tokens to unique contextual embeddings, by configuring weight matrices according to [Lemma H.3](#).
- Secondly, for the identical tokens within a different context, we utilize the tokenwise separateness guaranteed by [Lemma H.3](#) and [Lemma H.1](#) which shows Boltz preserves some separateness.

As a result, we prove that the self-attention function distinguishes input tokens $X_{:,k}^{(i)} = X_{:,l}^{(j)}$ such that $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$. This completes the proof. \square

Proof of Theorem H.1. We build our proof upon ([Kajitsuka and Sato, 2024](#)). We construct a self-attention layer that is a contextual mapping. There are mainly two things to prove. We first show that the attention later we constructed maps different tokens to unique ids. Secondly, we prove that the self-attention function distinguishes duplicate input tokens within different context. For the first part, we show that our self-attention layer satisfies:

$$\|\Psi\| = \left\| W_O \left(W_V X^{(i)} \right) \text{Softmax} \left[\left(W_K X^{(i)} \right)^\top \left(W_Q X_{:,k}^{(i)} \right) \right] \right\| < \frac{\epsilon}{4}, \quad (\text{H.5})$$

for $i \in [N]$ and $k \in [L]$. Since with (H.5), it is easy to show that

$$\begin{aligned} \left\| f^{(SA)} \left(X^{(i)} \right)_{:,k} - f^{(SA)} \left(X^{(j)} \right)_{:,l} \right\| &= \left\| X_{:,k}^{(i)} - X_{:,l}^{(j)} + \left(\Psi^{(i)} - \Psi^{(j)} \right) \right\| \\ &\geq \left\| X_{:,k}^{(i)} - X_{:,l}^{(j)} \right\| - \left\| \Psi^{(i)} - \Psi^{(j)} \right\| \\ &\geq \left\| X_{:,k}^{(i)} - X_{:,l}^{(j)} \right\| - \left\| \Psi^{(i)} \right\| - \left\| \Psi^{(j)} \right\| \\ &> \epsilon - \frac{\epsilon}{4} - \frac{\epsilon}{4} = \frac{\epsilon}{2}, \quad (\text{By } \epsilon\text{-separatedness of } X \text{ and H.5}) \end{aligned} \quad (\text{H.6})$$

for any $i, j \in [N]$ and $k, l \in [L]$ such that $X_{:,k}^{(i)} \neq X_{:,l}^{(j)}$. Now, we prove (H.5) by utilizing [Lemma H.3](#). We define the weight matrices as

$$\begin{aligned} W_K &= \sum_{i=1}^\rho p_i q_i^\top \in \mathbb{R}^{s \times d}, \\ W_Q &= \sum_{j=1}^\rho p'_j q'_j{}^\top \in \mathbb{R}^{s \times d}, \end{aligned}$$

where $p_i, p'_j \in \mathbb{R}^s$ and $q_i, q'_j \in \mathbb{R}^d$. In addition, let $\delta = 4 \ln n$ and $p_1, p'_1 \in \mathbb{R}^s$ be an arbitrary vector pair that satisfies

$$\left| p_1^\top p'_1 \right| = (|\mathcal{V}| + 1)^4 d \frac{\delta}{\epsilon \gamma_{\min}}. \quad (\text{H.7})$$

Then by [Lemma H.3](#), there are some unit vectors q_1, q'_1 such that we have,

$$\left| (W_K v_a)^\top (W_Q v_c) - (W_K v_b)^\top (W_Q v_c) \right| > \delta, \quad (\text{H.8})$$

for any $v_a, v_b, v_c \in \mathcal{V}$ with $v_a \neq v_b$. In addition, for the other two weight matrices $W_O \in \mathbb{R}^{d \times s}$ and $W_V \in \mathbb{R}^{s \times d}$, we set

$$W_V = \sum_{i=1}^{\rho} p_i'' q_i''^\top \in \mathbb{R}^{s \times d}, \quad (\text{H.9})$$

where $q'' \in \mathbb{R}^d$, $q_1'' = q_1$ and $p_i'' \in \mathbb{R}^s$ is some nonzero vector that satisfies

$$\|W_O p_i''\| = \frac{\epsilon}{4\rho\gamma_{\max}}, \quad (\text{H.10})$$

This can be accomplished, e.g., $W_O = \sum_{i=1}^{\rho} p_i''' p_i'''^\top$ for any vector p_i''' which satisfies $\|p_i'''\| = \epsilon / (4\rho^2\gamma_{\max}\|p_i''\|^2)$ for any $i \in [\rho]$. As a result, we now bound Ψ as:

$$\begin{aligned} \|\Psi\| &= \left\| W_O \left(W_V X^{(i)} \right) \text{Softmax} \left[\left(W_K X^{(i)} \right)^\top \left(W_Q X_{:,k}^{(i)} \right) \right] \right\| \\ &= \left\| \sum_{k'=1}^L s_{k'}^k W_O \left(W_V X^{(i)} \right)_{:,k'} \right\| \quad (\text{Denote } s_{k'}^k = \text{Softmax} \left[\left(W_K X^{(i)} \right)^\top \left(W_Q X_{:,k}^{(i)} \right) \right]_{k'}) \\ &= \sum_{k'=1}^L s_{k'}^k \left\| W_O \left(W_V X^{(i)} \right)_{:,k'} \right\| \\ &\leq \max_{k' \in [L]} \left\| W_O \left(W_V X^{(i)} \right)_{:,k'} \right\| \quad (\sum_{k'=1}^L s_{k'}^k = 1) \\ &= \max_{k' \in [L]} \left\| W_O \left(\sum_{i=1}^{\rho} p_i'' q_i''^\top \right) X_{:,k'}^{(i)} \right\| \quad (\text{By Lemma H.3}) \\ &= \sum_{i=1}^{\rho} \|W_O p_i''\| \cdot \max_{k' \in [L]} \left| q_i''^\top X_{:,k'}^{(i)} \right| \quad (\text{By (H.10)}) \\ &= \frac{\epsilon}{4\gamma_{\max}} \cdot \max_{k' \in [L]} \left\| X_{:,k'}^{(i)} \right\| \quad (\text{By (H.10) and } \|q_i''\| = 1) \\ &< \frac{\epsilon}{4}. \end{aligned}$$

Next, for the second part, we prove that with the weight matrices W_O, W_V, W_K, W_Q configured above, the attention layer distinguishes duplicate input tokens with different context, $X_{:,k}^{(i)} = X_{:,l}^{(j)}$ with $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$. We choose any $i, j \in [N]$ and $k, l \in [L]$ such that $X_{:,k}^{(i)} = X_{:,l}^{(j)}$ and $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$.

In addition, we define $a^{(i)}, a^{(j)}$ as

$$a^{(i)} = \left(W_K X^{(i)} \right)^\top \left(W_Q X_{:,k}^{(i)} \right) \in \mathbb{R}^n, \quad a^{(j)} = \left(W_K X^{(j)} \right)^\top \left(W_Q X_{:,l}^{(j)} \right) \in \mathbb{R}^n.$$

From (H.8) we have that $a^{(i)}$ and $a^{(j)}$ are tokenwise (γ, δ) -separated where γ is computed by

$$\begin{aligned} |a_{k'}^{(i)}| &= \left| \left(W_K X_{:,k'}^{(i)} \right)^\top \left(W_Q X_{:,k}^{(i)} \right) \right| \\ &= \left| \left(\sum_{i=1}^{\rho} p_i q_i^\top X_{:,k'}^{(i)} \right)^\top \left(\sum_{j=1}^{\rho} p_j' q_j'^\top X_{:,k}^{(i)} \right) \right| \\ &= \left| \left(\sum_{i=1}^{\rho} X_{:,k'}^{(i)\top} q_i p_i^\top \right) \left(\sum_{j=1}^{\rho} p_j' q_j'^\top X_{:,k}^{(i)} \right) \right| \end{aligned}$$

$$\begin{aligned}
&= \left| \sum_{i=1}^{\rho} \sum_{j=1}^{\rho} X_{:,k'}^{(i)\top} q_i p_i^\top p'_j q'_j{}^\top X_{:,k}^{(i)} \right| \\
&= \sum_{i=1}^{\rho} \sum_{j=1}^{\rho} \left| X_{:,k'}^{(i)\top} q_i \right| |p_i^\top p'_j| |q'_j{}^\top X_{:,k}^{(i)}| \\
&\leq (|\mathcal{V}| + 1)^4 d \frac{\delta}{\epsilon \gamma_{\min}} \gamma_{\max}^2. \quad (\text{By (H.7) and } \|q_i\| = \|q'_j\| = 1)
\end{aligned}$$

Therefore,

$$\gamma = (|\mathcal{V}| + 1)^4 d \frac{\delta \gamma_{\max}^2}{\epsilon \gamma_{\min}}.$$

Now, since $\mathcal{V}^{(i)} \neq \mathcal{V}^{(j)}$ and there is no duplicate token in $X^{(i)}$ and $X^{(j)}$ respectively, we use [Lemma H.1](#) and obtain that

$$\begin{aligned}
\left| \text{Boltz} \left(a^{(i)} \right) - \text{Boltz} \left(a^{(j)} \right) \right| &= \left| \left(a^{(i)} \right)^\top \text{Softmax} \left[a^{(i)} \right] - \left(a^{(j)} \right)^\top \text{Softmax} \left[a^{(j)} \right] \right| \quad (\text{H.11}) \\
&> \delta' \\
&= (\ln n)^2 e^{-2\gamma}.
\end{aligned}$$

As we assumed $X_{:,k}^{(i)} = X_{:,l}^{(j)}$, we have

$$\begin{aligned}
&\left| \left(a^{(i)} \right)^\top \text{Softmax} \left[a^{(i)} \right] - \left(a^{(j)} \right)^\top \text{Softmax} \left[a^{(j)} \right] \right| \quad (\text{H.12}) \\
&= \left| \left(X_{:,k}^{(i)} \right)^\top \left(W_Q \right)^\top W_K \left(X^{(i)} \text{Softmax} \left[a^{(i)} \right] - X^{(j)} \text{Softmax} \left[a^{(j)} \right] \right) \right| \\
&= \left| \left(X_{:,k}^{(i)} \right)^\top \left(\sum_{j=1}^{\rho} q'_j p'_j{}^\top \right) \left(\sum_{i=1}^{\rho} p_i q_i^\top \right) \left(X^{(i)} \text{Softmax} \left[a^{(i)} \right] - X^{(j)} \text{Softmax} \left[a^{(j)} \right] \right) \right| \\
&\quad (\text{By Lemma H.3}) \\
&= \sum_{i=1}^{\rho} \sum_{j=1}^{\rho} \left| q'_j{}^\top X_{:,k}^{(i)} \right| \cdot |p_j^\top p_i| \cdot \left| \left(q_i^\top X^{(i)} \right) \text{Softmax} \left[a^{(i)} \right] - \left(q_i^\top X^{(j)} \right) \text{Softmax} \left[a^{(j)} \right] \right| \\
&\leq \sum_{i=1}^{\rho} \gamma_{\max} \cdot (|\mathcal{V}| + 1)^4 \frac{\pi d}{8} \frac{\delta}{\epsilon \gamma_{\min}} \cdot \left| \left(q_i^\top X^{(i)} \right) \text{Softmax} \left[a^{(i)} \right] - \left(q_i^\top X^{(j)} \right) \text{Softmax} \left[a^{(j)} \right] \right|. \\
&\quad (\text{By (H.7)})
\end{aligned}$$

By combining [\(H.11\)](#) and [\(H.12\)](#), we have

$$\sum_{i=1}^{\rho} \left| \left(q_i^\top X^{(i)} \right) \text{Softmax} \left[a^{(i)} \right] - \left(q_i^\top X^{(j)} \right) \text{Softmax} \left[a^{(j)} \right] \right| > \frac{\delta'}{(|\mathcal{V}| + 1)^4 d \delta \gamma_{\max}} \epsilon \gamma_{\min}. \quad (\text{H.13})$$

Now we arrive at the lower bound of the difference between the self-attention outputs of $X^{(i)}$, $X^{(j)}$ as:

$$\begin{aligned}
&\left\| f_S^{(\text{SA})} \left(X^{(i)} \right)_{:,k} - f_S^{(\text{SA})} \left(X^{(j)} \right)_{:,l} \right\| \quad (\text{H.14}) \\
&= \left\| W_O \left(W_V X^{(i)} \right) \text{Softmax} \left[a^{(i)} \right] - W_O \left(W_V X^{(j)} \right) \text{Softmax} \left[a^{(j)} \right] \right\|
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^{\rho} \|W_O p_i''\| \cdot \left| \left(q_i''^\top X^{(i)} \right) \text{Softmax} \left[a^{(i)} \right] - \left(q_i''^\top X^{(j)} \right) \text{Softmax} \left[a^{(j)} \right] \right| \\
&\quad (W_V = \sum_{i=1}^{\rho} p_i'' q_i''^\top) \\
&> \frac{\epsilon}{4\gamma_{\max}} \frac{\delta'}{(|\mathcal{V}| + 1)^4} \frac{\epsilon\gamma_{\min}}{d\delta\gamma_{\max}}. \quad (\text{By (H.10) and (H.13)})
\end{aligned}$$

where $\delta = 4 \ln L$ and $\delta' = \ln^2(L) e^{-2\gamma}$ with $\gamma = (|\mathcal{V}| + 1)^4 d\delta\gamma_{\max}^2 / (\epsilon\gamma_{\min})$. Note that we are able to use (H.13) in the last inequality of (H.14) because (H.13) is guaranteed by q_1 , and we set $q_1'' = q_1$ when constructing W_V in (H.9). □

Theorem H.2 (Transformers with 1-Layer Self-Attention are Universal Approximators, Modified from Proposition 1 of (Kajitsuka and Sato, 2024)). Let $0 \leq p < \infty$ 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 $f' \in \mathcal{T}_R^{h,s,r}$ such that $d_p(f, f') < \epsilon$ holds

Proof of Theorem H.2. We restate the proof from (Kajitsuka and Sato, 2024) for completeness.

The proof consists of the following steps:

1. Approximate by Step Function: Given a permutation equivariant continuous function f on a compact set, there exists a Transformer $f' \in \mathcal{T}_R^{h,s,r}$ with one self-attention layer to approximate f by step function with arbitrary precision in terms of p -norm.
2. Quantization via f_1^{FF} : The first feed-forward network f_1^{FF} quantize the input domain, reducing the problem to memorization of finite samples.
3. Contextual Mapping $f^{(\text{SA})}$ and Memorization f_2^{FF} : According to Theorem H.1, we construct any-rank attention $f^{(\text{SA})}$ to be contextual mapping. Then use the second feed-forward f_2^{FF} to memorize the *context ID* with its corresponding label.

The details for the three steps are below.

1. Since f is a continuous function on a compact set, f has maximum and minimum values on the domain. By scaling with f_1^{FF} and f_2^{FF} , f is assumed to be normalized without loss of generality: That is for any $Z \in \mathbb{R}^{d \times L} \setminus [0, 1]^{d \times L}$, we have $f(Z) = 0$. For any $X \in [-1, 1]^{d \times L}$, the function $f(X)$ satisfies $-1 \leq f(X) \leq 1$.

Let $D \in \mathbb{N}$ be the granularity of a grid

$$\mathbb{G}_D = \{1/D, 2/D, \dots, 1\}^{d \times L} \subset \mathbb{R}^{d \times L}$$

such that a piece-wise constant approximation

$$\bar{f}(X) = \sum_{L \in \mathbb{G}_D} f(L) 1_{Z \in L + [-1/D, 0]^{d \times L}}$$

satisfies

$$d_p(f, \bar{f}) < \epsilon/3. \quad (\text{H.15})$$

Such a D always exists because of uniform continuity of f .

2. We use f_1^{FF} to quantize the input domain into \mathbb{G}_D .

We first define the following two terms for first feed-forward neural network to approximate.

2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023

- The quantize term ($\text{quant}_D^{d \times L} : \mathbb{R}^{d \times L} \rightarrow \mathbb{R}^{d \times L}$): Quantize $[0, 1]$ into $\{1/D, \dots, 1\}$, while it projects $\mathbb{R} \setminus [0, 1]$ to 0 by shifting and stacking step function.

$$\begin{aligned} & \sum_{t=0}^{D-1} \frac{\text{ReLU}[x/\delta - t/\delta D] - \text{ReLU}[x/\delta - 1 - t/\delta D]}{D} \\ & \approx \text{quant}_D(x) = \begin{cases} 0 & x < 0 \\ 1/D & 0 \leq x < 1/D \\ \vdots & \vdots \\ 1 & 1 - 1/D \leq x \end{cases}. \end{aligned} \quad (\text{H.16})$$

- The penalty term (penalty): Identify whether an input sequence is in $[0, 1]^{d \times L}$. This is defined by

$$\begin{aligned} & \text{ReLU}[(x-1)/\delta] - \text{ReLU}[(x-1)/\delta - 1] - \text{ReLU}[-x/\delta] - \text{ReLU}[-x/\delta - 1] \\ & \approx \text{penalty}(x) = \begin{cases} -1 & x \leq 0 \\ 0 & 0 < x \leq 1 \\ -1 & 1 < x \end{cases}. \end{aligned} \quad (\text{H.17})$$

Combining these components together, the first feed-forward neural network layer f_1^{FF} approximates the following function:

$$\bar{f}_1^{(\text{FF})}(X) = \text{quant}_D^{d \times L}(X) + \sum_{t=1}^d \sum_{k=1}^L \text{penalty}(X_{t,k}) \quad (\text{H.18})$$

Note that this function quantizes inputs in $[0, 1]^{d \times L}$ with granularity D , while every element of the output is non-positive for inputs outside $[0, 1]^{d \times L}$. In particular, the norm of the output is upper-bounded by

$$\max_{X \in \mathbb{R}^{d \times L}} \|f_1^{\text{FF}}(X)_{:,k}\| = \underbrace{dL}_{\text{Total number of elements in } X} \times \underbrace{\sqrt{d}}_{\text{Maximum Euclidean norm in } d\text{-dimensional space}} \quad (\text{H.19})$$

for any $k \in [L]$.

- Let $\tilde{\mathbb{G}}_D \subset \mathbb{G}_D$ be a sub-grid

$$\tilde{\mathbb{G}}_D = \{G \in \mathbb{G}_D \mid \forall k, l \in [L], G_{:,k} \neq G_{:,l}\},$$

and consider memorization of $\tilde{\mathbb{G}}_D$ with its labels given by $f(G)$ for each $G \in \tilde{\mathbb{G}}_D$. Using our modified any-rank attention is contextual mapping in [Theorem H.1](#) allows us to construct a self-attention $f^{(\text{SA})}$ to be a contextual mapping for such input sequences, because $\tilde{\mathbb{G}}_D$ can be regarded as tokenwise $(1/D, \sqrt{d}, 1/D)$ -separated input sequences. By taking sufficiently large granularity D of \mathbb{G}_D , the number of cells with duplicate tokens, that is, $|\mathbb{G}_D \setminus \tilde{\mathbb{G}}_D|$ is negligible.

From the way the self-attention $f^{(\text{SA})}$ is constructed, we have

$$\|f^{(\text{SA})}(X)_{:,k} - X_{:,k}\| < \frac{1}{4\sqrt{d}D} \max_{k' \in [L]} \|X_{:,k'}\|$$

for any $k \in [L]$ and $X \in \mathbb{R}^{d \times L}$.

If we take large enough D , every element of the output for $X \in \mathbb{R}^{d \times L} \setminus [0, 1]^{d \times L}$ is upper-bounded by

$$f^{(\text{SA})} \circ f_1^{\text{FF}}(X)_{t,k} < \frac{1}{4D} \quad (\forall t \in [d], k \in [L]),$$

while the output for $X \in [0, 1]^{d \times L}$ is lower-bounded by

$$f^{(\text{SA})} \circ f_1^{\text{FF}}(X)_{t,k} > \frac{3}{4D} \quad (\forall t \in [d], k \in [L]).$$

Finally, we construct bump function of scale $R > 0$ to map each input sequence $L \in \tilde{\mathbb{G}}_D$ to its labels $f(L)$ and for input sequence outside the range $X \in (-\infty, 1/4D)^{d \times L}$ to 0 using the second feed-forward f_2^{FF} . Precisely, bump function of scale $R > 0$ is given by

$$\begin{aligned} \text{bump}_R(x) = \frac{f(L)}{dL} \sum_{t=1}^d \sum_{k=1}^L & (\text{ReLU}[R(X_{t,k} - G_{t,k}) - 1] - \text{ReLU}[R(Z_{t,k} - G_{t,k})] \\ & + \text{ReLU}[R(Z_{t,k} - G_{t,k}) + 1]) + \text{ReLU}[R(G_{t,k} - Z_{t,k})] \end{aligned} \quad (\text{H.20})$$

for each input sequence $G \in \tilde{\mathbb{G}}_D$ and add up these functions to implement f_2^{FF} .

In addition, the value of $f_2^{(\text{FF})}$ is always bounded: $0 \leq f_2^{(\text{FF})} \leq 1$. Thus, by taking sufficiently small $\delta > 0$ to quantize the step function, we have

$$d_p \left(f_2^{(\text{FF})} \circ f^{(\text{SA})} \circ f_1^{(\text{FF})}, f_2^{(\text{FF})} \circ f^{(\text{SA})} \circ \bar{f}_1^{(\text{FF})} \right) < \frac{\epsilon}{3}. \quad (\text{H.21})$$

Taking large enough D to make duplicate tokens negligible, we have

$$d_p \left(f_2^{(\text{FF})} \circ f^{(\text{SA})} \circ \bar{f}_1^{(\text{FF})}, \bar{f} \right) < \frac{\epsilon}{3}. \quad (\text{H.22})$$

Combining estimation of step function (H.15), estimation of quantization (H.21) and estimation of duplicate tokens (H.22) together, we get the approximation error of the any-rank Transformer as

$$d_p \left(f_2^{(\text{FF})} \circ f^{(\text{SA})} \circ \bar{f}_1^{(\text{FF})}, f \right) < \epsilon. \quad (\text{H.23})$$

This completes the proof. \square

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

Corollary H.2.1 (Universal Approximation of Transformers). From [Theorem H.2](#), for any permutation 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 approximate f such that $d_p(f, f') < \epsilon$.

Remark H.3. We remark that $\mathcal{T}_R^{1,1,4}$ belongs to the considered transformer network function class [Definition 2.2](#).

We establish in [Corollary H.2.1](#) the minimal transformer configuration required to achieve universal approximation for compactly supported functions. We remark that this configuration is minimally sufficient but not necessary. More complex configurations can also achieve transformer universality, as reported in (Hu et al., 2024; Kajitsuka and Sato, 2024; Yun et al., 2020).

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

H.2 PARAMETER NORM BOUNDS FOR TRANSFORMER APPROXIMATION

In the analysis of the approximation ability of transformers in (Kajitsuka and Sato, 2024), universal approximation is ensured by using a sufficiently large granularity D , a sufficiently small δ in $f_1^{(\text{FF})}$, and an appropriate scaling factor R in $f_2^{(\text{FF})}$. 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.

Lemma H.4 (Order of Granularity and Scaling Factor). Consider the universal approximation theorem for transformers in Theorem H.2. The order for the granularity and the scaling factor follows $D = \mathcal{O}(\epsilon^{-1/d})$ and $R = \mathcal{O}(D)$, and the parameter δ for the first feed-forward layer in (H.16) follows $\delta = o(D^{-1})$.

Proof. We investigate the more precise choice of D , R , and δ respectively.

- **Bound on Scaling Factor in $f_2^{(\text{FF})}$.**

First, we need to ensure that $R > 0$ is large enough such that it maps input $Z \in (-\infty, \frac{1}{4D})^{d \times L}$ to zero.

Because we have $Z_{t,k} - L_{t,k} \leq -\frac{3}{4D}$, we obtain the desired result from (H.20) by taking $R = \mathcal{O}(D)$ such that three $\text{ReLU}(\cdot)$ output zero.

Second, we need to ensure that $R > 0$ is large enough such that it maps $L \in \tilde{\mathbb{G}} \subset (\frac{3}{4D}, \infty)^{d \times L}$ to the corresponding label $f(L)$.

From (H.20), we achieve this by selecting proper R such that

$$\sum_{t=1}^d \sum_{k=1}^L \text{ReLU}[RS - 1] - \text{ReLU}[RS] + \text{ReLU}[RS + 1] \text{ReLU}[-RS] = dL,$$

where $S := Z_{t,k} - L_{t,k} = \mathcal{O}(D^{-1})$.

For any $S \in \mathbb{R}$, we take $R = \mathcal{O}(D)$ such that $|RS| \leq 1$.

- **Bound on Granularity D .**

In (Kajitsuka and Sato, 2024), there are $\mathcal{O}(D^{-d}|\mathbb{G}_D|)$ omitted duplicated input. Clearly, by taking sufficiently large granularity $|\mathbb{G}_D \setminus \tilde{\mathbb{G}}_D|$ becomes negligible, but here we aim to evaluate the corresponding order of D .

First, by the extreme value theorem, the continuous function f on $[0, 1]^{d \times L}$ here is bounded by some constant, denoted by B .

Second, the total omitted points are $\mathcal{O}(D^{d(L-1)})$.

Third, the probability for each point in \mathbb{G}_D is $1/D^{dL}$.

Therefore, the corresponding error is bounded by $\mathcal{O}(D^{-d/p})$. Since we require error to be bounded $\epsilon/3$, setting $D = \mathcal{O}(\epsilon^{-p/d})$ for some constant $p > 0$ guarantees the result. We provide the detailed derivations as follows.

We follow (Kajitsuka and Sato, 2024) considering Lipschitz (under p -norm) function class of continuous sequence-to-sequence. This consideration is practical as realistic input of transformer blocks are vector embedding in Euclidean space. Let $f(\cdot) : [0, 1]^{d \times L} \rightarrow [0, 1]^{d \times L}$ be the target function and $\tilde{f}(\cdot)$ be the piece-wise constant approximation of regularity D . Recall the p -norm

3132 difference between two function $f(\cdot)$ and $\bar{f}(\cdot)$. (H.15) gives

$$\begin{aligned}
 3133 \quad d_p(f, \bar{f}) &= \left(\int \|f(x) - \bar{f}(x)\|^p dx \right)^{1/p} \\
 3134 &= \mathcal{O}(D^{dL-d}) \cdot (B^p(1/D)^{dL})^{1/p} \\
 3135 &= \mathcal{O}(D^{(dL-d)/p}) \cdot \mathcal{O}(D^{-dL/p}) \\
 3136 &= \mathcal{O}(D^{-d/p}).
 \end{aligned}$$

3137 Here, $\mathcal{O}(D^{-d/p}) = \epsilon$ implies $D = \mathcal{O}(\epsilon^{-p/d})$ for some constant $p > 0$. For simplicity, we use
 3138 $D = \mathcal{O}(\epsilon^{-1/d})$ in our analysis without loss of generality.

3139 • **Bound on Parameter δ in $f_1^{(\text{FF})}$.**

3140 In the quantization operation realized by the network, we need to ensure the error within region
 3141 $(i/D, i/D + \delta)$ does not affect the desired interval $(i/D, (i+1)/D)$ for $i \in [D]$.

3142 Thus, we need $\delta = o(1/D)$.

3143 This completes the proof. \square

3144 Building upon Lemma H.4, we extend the results to derive explicit parameter bounds for matrices
 3145 regarding the transformer-based universal approximation framework. That is, we ensure a more
 3146 precise quantification of parameter constraints across the architecture.

3147 **Lemma H.5** (Transformer Matrices Bounds). Consider an input sequence $Z \in [0, 1]^{d \times L}$. Let $f(Z) :$
 3148 $[0, 1]^{d \times L} \rightarrow \mathbb{R}^{d \times L}$ be any permutation equivariant and continuous sequence-to-sequence function
 3149 on compact support $[0, 1]^{d \times L}$. For the transformer network $f' \in \mathcal{T}_R^{r,h,s}$ defined in Definition 2.4 to
 3150 approximate f within ϵ precision, i.e., $d_p(f, f') < \epsilon$, the following parameter bounds must hold for
 3151 $d \geq 1$ and $L \geq 2$:

$$\begin{aligned}
 3152 \quad \|W_Q\|_2 &= \|W_K\|_2 = \mathcal{O}(d \cdot \epsilon^{-(\frac{2dL+1}{d})})(\log L)^{\frac{1}{2}}; \\
 3153 \quad \|W_Q\|_{2,\infty} &= \|W_K\|_{2,\infty} = \mathcal{O}(d^{\frac{3}{2}} \cdot \epsilon^{-(\frac{2dL+1}{d})})(\log L)^{\frac{1}{2}}; \\
 3154 \quad \|W_O\|_2 &= \mathcal{O}(\sqrt{d}\epsilon^{\frac{1}{d}}); \|W_O\|_{2,\infty} = \mathcal{O}(\epsilon^{\frac{1}{d}}); \\
 3155 \quad \|W_V\|_2 &= \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \\
 3156 \quad \|W_1\|_2 &= \mathcal{O}(d\epsilon^{-\frac{1}{d}}), \|W_1\|_{2,\infty} = \mathcal{O}(\sqrt{d}\epsilon^{-\frac{1}{d}}); \\
 3157 \quad \|W_2\|_2 &= \mathcal{O}(d\epsilon^{-\frac{1}{d}}); \|W_2\|_{2,\infty} = \mathcal{O}(\sqrt{d}\epsilon^{-\frac{1}{d}}); \\
 3158 \quad \|E^\top\|_{2,\infty} &= \mathcal{O}(d^{\frac{1}{2}}L^{\frac{3}{2}}).
 \end{aligned}$$

3159 For the case $L = 1$, the parameter bounds remain valid with the substitution of $\log L$ with 1.

3160 *Proof.* For the self-attention layer, we denote the separatedness of the input tokens by $(\gamma_{\min}, \gamma_{\max}, \epsilon_s)$
 3161 and the separatedness of the output tokens by (γ, δ_s) . Moreover, in (H.16) we denote the parameter
 3162 taken in f_1^{FF} corresponding to the granularity by δ_{f_1} .

3163 • **Bounds for W_Q and W_K in $f^{(\text{SA})}$.**

3164 From the universal approximation theorem of transformer Theorem H.2, with $p_i, p'_i \in \mathbb{R}^s$ and
 3165 q_i, q'_i , being any unit vectors in \mathbb{R}^d , we construct rank ρ matrix W_Q and W_K as

$$3166 \quad W_K = \sum_{i=1}^{\rho} p_i q_i^\top \in \mathbb{R}^{s \times d},$$

3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239

$$W_Q = \sum_{i=1}^{\rho} p'_i q_i{}^\top \in \mathbb{R}^{s \times d},$$

with the identity $p_i{}^\top p'_i = (|\mathcal{V}| + 1)^4 d \delta_s / (\epsilon_s \gamma_{\min})$. With this, we have the bound for p_i, p'_i :

$$\|p_i\| = \mathcal{O}\left(|\mathcal{V}|^2 \sqrt{d \frac{\delta_s}{\epsilon_s \gamma_{\min}}}\right), \quad \|p'_i\| = \mathcal{O}\left(|\mathcal{V}|^2 \sqrt{d \frac{\delta_s}{\epsilon_s \gamma_{\min}}}\right). \quad (\text{H.24})$$

Summing over the set of $p_i{}^\top p'_i$ for $i = 1, \dots, \rho$, we obtain the bound for rank ρ matrix W_Q and W_K

$$\begin{aligned} \|W_Q\|_2 &= \sup_{\|x\|_2=1} \|W_Q x\|_2 \leq C_Q = \mathcal{O}\left(\sqrt{\rho} |\mathcal{V}|^2 \sqrt{d \frac{\delta_c}{\epsilon_c \gamma_{\min}}}\right), \\ \|W_Q\|_{2,\infty} &= \max_{1 \leq i \leq d} \|(W_Q)_{(i,:)}\|_2 \leq C_Q^{2,\infty} = \mathcal{O}\left(\rho |\mathcal{V}|^2 \sqrt{d \frac{\delta_s}{\epsilon_s \gamma_{\min}}}\right), \\ \|W_K\|_2 &= \sup_{\|x\|_2=1} \|W_K x\|_2 \leq C_K = \mathcal{O}\left(\sqrt{\rho} |\mathcal{V}|^2 \sqrt{d \frac{\delta_s}{\epsilon_s \gamma_{\min}}}\right), \\ \|W_K\|_{2,\infty} &= \max_{1 \leq i \leq d} \|(W_K)_{(i,:)}\|_2 \leq C_K^{2,\infty} = \mathcal{O}\left(\rho |\mathcal{V}|^2 \sqrt{d \frac{\delta_s}{\epsilon_s \gamma_{\min}}}\right), \end{aligned}$$

where $\rho \leq s$ and the head size $s \leq d$.

After the first step quantization, we obtain vocabulary bounds $|\mathcal{V}| = \mathcal{O}(D^{dL})$ and output sequences with $(1/D, \sqrt{d}, 1/D)$ tokenwise separatedness. Also, in [Theorem H.2](#) we take $\delta_s = 4 \log L$ so that $f^{(\text{SA})}$ is a contextual mapping.

Next, by [Lemma H.4](#), we need $D = \mathcal{O}(\epsilon^{1/(dL)})$ for [Theorem H.2](#) to hold.

Combining all the components, we have the bounds for W_Q and W_K

$$\begin{aligned} \|W_Q\|_2, \|W_K\|_2 &= \mathcal{O}\left(d D^{2dL+1} (\log L)^{\frac{1}{2}}\right) = \mathcal{O}\left(d \epsilon^{\frac{2dL+1}{dL}} (\log L)^{\frac{1}{2}}\right), \\ \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} &= \mathcal{O}\left(d^{\frac{3}{2}} D^{2dL+1} (\log L)^{\frac{1}{2}}\right) = \mathcal{O}\left(d^{\frac{3}{2}} \epsilon^{\frac{2dL+1}{dL}} (\log L)^{\frac{1}{2}}\right) \end{aligned}$$

• **Bounds for W_O and W_V in $f^{(\text{SA})}$.**

Following the construction of W_Q and W_K in [Theorem H.2](#), we have the relation for W_V and W_O as

$$\begin{aligned} W_V &= \sum_{i=1}^{\rho} p''_i q_i{}^\top \in \mathbb{R}^{s \times d}, \\ W_O &= \sum_{i=1}^{\rho} p'''_i p''_i{}^\top \in \mathbb{R}^{d \times s}, \end{aligned}$$

with the identity $\|p'''_i\| \lesssim \epsilon_s / (4\rho \gamma_{\max} \|p''_i\|)$ from [\(H.10\)](#), and $p''_i \in \mathbb{R}^s$ is any nonzero vector.

Along with the $(\gamma_{\min} = 1/D, \gamma_{\max} = \sqrt{d}, \epsilon_s = 1/D)$ separateness and taking $D = \mathcal{O}(\epsilon^{1/(dL)})$, we have the following bounds for W_V and W_O :

$$\|W_V\|_2 = \sup_{\|x\|_2=1} \|W_V x\|_2 \leq C_V = \mathcal{O}(\sqrt{\rho}),$$

$$\begin{aligned} 3240 \quad \|W_V\|_{2,\infty} &= \max_{1 \leq i \leq d} \|(W_V)_{(i,:)}\|_2 \leq C_V^{2,\infty} = \mathcal{O}(\rho), \\ 3241 \\ 3242 \quad \|W_O\|_2 &= \sup_{\|x\|_2=1} \|W_O x\|_2 \leq C_O = \mathcal{O}(\sqrt{\rho} \cdot \rho^{-1} \cdot \gamma_{\max}^{-1} \cdot \epsilon_s) = \mathcal{O}\left(d^{-1} \epsilon^{-\frac{1}{dL}}\right) \\ 3243 \\ 3244 \quad \|W_O\|_{2,\infty} &= \max_{1 \leq i \leq s} \|(W_O)_{(i,:)}\|_2 \leq C_O^{2,\infty} = \mathcal{O}(\rho \cdot \rho^{-1} \cdot \gamma_{\max}^{-1} \cdot \epsilon_s) = \mathcal{O}\left(d^{-\frac{1}{2}} \epsilon^{-\frac{1}{dL}}\right). \\ 3245 \\ 3246 \end{aligned}$$

3247 Note that we use the fact $\max \rho = d$ in the last two lines.

3248 • **Bounds for W_1 in f_1^{FF} .**

3249 In order to approximate the quantization in [Theorem H.2](#), we set up f_1^{FF} as in [\(H.16\)](#) where every entry of W_1 in the layer is bounded by $\mathcal{O}(1/\delta)$. Therefore we have

$$3250 \quad \|W_1\|_{2,\infty} \leq C_{F_1}^{2,\infty} = \mathcal{O}\left(\frac{\sqrt{d}}{\delta}\right), \quad (\text{H.25})$$

$$3251 \quad \|W_1\|_2 \leq \|W_1\|_F \leq C_{F_1} = \mathcal{O}\left(\frac{d}{\delta}\right), \quad (\text{H.26})$$

3252 where the bound for δ is given from [Lemma H.4](#). We set $\delta = \nu D^{-1}$ for some $\nu \in (0, 1)$ such that we have the bounds $\mathcal{O}(\sqrt{d}\epsilon^{1/(dL)})$ and $\mathcal{O}(d\epsilon^{1/(dL)})$ respectively.

3253 • **Bounds on W_2 in f^{FF} .**

3254 The bounds for $\|W_2\|_2, \|W_2\|_{2,\infty}$ in [\(H.20\)](#) follow the same argument as for W_1 , with the replacement of the largest element with the scaling factor R . So we have

$$3255 \quad \|W_2\|_{2,\infty} \leq C_{F_2}^{2,\infty} = \mathcal{O}(\sqrt{d}R), \quad (\text{H.27})$$

$$3256 \quad \|W_2\|_2 \leq C_{F_2} = \mathcal{O}(dR). \quad (\text{H.28})$$

3257 Again, by [Lemma H.4](#), we take $R = \mathcal{O}(D) = \mathcal{O}(\epsilon^{1/(dL)})$ such that we have the bounds $\mathcal{O}(\sqrt{d}\epsilon^{1/(dL)})$ and $\mathcal{O}(d\epsilon^{1/(dL)})$ respectively.

3258 • **Bounds on Positional Encoding Matrix E .**

3259 For $\|E^\top\|_2, \|E^\top\|_{2,\infty}$, following [\(Kajitsuka and Sato, 2024\)](#), it suffices to set the positional encoding:

$$3260 \quad E = \begin{pmatrix} 2\gamma_{\max} & 4\gamma_{\max} & \cdots & 2L\gamma_{\max} \\ \vdots & \vdots & \ddots & \vdots \\ 2\gamma_{\max} & 4\gamma_{\max} & \cdots & 2L\gamma_{\max} \end{pmatrix}.$$

3261 Since the ℓ_2 norm over every row is identical, it suffices to derive

$$3262 \quad \|E^\top\|_{2,\infty} = \left(\sum_{i=1}^L (2i\gamma_{\max})^2\right)^{\frac{1}{2}} = \left(4\gamma_{\max}^2 \frac{L(L+1)(2L+1)}{6}\right)^{\frac{1}{2}} = \mathcal{O}\left(\gamma_{\max} L^{\frac{3}{2}}\right).$$

3263 Recall that we have the relation $\gamma_{\max} = \sqrt{d}$ in the self-attention layer. Therefore, we have the following bound for encoding matrix E :

$$3264 \quad \|E^\top\|_{2,\infty} \leq C_E = \mathcal{O}(d^{1/2} L^{3/2}). \quad (\text{H.29})$$

3265 This completes the proof. \square

I PROOF OF THEOREM 3.1

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

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

I.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 I.1 (Approximating Clipped Multi-Index Gaussian Integral, Lemma A.8 of (Fu et al., 2024b)). Assume 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

$$\int_{\mathbb{R}^{d_x} \setminus B_x} \left| \left(\frac{\alpha_t x_0 - x}{\sigma_t} \right)^\kappa \right| \cdot p(x_0|y) \cdot \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|\alpha_t x_0 - x\|^2}{2\sigma_t^2}\right) dx_0 \leq \epsilon, \quad (\text{I.1})$$

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

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

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

Based on Lemma I.1, we have the following considerations:

- For each $x \in \mathbb{R}^{d_x}$, consider a bounded domain

$$B_{x,N} := \underbrace{\left[\frac{x - \sigma_t C(0, d_x) \sqrt{\beta \log N}}{\alpha_t}, \frac{x + \sigma_t C(0, d_x) \sqrt{\beta \log N}}{\alpha_t} \right]}_{\text{(I)}} \cap \underbrace{\left[-C(0, d_x) \sqrt{\beta \log N}, C(0, d_x) \sqrt{\beta \log N} \right]^{d_x}}_{\text{(II)}}, \quad (\text{I.2})$$

where $C(0, d_x)$ is some positive constant depending on d_x and N . Here, we pick $n = 0$ for $C(n, d_x)$ to reduce (I.1) to

$$p_t(x|y) = \int_{\mathbb{R}^{d_x} \setminus B_{x,N}} p(x_0|y) \cdot \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|\alpha_t x_0 - x\|^2}{2\sigma_t^2}\right) dx_0 \leq \epsilon = 1/N.$$

This motivates a polynomial expansion of (3.1) on $B_{x,N}$ with precision $1/N$.

- Uniformly discretize each dimension of $B_{x,N}$ into N segments. Note that while not necessary, it is possible to pick a $C(0, d_x)$ such that grids in $B_{x,N}$ are non-overlapping.
- Uniformly discretize each dimension of $[0, 1]^{d_y}$ into N segments of length $1/N$.

This discretization of domains leads to $N^{d_x+d_y}$ hypercubes on bounded domain $B_{x,N} \times [0, 1]^{d_y}$.

Remark I.2. For any $x \in \mathbb{R}^{d_x}$, we shorthand (I.2) with

$$B_{x,N} = \left[-C_x \sqrt{\log N}, C_x \sqrt{\log N}\right]^{d_x}, \quad (\text{I.3})$$

where C_x summarize all factors except $\sqrt{\log N}$ in all dimensions of $x \in \mathbb{R}^{d_x}$. Moreover, when content is clear, we suppress the notation dependence on d_x for (I.3). Namely, we use the notation $B_{x,N} = [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]$ and $B_{x,N} = [-C_x \sqrt{\log N}, C_x \sqrt{\log N}]^{d_x}$ interchangeably.

Remark I.3. Lemma I.1 ensures that we can approximate the Gaussian integral of any polynomial function of the form (I.1) on \mathbb{R}^{d_x} with the same integral on B_x to an arbitrary precision $0 < \epsilon < 1/e$. This motivate us to approximate functions on \mathbb{R}^{d_x} with polynomials evaluated at $x \in \mathbb{R}^{d_x}$ on $B_{x,N}$. A natural choice is through Taylor expansion around $x \in \mathbb{R}^{d_x}$, as the Hölder class assumption guarantees local smoothing behavior for our error analysis.

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

Definition I.1 (Normalization of $B_{x,N}$). Consider the clipping in Lemma I.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 I.4. The purpose of Definition I.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 I.5. 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 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 $M(x'_0, y) \in \mathcal{H}([0, 1]^{d_x+d_y}, BR_B^{k_1})$. This follows from the fact that $p(\cdot|y)$ is k_1 -time differentiable so that the radius scale by a factor of $R_B^{k_1}$.

Lemma I.2 (Diffused Local Polynomial, Modified from (Fu et al., 2024a)). Assume 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 in (I.3):

$$f_1(x, y, t) = \sum_{v \in [N]^d, w \in [N]^{d_y} \|\mathbf{n}_x\|_1 + \|\mathbf{n}_y\|_1 \leq k_1} \sum_{\substack{R_B^{\|\mathbf{n}_x\|} \\ n_x! n_y!}} \frac{\partial^{n_x+n_y} p}{\partial x^{n_x} \partial y^{n_y}} \Big|_{x=R_B(\frac{v}{N}-\frac{1}{2}), y=\frac{w}{N}} \Phi_{n_x, n_y, v, w}(x, y, t), \quad (\text{I.4})$$

where

- $\phi(\cdot)$ is the trapezoid function.
- $g(x, n_x, v, k_2) := \frac{1}{\sigma_t \sqrt{2\pi}} \int \left(\frac{x_0}{R} + \frac{1}{2} - \frac{v}{N} \right)^{n_x} \frac{1}{k_2!} \left(-\frac{|x - \sigma_t x_0|}{2\sigma_t^2} \right)^{k_2} dx_0$.
- $\Phi_{n_x, n_y, v, w}(x, y, t) := \left(y - \frac{w}{N} \right)^{n_y} \prod_{j=1}^{d_y} \phi \left(3N(y[j] - \frac{w}{N}) \right) \prod_{i=1}^{d_x} \sum_{k_2 < p} g(x[i], n_x[i], v[i], k_2)$.

Remark I.6. 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 I.7. The Hölder space assumption in [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.

Proof Sketch. We provide the proof overview of [Lemma I.2](#). with the following three steps.

Step A: Clip $\mathbb{R}^{d_x} \times [0, 1]^{d_y}$.

We clip the domain $\mathbb{R}^{d_x} \times [0, 1]^{d_y}$ into closed and bounded region $B_{x,N}$.

Step B: Replace $p(x_0|y)$ with k_1 -order Taylor Polynomials.

We discretize $[0, 1]^{d_x+d_y}$ into $N^{d_x+d_y}$ hypercubes. We apply Taylor expansion to each grid point. For areas not located on any grid point, we construct a trapezoid function and an indicator function to control the approximation error.

Step C: Replace $\exp(\cdot)$ with k_2 -order Taylor Polynomials.

We apply Taylor expansion to approximate regions within $B_{x,N}$ for $\exp(\cdot)$. Note that we leverage the explicit form of the exponential function to achieve accurate approximation without additional discretization as in previous step.

Step D: Altogether, the Diffused Local Polynomials.

We combine these 4 steps and construct *the diffused local polynomial* (I.4). \square

Proof of Lemma I.2. We demonstrate details regarding the three steps.

• **Step A: Clip** $\mathbb{R}^{d_x} \times [0, 1]^{d_y}$.

We take $\kappa[i] = 0$ for $i = [d_x]$ and set $\epsilon = N^{-\beta}$ in [Lemma I.1](#). This gives closed and bounded domain $B_{x,N}$ specified in (I.3) and clipping-induced error:

$$\left| p_t(x|y) - \int_{B_{x,N}} p(x_0|y) \cdot \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|\alpha_t x_0 - x\|^2}{2\sigma_t^2}\right) dx_0 \right| \leq N^{-\beta}. \quad (\text{I.5})$$

• **Step B: Replace $p_0(x_0|y)$ with k_1 -order Taylor Expansion.**

We construct a approximator $Q(x'_0, y)$ for $M(x'_0, y)$ with domain $[0, 1]^{d_x+d_y}$.⁵ At the end of this step, we reset $x'_0 = x_0/R_B + 1/2$ in $Q(x'_0, y)$ as the final approximator of $p(x_0|y)$.

⁵Recall $R_B := (2C(0, d)\sqrt{\beta \log N})$, $x'_0 := x_0/R_B + 1/2$, and $M(x'_0, y) := p(R_B(x'_0 - 1/2)|y)$ from [Definition I.1](#).

– **Step B.1: Discretize** $[0, 1]^{d_x+d_y}$.

We uniformly discretize $[0, 1]^{d_x+d_y}$ into grid points $[0, 1/N, 2/N, \dots, (N-1)/N, 1]^{d_x+d_y}$.

– **Step B.2: Implement Taylor Expansion.**

We construct the k_1 -order Taylor polynomial $P_{v,w}(x, y)$ at point $(v/N, w/N)$ for $M(x'_0, y)$:⁶

$$P_{v,w}(x'_0, y) := \sum_{\|n_x\|_1 + \|n_y\|_1 \leq k_1} \frac{1}{n_x! n_y!} \frac{\partial^{n_x+n_y} M}{\partial x^{n_x} \partial y^{n_y}} \Bigg|_{x'_0 = \frac{v}{N}, y = \frac{w}{N}} \left(x'_0 - \frac{v}{N}\right)^{n_x} \left(y - \frac{w}{N}\right)^{n_y}. \quad (\text{I.6})$$

For x'_0 and y not located on any grid point, we construct an indicator function that ensures $\|x'_0 - v/N\|_\infty < 1/N$ and $\|y - w/N\|_\infty < 1/N$ in the next step. For now, we assume these conditions hold.

To analyze the error, we expand the target function $M(x'_0, y)$. By Taylor's theorem, there exist $\theta_x \in [0, 1]^{d_x}$ and $\theta_y \in [0, 1]^{d_y}$ such that

$$\begin{aligned} M(x'_0, y) &= \sum_{\|n_x\|_1 + \|n_y\|_1 < k_1} \frac{1}{n_x! n_y!} \cdot \frac{\partial^{n_x+n_y} M}{\partial x_0^{n_x} \partial y^{n_y}} \Bigg|_{x'_0 = \frac{v}{N}, y = \frac{w}{N}} \left(x'_0 - \frac{v}{N}\right)^{n_x} \left(y - \frac{w}{N}\right)^{n_y} \\ &+ \sum_{\|n_x\|_1 + \|n_y\|_1 = k_1} \frac{1}{n_x! n_y!} \cdot \frac{\partial^{n_x+n_y} M}{\partial x_0^{n_x} \partial y^{n_y}} \Bigg|_{x'_0 = x_1, y = y_1} \left(x'_0 - \frac{v}{N}\right)^{n_x} \left(y - \frac{w}{N}\right)^{n_y}, \end{aligned}$$

where $x_1 = (1 - \theta_x)v/N + \theta_x x'_0$ and $y_1 = (1 - \theta_y)w/N + \theta_y y$. This ensures x_1 lies between x'_0 and v/N , and y_1 lies between y and w/N .

Note that the difference between $P_{v,w}(x'_0, y)$ and $M(x'_0, y)$ stems from the different value taken in $\partial^{n_x+n_y} M / (\partial x_0^{n_x} \partial y^{n_y})$ for all terms in the series with $\|n_x\|_1 + \|n_y\|_1 = k_1$.

To study the error, let $z = (x'_0, y)$ and recall from the definition of Hölder norm (Definition 3.1):

$$\max_{\alpha: \|\alpha\|_1 = k_1} \sup_{z \neq z'} \frac{|\partial^{k_1} M(z) - \partial^{k_1} M(z')|}{\|z - z'\|_\infty^\gamma} < \|M(x'_0, y)\|_{\mathcal{H}^\beta([0,1]^{d_x+d_y})} < R_B^{k_1} B. \quad (\text{I.7})$$

We rewrite the error as

$$\begin{aligned} &|P_{v,w}(x'_0, y) - M(x'_0, y)| \\ &\leq \sum_{\|n_x\|_1 + \|n_y\|_1 = k_1} \frac{1}{n_x! n_y!} \left(x'_0 - \frac{v}{N}\right)^{n_x} \left(y - \frac{w}{N}\right)^{n_y} \underbrace{\left| \frac{\partial^{n_x+n_y} M}{\partial x_0^{n_x} \partial y^{n_y}} \Bigg|_{x'_0 = x_1, y = y_1} - \frac{\partial^{n_x+n_y} M}{\partial x_0^{n_x} \partial y^{n_y}} \Bigg|_{x'_0 = \frac{v}{N}, y = \frac{w}{N}} \right|}_{\text{Apply Hölder Regularity}} \\ &\leq \sum_{\|n_x\|_1 + \|n_y\|_1 = k_1} \frac{1}{n_x! n_y!} \left(x'_0 - \frac{v}{N}\right)^{n_x} \left(y - \frac{w}{N}\right)^{n_y} \underbrace{\|M(x'_0, y)\|_{\mathcal{H}^\beta([0,1]^{d_x+d_y})}}_{(\text{I.7})} \underbrace{\left\| [\theta_x x'_0, \theta_y y] - \frac{1}{N} [\theta_x v, \theta_y w] \right\|_\infty^\gamma}_{\text{Controlled by indicator function (I.8)}} \\ &\leq \sum_{\|n_x\|_1 + \|n_y\|_1 = k_1} \frac{BR_B^{k_1}}{n_x! n_y! N^{\|n_x\|_1 + \|n_y\|_1 + \gamma}} = \frac{BR_B^{k_1} (d_x + d_y)^{k_1}}{N^\beta k_1!}. \end{aligned}$$

– **B.3: Control Error for the Off-Grid Regions.**

⁶Please see Remarks I.4 and I.5 for details.

For regions not located on any grid point $(v/N, w/N)$, we construct an indicator function $\psi(x'_0, y)$ to ensure that our Taylor approximation at $(v/N, w/N)$ does not deviate from (x'_0, y) by more than $1/N$ in ℓ_∞ distance.

Specifically, we define

$$\psi_{v,w}(x'_0, y) := \mathbb{1} \left\{ x'_0 \in \left(\frac{v-1}{N}, \frac{v}{N} \right] \right\} \prod_{j=1}^{d_y} \phi \left(3N \left(y[j] - \frac{w}{N} \right) \right), \quad (\text{I.8})$$

where $\phi(\cdot)$ is the trapezoid function:

$$\phi(\tau) = \begin{cases} 1, & |\tau| < 1 \\ 2 - |\tau|, & |\tau| \in [1, 2] \\ 0, & |\tau| > 2. \end{cases}$$

Note that, $\psi_{v,w}$ is nonzero if and only if $x'_0 \in [(v-1)/N, v/N]$ and $y[j] \in [(w[j] - 2/3)/N, (w[j] + 2/3)/N]$ for $j \in [d_y]$. This guarantees $\|x'_0 - v/N\|_\infty \leq 1/N$ and $\|y - w/N\|_\infty \leq 1/N$.

– **Step B.4: Construct the Final Approximator for $p(x_0|y)$.**

Combining (I.6) and (I.8), we obtain an approximator of the form:

$$Q(x'_0, y) = \sum_{v,w} \psi_{v,w}(x, y) P_{v,w}(x'_0, y).$$

Since for all $x \in (0, 1]^{d_x}$ and $y \in [0, 1]^{d_y}$ the indicator function $\psi_{v,w}(x'_0, y)$ sums to 1, it holds:

$$|M(x'_0, y) - Q(x'_0, y)| \leq \frac{BR^{k_1}(d_x + d_y)^{k_1}}{k_1!N^\beta}. \quad (\text{I.9})$$

We conclude this step with the approximator $Q(x'_0, y) = Q(x_0/R_B + 1/2, y)$ for $p(x_0|y)$.

• **Step C: Replace $\exp(\cdot)$ with k_2 -order Taylor Expansion.**

Recall that we set $B_{x,N}$ as

$$B_{x,N} = \left[\frac{x - \sigma_t C(0, d_x) \sqrt{\beta \log N}}{\alpha_t}, \frac{x + \sigma_t C(0, d_x) \sqrt{\beta \log N}}{\alpha_t} \right] \\ \cap \left[-C(0, d_x) \sqrt{\beta \log N}, C(0, d_x) \sqrt{\beta \log N} \right]^{d_x}.$$

This gives $|(x[i] - \alpha_t x_0[i])/\sigma_t| \leq C(0, d_x) \sqrt{\beta \log N}$ for any $i \in [d_x]$ and $x_0 \in B_{x,N}$.

Furthermore, we have

$$\|(x - \alpha_t x_0)/\sigma_t\|^2 = \sum_{i=1}^{d_x} |(x[i] - \alpha_t x_0[i])/\sigma_t|^2 \leq d_x \cdot \left(C(0, d_x) \sqrt{\beta \log N} \right)^2. \quad (\text{I.10})$$

From this fact, we implement the k_2 -order Taylor expansion to $\exp\left(-\|(x - \alpha_t x_0)/\sigma_t\|^2/2\right)$:

$$\left| \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) - \sum_{k_2 < u} \frac{1}{k_2!} \left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right)^{k_2} \right| \quad (\text{By Taylor theorem})$$

$$\begin{aligned}
&\leq \frac{1}{u!2^u} \left(\left\| \frac{x - \alpha_t x_0}{\sigma_t} \right\|^2 \right)^u \\
&= \frac{1}{u!2^u} \left(\sum_{i=1}^{d_x} |(x[i] - \alpha_t x_0[i]) / \sigma_t|^2 \right)^u \\
&\leq \frac{1}{u!2^u} \left(d_x \cdot \left(C(0, d) \sqrt{\beta \log N} \right)^2 \right)^u.
\end{aligned}$$

for all $x_0 \in B_{x,N}$, and u is a positive real number.

Following the choice of u from (Fu et al., 2024b), by utilizing the inequality $u! \geq (u/3)^u$ for $u \geq 3$ and setting

$$u := \max \left(\frac{2}{3} C^2(0, d) \beta^2 e \log N, \beta \log N + \log d_x \right),$$

we further write the bound as:

$$\left| \exp \left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2} \right) - \sum_{k_2 < u} \frac{1}{k_2!} \left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2} \right)^{k_2} \right| \lesssim N^{-\beta}. \quad (\text{I.11})$$

• **Step D: The Diffused Local Polynomial.**

Substituting $p(x_0|y)$ and $\exp(\cdot)$ with their respective approximator in (I.9) and (I.11), we obtain the following expression:

$$f_1(x, y, t) = \frac{1}{\sigma_t^{d_x} (2\pi)^{\frac{d_x}{2}}} \int_{B_{x,N}} Q \left(\frac{x_0}{R_B} + \frac{1}{2}, y \right) \sum_{k_2 < u} \frac{1}{k_2!} \left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2} \right)^{k_2} dx_0. \quad (\text{I.12})$$

We term f_1 as *diffused local polynomial*, following (Fu et al., 2024a).

⁷Rearranging (I.12), we obtain the form

$$f_1(x, y, t) = \sum_{v \in [N]^d, w \in [N]^{d_y}} \sum_{\|n_x\|_1 + \|n_y\|_1 \leq k_1} \frac{R_B^{\|n_x\|}}{n_x! n_y!} \frac{\partial^{n_x + n_y} f}{\partial x^{n_x} \partial y^{n_y}} \Big|_{x = \frac{v}{N}, y = \frac{w}{N}} \Phi_{n_x, n_y, v, w}(x, y, t), \quad (\text{I.13})$$

where

$$\begin{aligned}
- g(x, n_x, v, k_2) &:= \frac{1}{\sigma_t \sqrt{2\pi}} \int \left(\frac{x_0}{R} + \frac{1}{2} - \frac{v}{N} \right)^{n_x} \frac{1}{k_2!} \left(-\frac{\|x - \sigma_t x_0\|^2}{2\sigma_t^2} \right)^{k_2} dx_0. \\
- \Phi_{n_x, n_y, v, w}(x, y, t) &:= \left(y - \frac{w}{N} \right)^{n_y} \prod_{j=1}^{d_y} \phi \left(3N(y[j] - \frac{w}{N}) \right) \prod_{i=1}^{d_x} \sum_{k_2 < p} g(x[i], n_x[i], v[i], k_2).
\end{aligned}$$

This completes the proof. \square

We specifies the error from the approximation of p_t and ∇p_t with f_1 and f_2 in Lemmas I.3 and I.4.

Lemma I.3 (Approximation of $p_t(x|y)$ by Polynomials, Lemma A.4 of (Fu et al., 2024b)). Assume 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

⁷Further details regarding the derivation are in (Fu et al., 2024b, Appendix A.4).

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 I.4 (Approximation of $\nabla \log p_t(x|y)$ by Polynomials, Lemma A.6 of (Fu et al., 2024b)). Assume [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 $f_2 := (f_2[1], \dots, f_2[d_x])^\top \in \mathbb{R}^{d_x}$ with local diffused polynomial $f_2[i]$ such that

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

where each $f_2[i]$ contains at most $N^{d_x+d_y} (d_x + d_y)^{k_1}$ monomials.

We have finished the approximation of p_t and ∇p_t with diffused local polynomial f_1 and f_2 .

Step 3. Approximate Diffused Local Polynomials and Algebraic Operators with Transformers.

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 necessary in final score approximation. Lastly, we present the incorporation of these components in [Lemma I.13](#) with a unified transformer architecture and corresponding parameter configuration.

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

We invoke the universal approximation theorem of transformer ([Theorem H.2](#)). We utilize network consisting of one transformer block and one feed-forward layer (see [Figure 1](#) and [Definition 2.2](#)).

Lemma I.5 (Approximate Scalar Polynomials with Transformers). Assume [Assumption 3.1](#). Consider the diffused local polynomial f_1 in [Lemma I.3](#). 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

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

The parameter bounds in the Transformer network class satisfy

$$\begin{aligned} \|W_Q\|_2, \|W_K\|_2 &= \mathcal{O}\left(d\epsilon^{-\frac{2dL+4d+1}{d}} (\log L)^{\frac{1}{2}}\right); \\ \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} &= \mathcal{O}\left(d^{\frac{3}{2}} \epsilon^{-\frac{2dL+4d+1}{d}} (\log L)^{\frac{1}{2}}\right); \\ \|W_V\|_2 &= \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \\ \|W_O\|_2 &= \mathcal{O}\left(\sqrt{d}\epsilon^{\frac{1}{d}}\right); \|W_O\|_{2,\infty} = \mathcal{O}\left(\epsilon^{\frac{1}{d}}\right); \\ \|W_1\|_2 &= \mathcal{O}\left(d\epsilon^{-\frac{1}{d}} \cdot \log N\right); \|W_1\|_{2,\infty} = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}} \cdot \log N\right); \\ \|W_2\|_2 &= \mathcal{O}\left(d\epsilon^{-\frac{1}{d}}\right); \|W_2\|_{2,\infty} = \mathcal{O}\left(\sqrt{d}\epsilon^{-\frac{1}{d}}\right); \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}} L^{\frac{3}{2}}\right). \end{aligned}$$

Proof of Lemma I.5. We first skip the embedded dimension of y and t for the following proof without loss of generality. We put it back at the end of the derivation, by replacing L with $L + 2$.

To implement a sequence-to-sequence model for approximating a function that outputs a scalar, we define a trivial function for converting the scalar target into a sequence represented by matrices.

To begin with, for $x \in \mathbb{R}^{d_x}$ and $f_1 : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$, we define a trivial function:

$$F_1(x) := \underbrace{(\alpha_1 f_1(x), \alpha_2 f_1(x), \dots, \alpha_{d_x-1} f_1(x), f_1(x))^\top}_{(\text{padding } d_x - 1 \text{ elements})} \in \mathbb{R}^{d_x},$$

for any set of non-repeated constants $\{\alpha_i\}_{i=1}^{d_x-1} \in \mathbb{R} \setminus \{1\}$.

3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725

By *trivial*, we mean that F_1 transforms $f_1(x) \in \mathbb{R}$ into a vector $F_1(x) \in \mathbb{R}^{d_x}$ where only the last entry is meaningful.

In order to apply the universal approximation of transformers in [Theorem H.2](#), we show the uniform continuity of F_1 as follows.

– **Step A: Uniform Continuity.**

For different input x, x' , we start by writing

$$\begin{aligned} \|F_1(x) - F_1(x')\|_p &= \left\{ |f(x) - f(x')|^p + \sum_{i=1}^{d_x-1} |\alpha_i f(x) - \alpha_i f(x')|^p \right\}^{1/p} \\ &= \left\{ |f(x) - f(x')|^p \left(1 + \sum_{i=1}^{d_x-1} |\alpha_i|^p \right) \right\}^{1/p} \\ &= \eta |f(x) - f(x')|, \end{aligned}$$

where $\eta = \left(1 + \sum_{i=1}^{d_x-1} |\alpha_i|^p \right)^{1/p} \in \mathbb{R}_+$.

Next, we utilize the fact that the diffused local polynomials f_1 is continuous on compact support. That is, for all $\epsilon > 0$, there exists $\delta > 0$ such that for all x and x' , if $\|x - x'\|_\infty < \delta$, then $|f_1(x) - f_1(x')| < \epsilon$.

From this fact, by taking $\epsilon = \epsilon' / \eta$, we have that for all $\epsilon' > 0$, there exists $\delta' > 0$ such that for all x and x' , if $\|x - x'\|_\infty < \delta'$, then $|f_1(x) - f_1(x')| < \epsilon' = \epsilon \eta$.

This gives $\|F_1(x) - F_1(x')\|_p \leq \epsilon'$ and therefore we obtain the uniform continuity for F_1 .

Also, the reshape layer $R(\cdot)$ that converts $x \in \mathbb{R}^{d_x}$ into sequential input $R(x) \in \mathbb{R}^{d \times L}$ does not harm this continuity due to its linearity. Therefore, the map $R \circ F_1(x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d \times L}$ is also uniformly continuous.

– **Step B: Universal Approximation.**

We apply [Theorem H.2](#) that guarantees for any $\epsilon_{f_1} > 0$, there exists one transformer block and one feed-forward layer such that

$$\|R \circ F_1 - f^{h,s,r} \circ f^{\text{FF}} \circ R\|_p \leq \epsilon_{f_1}.$$

Adding a reverse reshape layer, we have $\mathcal{T}_{f_1} = R^{-1} \circ f^{h,s,r} \circ f^{\text{FF}} \circ R$ with $\|F_1 - \mathcal{T}_{f_1}\|_p \leq \epsilon_{f_1}$.

Next, observe that

$$|\mathcal{T}_{f_1}[d_x] - f_1| \leq \left\{ \sum_{i=1}^{d_x} |\mathcal{T}_{f_1}[i] - \alpha_i f_1|^p \right\}^{1/p} = \|\mathcal{T}_{f_1} - F_1\|_p \leq \epsilon_{f_1}, \quad (\text{I.14})$$

with $\alpha_{d_x} = 1$. [\(I.14\)](#) completes the proof of the approximation error.

– **Step C: Parameter Bounds.**

To establish the approximation [\(I.14\)](#), we need the parameter bounds in [Lemma H.5](#) to hold. This requires transforming the input domain from $[-C_x \sqrt{\log N}, C_x \sqrt{\log N}]$ to normalized compact support $[0, 1]$ for all dimensions (i.e., $x[i]$ for all $i \in [d_x]$.)

Recall that [\(H.25\)](#), we have bound for W_1 :

$$\|W_1\|_{2,\infty} = \mathcal{O}(\sqrt{dD}) = \mathcal{O}(\sqrt{d}\epsilon^{-dL}), \quad (\text{I.15})$$

$$\|W_1\|_2 = \mathcal{O}(dD) = \mathcal{O}(d\epsilon^{-dL}), \quad (\text{I.16})$$

that is, the bounds on each element in W_1 scales up as the granularity increases. Because for a fixed precision level, the granularity is proportional to the length of the interval in each dimension of the input domain, we conclude that $\|W_1\|_2 = \mathcal{O}(d\epsilon^{-dL} \log N)$ and $\|W_1\|_{2,\infty} = \mathcal{O}(\sqrt{d}\epsilon^{-dL} \log N)$.

The rest of bounds for each operation follows [Lemma H.5](#). Lastly, we incorporate the embedded dimensions of y and t by replacing L with $L + 2$ (see [Figure 1](#)).

This completes the proof. \square

Similarly, we have the corresponding $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$ for the approximation of $f_2(x, y, t)$.

Lemma I.6 (Approximate Vector-Valued Polynomials with Transformers). Assume [Assumption 3.1](#) and consider $f_2(x, y, t) \in \mathbb{R}^{d_x}$ with every entry $f_2[1], \dots, f_2[d_x]$ is a local diffused polynomial defined in [Lemma I.2](#). For any $\epsilon > 0$, there exists a transformer $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$ such that

$$\|f_2(x, y, t) - \mathcal{T}_{f_2}\|_\infty \leq \epsilon,$$

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 parameter bounds in the transformer network class follows [Lemma I.5](#).

Proof of Lemma I.6. Since each entry of the diffused local polynomials in f_2 is continuous on compact support, $f_2 \in \mathbb{R}^{d_x}$ is uniformly continuous by the same argument as in the proof of [Lemma I.5](#).

Similarly, by [Theorem H.2](#), for any $\epsilon_{f_2} > 0$, there exists a transformer block and a feed-forward layer such that $\|R \circ f_2 - f^{h,s,r} \circ f^{\text{FF}} \circ R\|_p \leq \epsilon_{f_2}$.

By adding the reversed reshape layer, we obtain $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$, satisfying $\|f_2 - \mathcal{T}_{f_2}\|_p \leq \epsilon_{f_2}$.

Then we have,

$$|\mathcal{T}_{f_2}[j] - f_2[j]| \leq \left\{ \sum_{j=1}^{d_x} |\mathcal{T}_{f_2}[j] - f_2[j]|^p \right\}^{1/p} \leq \epsilon_{f_2}$$

for all $j = 1, \dots, d_x$. Thus the result with ℓ_∞ bound also holds.

The network configuration follows the argument as in the proof of [Lemma I.5](#).

This completes the proof. \square

So far, we have obtained approximation results for f_1 and f_2 . To complete the full approximation of the score decomposition $\nabla \log p = \frac{\nabla p}{p}$, we still need to approximate several key algebraic operators, including the product ([Lemma I.8](#)), inverse ([Lemma I.9](#))...etc.

We establish their approximations as follows.

• **Step 3.2: Approximate Algebraic Operators with Transformers.**

We give transformer approximation theory for the clipping operator, the inverse operator, the product operator, and functions that evolve with time t :

- Clipping operation ([Lemma I.7](#))
- Product operation ([Lemma I.8](#))
- Inverse operation ([Lemma I.9](#))
- Mean $\alpha_t = \exp(-t/2)$ ([Lemma I.10](#))

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

In their works, the general network structure is as follows.

Definition I.2. A family of fully-connected neural networks with length L , width W , sparsity constraint S , and norm constraint B is defined as:

$$\Phi(L, W, S, B) := A^{(L)}\text{ReLU}(\cdot) + b^{(L)} \circ \dots \circ A^{(1)}x + b^{(1)},$$

where $A^{(i)}$ and $b^{(i)}$ represent the matrix operator and bias in the i -th layer. Specifically:

- Length: $L \in \mathbb{R}$ denotes the number of hidden layers plus one.
- Width: $W \in \mathbb{N}^{L+1}$ is a vector representing the output dimension of each layer.
- Sparsity Constraint: $\sum_{i=1}^L \|A^{(i)}\|_{0,0} + \|b^{(i)}\|_0 \leq S$ specifies the maximum number of non-zero terms.
- Norm Constraint: $\max_{1 \leq i \leq L} \|A^{(i)}\|_{\infty, \infty} \vee \|b^{(i)}\|_{\infty} \leq B$ specifies the upper bound on the parameter norms.

Here \vee denotes the maximum of two values.

Remark I.8 (Generalization ReLU Networks with Transformers). Transformers are more general network class that encompasses ReLU-based networks defined in Definition I.2. By setting all self-attention layers in the transformer to identity maps, we recover the ReLU feed-forward network structure. Therefore, our work on approximating with transformers extends previous works Fu et al. (2024b); Oko et al. (2023) by incorporating the flexibility of self-attention mechanisms.

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

Lemma I.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

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

with

$$L = 2, \quad W = (d, 2d, d)^\top, \quad S = 7d, \quad B = \max_{1 \leq i \leq d} \max(|a[i]|, b[i]). \quad (\text{I.17})$$

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

Proof. It suffices to show the result for i -th coordinate, and implement the parallelization to complete the proof that holds for the entire vector $\phi_{\text{clip}}(x; a, b)$.⁸ The clipping operation yields the middle among $a[i]$, $b[i]$ and the input $x[i]$. Following (Oko et al., 2023), we achieve the task by setting:

$$\min(b[i], \max(x[i], a[i])) = \text{ReLU}(x[i] - a[i]) - \text{ReLU}(x[i] - b[i]) + a[i].$$

Note that the RHS is realized by the network with one hidden layer:

$$(1, -1)\text{ReLU} \left((1, 1)x[i] + \begin{pmatrix} -a[i] \\ -b[i] \end{pmatrix} \right) + a[i],$$

⁸For a more detailed description regarding parallelization please see Appendix F of (Oko et al., 2023).

with 7 non-zero parameters, and the scale of parameter is $\max(|a[i]|, b[i])$. So there exists $\phi_{\text{clip}}(x[i]; a[i], b[i]) \in \Phi(2, (1, 2, 1)^\top, 7, \max(|a[i]|, b[i]))$ executing the clipping operation. Then the proof is complete by the parallelization for all the components $i = 1, \dots, d$.

This completes the proof. \square

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

Lemma I.8 (Approximation of the Product Operator with Transformer.). Let $m \geq 2$ and $C \geq 1$. 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 $\|x - x'\|_\infty \leq \epsilon_{\text{error}}$, it holds

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

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

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

Proof. We build our proof on (Oko et al., 2023, Lemma F.6).

Unlike approximation for input $x \in [-C_x\sqrt{\log N}, C_x\sqrt{\log N}]^{d_x}$ in Lemma I.5, the input dimension for the product operator is sufficiently smaller so that we skip the reshape layer by setting R and R^{-1} as identity map.

Next, let $f(x) = \prod_{i=1}^m x[i]$, and define a trivial function $F(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}^{1 \times m}$ as

$$F(x) := \underbrace{(\alpha_1 f(x), \alpha_2 f(x), \dots, \alpha_{m-1} f(x))}_{\text{(padding } m-1 \text{ elements)}}, f(x) \in \mathbb{R}^{1 \times m}.$$

The idea of padding a scalar into a row vector again stems from the purpose of utilizing sequence-to-sequence model to approximate functions that output a scalar.

By the same argument as in the proof of Lemma I.5, the uniform continuity of f guarantees the uniform continuity of F with respect to the L_p norm.

By Theorem H.2, for any $\epsilon > 0$, there exist $\mathcal{T}_{\text{mult}} \in \mathcal{T}_R^{h,s,r}$ with R, R^{-1} being identity map such that

$$\|\mathcal{T}_{\text{mult}} - F\|_p \leq \epsilon.$$

Clearly, $|\mathcal{T}_{\text{mult}}[m] - F[m]| \leq \|\mathcal{T}_{\text{mult}} - F\|_p \leq \epsilon$.

To extend the input to $x' \in \mathbb{R}^m$ with $\|x - x'\| \leq \epsilon_{\text{error}}$, we adopt Lemma I.7 and write

$$\begin{aligned} & \left| C^m \mathcal{T}_{\text{mult}}(\phi_{\text{clip}}(x'; -C, C)/C) - \prod_{i=1}^m x[i] \right| \\ & \leq \left| C^m \mathcal{T}_{\text{mult}}(\phi_{\text{clip}}(x'; -C, C)/C) - \prod_{i=1}^m \min(C, \max(x'[i], -C)) \right| + \left| \prod_{i=1}^m \min(C, \max(x'[i], -C)) - \prod_{i=1}^m x[i] \right| \\ & \leq C^m C^{-m} \epsilon + C^{m-1} \sum_{i=1}^m |x[i] - \min(C, \max(x'[i], -C))| \\ & = \epsilon + mC^{m-1}\epsilon_{\text{error}}. \end{aligned}$$

Further details regarding the product approximation are in Appendix F.2 of (Oko et al., 2023).

For the parameter bounds, following the same argument in the proof of Lemma I.5, it suffices to take $\mathcal{O}(C\epsilon^{-1})$ for W_1 . The rest of bounds for each operation follows Lemma H.5 with $d = 1$ and $L = m$.

This completes the proof. \square

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

Lemma I.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

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

The parameter bounds in the Transformer network class satisfy

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

Proof. We build our proof on (Oko et al., 2023, Lemma F.7). For any $\epsilon_{\text{rec}} \in (0, 1)$, since $1/x$ is continuous on $x \in [\epsilon_{\text{rec}}, \epsilon_{\text{rec}}^{-1}]$, by Theorem H.2, there exist a transformer $\mathcal{T}_{\text{rec}} \in \mathcal{T}_R^{h,s,r}$ such that

$$\left| \mathcal{T}_{\text{rec}} - \frac{1}{x} \right| \leq \epsilon_{\text{rec}}.$$

Extending to network with input $x' \in \mathbb{R}$, the sensitivity analysis follows:

$$\left| \mathcal{T}_{\text{rec}}(x') - \frac{1}{x} \right| \leq \left| \mathcal{T}_{\text{rec}}(x') - \frac{1}{\max(x', \epsilon)} \right| + \left| \frac{1}{x} - \frac{1}{\max(x', \epsilon)} \right|.$$

This yields the result.

For the parameter bounds, by the same discussion in the proof of Lemma I.8, we scale W_1 up by ϵ_{rec} such that the quantization in (H.16) works on normalized $[0, 1]$. The rest of the bounds follow Lemma H.5.

This completes the proof. \square

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

Lemma I.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

$$|\mathcal{T}_\alpha(t) - \alpha_t| \leq \epsilon_\alpha.$$

The parameter bounds in the Transformer network class satisfy

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

Proof. We build our proof on (Fu et al., 2024b, Lemma F.8). The proof consists of four steps.

3942 – **Step A: Approximate $\exp(\cdot)$ with Taylor polynomial for $t \in [0, T]$.**

3943 By Taylor theorem, there exist some $\theta \in [0, T]$ such that

$$3944 \exp\left(-\frac{t}{2}\right) = \sum_{i=0}^{s-1} \frac{(-1)^i}{i!} \left(\frac{t}{2}\right)^i + \frac{(-1)^s}{s!} \left(\frac{\theta}{2}\right)^s \exp\left(-\frac{\theta}{2}\right).$$

3949 We further bound the error from the remainder by

$$3951 \left| \exp\left(-\frac{t}{2}\right) - \sum_{i=0}^{s-1} \frac{(-1)^i}{i!} \left(\frac{t}{2}\right)^i \right| \leq \frac{T^s}{2^s s!}, \quad (\text{I.18})$$

3954 with T and s to be chosen later.

3956 – **Step B: Approximate Taylor polynomial with transformer for $t \in [0, T]$.**

3957 We take t as a sequence with length 1 and one-dimensional token.

3958 For $t \in [0, T]$, Taylor polynomial is a continuous function with compact support.

3959 Therefore, by [Theorem H.2](#), for any ϵ there exist a transformer $\mathcal{T}'_\alpha \in \mathcal{T}_R^{h,s,r}$ such that

$$3962 \left| \mathcal{T}'_\alpha - \sum_{i=1}^{s-1} \frac{(-1)^i}{i!} \left(\frac{t}{2}\right)^i \right| \leq \epsilon. \quad (\text{I.19})$$

3966 – **Step C: Extend the two approximation results from Step 1. and Step 2. to $t > T$.**

3967 We define \mathcal{T}_α as

3968 (i) $\mathcal{T}_\alpha(t) = \mathcal{T}'_\alpha(t)$ for $t \in [0, T]$.

3969 (ii) $\mathcal{T}_\alpha(t) = \mathcal{T}'_\alpha(T)$ for $t \geq T$.

3970 Next, we bound the error for $t > T$ by

$$3973 \left| \exp\left(-\frac{t}{2}\right) - \mathcal{T}_\alpha(t) \right| \leq \left| \exp\left(-\frac{T}{2}\right) - \exp\left(-\frac{t}{2}\right) \right| + \left| \mathcal{T}_\alpha(t) - \exp\left(-\frac{T}{2}\right) \right|. \quad (\text{I.20})$$

3974 – **Step D: Select T , s and transformer approximation error such that the result holds for all $t \geq 0$.**

3975 For any $\epsilon_\alpha > 0$, we ensure $|\mathcal{T}_\alpha - \exp(-t/2)| \leq \epsilon_\alpha$ holds for all $t \geq 0$.

3976 To achieve this, apply Stirling formula to [\(I.18\)](#) and set $s = eT$, $T = 2 \log 3\epsilon_\alpha^{-1}$, we have

$$3982 \left| e^{-\frac{t}{2}} - \sum_{i=0}^{s-1} \frac{(-1)^i}{i!} \left(\frac{t}{2}\right)^i \right| \leq \left(\frac{1}{2}\right)^{eT} = \left(\frac{\epsilon_\alpha}{3}\right)^{\frac{2e}{\log_2 e}} \leq \frac{\epsilon_\alpha}{3}.$$

3986 Next we set the transformer error $\epsilon = \epsilon_\alpha/3$. Combining [\(I.18\)](#) and [\(I.19\)](#), for $t \in [0, T]$ we obtain

$$3989 \left| \mathcal{T}_t - \exp\left(-\frac{t}{2}\right) \right| \leq \frac{2}{3}\epsilon_\alpha.$$

3992 Furthermore, since $\exp(-T/2) = \epsilon_\alpha/3$, [\(I.20\)](#) becomes

$$3994 \left| \exp\left(-\frac{t}{2}\right) - \mathcal{T}_\alpha(t) \right| \leq \frac{\epsilon_\alpha}{3} + \frac{2\epsilon_\alpha}{3} = \epsilon_\alpha.$$

For the parameter bounds, by the same argument as in the proof of [Lemma I.5](#), we normalize the domain from $[0, T]$ to $[0, 1]$ for the quantization, and then the rest of the step follows [Theorem H.2](#).

This results in parameter bound $\mathcal{O}(\log \epsilon_\alpha^{-1} \epsilon_\alpha^{-\frac{1}{d}})$ for $\|W_1\|_2$ and $\|W_1\|_{2,\infty}$, and the rest of the bounds follow the result in [Lemma H.5](#) with $d = 1$ and $L = 1$.

This completes the proof. \square

Lemma I.11 (Approximation of $\sigma_t = \sqrt{1 - e^{-t}}$ with transformer). For any $\sigma_\sigma \in (0, 1)$, 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

$$|\mathcal{T}_\sigma(t) - \sigma_t| \leq \epsilon_\sigma.$$

The parameter bounds in the transformer network class satisfy

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

Proof. We follow the proof structure of ([Fu et al., 2024b](#), Lemma F.10).

Since $f(t) = \sqrt{1 - e^{-t}}$ with $t \in [t_0, T]$ is a continuous on compact domain. The first part of the proof is complete by applying [Theorem H.2](#).

For the parameter bounds, we take $\mathcal{O}(T\epsilon_\sigma^{-1})$ for $\|W_1\|_2$ and $\|W_1\|_{2,\infty}$ in the first feed-forward layer. This follows from the argument in the proof of [Lemma I.5](#).

The rest of the bounds follow [Lemma H.5](#) with $d = 1$ and $L = 1$

This completes the proof. \square

We have finished the approximation of every key component for the proof of [Theorem 3.1](#). We now proceed to the detailed assembly and integration of these components to finalize the proof.

• Step 3.3: Unified Transformer-Based Score Function Approximation.

First, we establish a theoretical upper bound for transformer model output by analyzing the upper bound of the score function in ℓ_∞ distance under [Assumption 3.1](#) as follows.

– Bound on $p_t(x|y)$:

Recall that the conditional distribution at time t has the form:

$$p_t(x|y) = \frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \int p(x_0|y) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) dx_0.$$

Applying the light tail property in [Assumption 3.1](#), the upper bound follows:

$$p_t(x|y) \leq \frac{C_1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \int \exp\left(-\frac{C_2 \|x_0\|^2}{2}\right) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) dx_0. \quad (\text{I.21})$$

On the other hand, the lower bound follows:

$$p_t(x|y) \geq \frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \int_{\|x_0\| \leq 1} p(x_0|y) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) dx_0. \quad (\text{I.22})$$

– **Bound on $\nabla p_t(x|y)$:** The first element of the gradient has the form:

$$|(\nabla p_t)[1]| = \frac{1}{\sigma_t^2 (2\pi)^{\frac{d}{2}}} \cdot \left| \int \left(\frac{x[1] - \alpha_t x_0[1]}{\sigma_t^2} \right) p(x_0|y) \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right) dx_0 \right|. \quad (\text{I.23})$$

The ℓ_∞ bound on ∇p_t follows by applying light tail property to each coordinate as in (I.21).

Combining (I.21), (I.22) and (I.23), we provide the ℓ_∞ bounds on the score.

Lemma I.12 (Bounds on Score, Lemma A.10 of (Fu et al., 2024b)). Assume [Assumption 3.1](#). There exists a constant K such that

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

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 I.13 (Approximation Score Function with Transformer on Supported Domain). Assume [Assumption 3.1](#). Consider $t \in [N^{-C_\sigma}, C_\alpha \log N]$, for constant C_σ, C_α , 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, β, 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

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

Proof of Lemma I.13. Our poof follows the structure of Fu et al. (2024b, Proposition A.3).

Recall that from [Lemma I.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 ϵ_{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 I.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,

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

From [Lemma I.5](#), we have

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

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

We further write

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

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

$$\|\nabla \log p_t - f_3\|_\infty \lesssim \frac{B}{\sigma_t^2 p_t} N^{-\beta} (\log N)^{\frac{d_x+k_1+1}{2}}. \quad (\text{I.24})$$

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

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 [\(I.24\)](#).

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 I.5](#), [Lemma I.6](#), and [Lemma I.11](#) for f_1 , f_2 , and σ_t respectively. This gives error ϵ_{f_1} , ϵ_{f_2} and $\epsilon_{\sigma,1}$, and we address the clipping of f_1 in later paragraph.

Next, We utilize $\mathcal{T}_{\text{rec},1}$ and $\mathcal{T}_{\text{rec},2}$ in [Lemma I.9](#) for the approximation of the inverse of f_1 and σ_t .

This gives error

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

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 I.8](#) for the approximation of the product of f_1^{-1} , f_2 and σ_t^{-1} .

This gives error of

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

and K_2 is a positive constant. From [Lemma I.8](#) we require that $[-K_2, K_2]$ covers the domain for all of f_1^{-1} , f_2 and f_σ^{-1} .

To be more specific, we reiterate three facts that determines the choice of K_2 .

- Recall that in the **Step A.**, we set $f_{1,\text{clip}} = \{f_1, \epsilon_{\text{low}}\}$.
- [Lemma I.12](#) states $K(C_x \sqrt{d_x \log N} + 1)/\sigma_t^2$ is the ℓ_∞ bound on the score.
- The maximum value of σ_t^{-1} happens at $t = t_0$.

As a result, we set K_2 as

$$K_2 = \max \left(\frac{1}{\epsilon_{\text{low}}}, \frac{K}{\sigma_{t_0}} (C_x \sqrt{d_x \log N} + 1), \frac{1}{\sigma_{t_0}} \right).$$

By the earlier choice of ϵ_{low} , we have $\epsilon_{\text{low}}^{-1} = \mathcal{O}(N^\beta \log N^{-(d_x+k_1)/2})$, and next we expand σ_{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

$$K_2 = \mathcal{O} \left(N^{\beta+C_\sigma} \log^{-\frac{d_x+\beta}{2}} N \right), \quad (\text{I.25})$$

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 I.11](#) for the approximation of σ_t , and this gives error $\epsilon_{\sigma,2}$.

Next, we utilize $\mathcal{T}_{\text{rec},3}$ in [Lemma I.8](#) for the approximation of the inverse of σ_t .

This gives error

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

Next, we utilize $\mathcal{T}_{\text{mult},2}$ for the approximation of the square of σ_t^{-1} .

This gives error of

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

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 $\mathcal{O}(\sigma_t^{-1})$:

$$K_1 = \mathcal{O}(N^{C_\sigma}). \quad (\text{I.26})$$

4212 * **Step B.3: Error Bound on Every Approximation Combined.**

4213 Combining **Step B.1** and **Step B.2**, the total error is bounded by

4214
4215
4216
$$\epsilon_{\text{score}} \leq \max \left(\epsilon_{\text{mult},2} + 2K_1 \left(\epsilon_{\text{rec},3} + \frac{\epsilon_{\sigma,2}}{\epsilon_{\text{rec},3}^2} \right), \epsilon_{\text{mult},1} + 3K_2^2 \epsilon_1 \right).$$

4217
4218
4219 The goal is to guarantee the final error $\epsilon_{\text{score}} \leq N^{-\beta}$ such that it matches the order of the
4220 approximation error in **Step A**. We list all the error choice to achieve the goal.⁹

4221 · **For the Error of the First Two Inverse Operators:**

4222
4223
$$\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2} = \mathcal{O} \left(N^{-(3\beta+2C_\sigma)} (\log N)^{(d_x+\beta)} \right).$$

4224
4225
4226 · **For the Error of the Last Inverse Operator:**

4227
4228
$$\epsilon_{\text{rec},3} = \mathcal{O} \left(N^{-(\beta+2C_\sigma)} \right).$$

4229
4230 · **For the Error of f_1 :**

4231
4232
$$\epsilon_{f_1} = \mathcal{O} \left(N^{-(9\beta+6C_\sigma)} (\log N)^{3(d_x+\beta)} \right).$$

4233
4234 · **For the Error of f_2 :**

4235
4236
$$\epsilon_{f_2} = \mathcal{O} \left(N^{-(3\beta+2C_\sigma)} (\log N)^{(d_x+\beta)} \right).$$

4237
4238 · **For the Error of the First Variance:**

4239
4240
$$\epsilon_{\sigma,1} = \mathcal{O} \left(N^{-(9\beta+6C_\sigma)} (\log N)^{3(d_x+\beta)} \right).$$

4241
4242 · **For the Error of the Second Variance:**

4243
4244
$$\epsilon_{\sigma,2} = \mathcal{O} \left(N^{-(7\beta+5C_\sigma)} (\log N)^{2(d_x+\beta)} \right).$$

4245
4246 · **For the Error of the Two Product Operators:**

4247
4248
$$\epsilon_{\text{mult},1}, \epsilon_{\text{mult},2} = \mathcal{O}(N^{-\beta}).$$

4249
4250 The above error choice renders $\epsilon_{\text{score}} \leq N^{-\beta}$.

4251 Therefore we conclude that there exist a transformer $\mathcal{T}_{\text{score}} \in \mathcal{T}_R^{h,s,r}$ such that

4252
4253
4254
$$\|\mathcal{T}_{\text{score}}(x, y, t) - f_3(x, y, t)\|_\infty \leq N^{-\beta}. \quad (\text{I.27})$$

4255
4256 Combining (I.24) and (I.27) we obtain

4257
4258
$$\|\nabla \log p_t - \mathcal{T}_{\text{score}}(x, y, t)\|_\infty \lesssim \frac{1}{p_t} \frac{B}{\sigma_t^2} N^{-\beta} (\log N)^{\frac{d_x+k_1+1}{2}}.$$

4259
4260
4261
4262
4263
4264
4265 ⁹Further details regarding the choice of each one of ϵ are in Appendix F.4 of (Fu et al., 2024b).

4266 We have completed the first part of the proof. We next give the norm bounds for the transformer
 4267 parameters. Specifically, we select the parameter bounds that are consistent across all operations.
 4268 including [Lemma I.5](#), [Lemma I.6](#), [Lemma I.8](#), [Lemma I.9](#) and [Lemma I.11](#).

4270 – **Step C: Transformer Parameter Bound.**

4271 Our result highlights the influence of N under varying d_x . Therefore, for the transformer
 4272 parameter bounds, we keep terms with d_x, d, L appearing in the exponent of N and $\log N$.

4273 Note that the following parameter selection is based on high-dimensional case where $\log N$ term
 4274 dominates N term.

4275 * **Parameter Bound on W_Q and W_K .**

4276 Given error ϵ , the bound on each operation follows:

4277 · **For ϵ_{f_1} :** By [Lemma I.5](#), we have

$$4280 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)\cdot\frac{2dL+4d+1}{d}} \cdot (\log N)^{-3(d_x+\beta)\cdot\frac{2dL+4d+1}{d}}\right).$$

4283 · **For ϵ_{f_2} :** By [Lemma I.6](#), we have

$$4286 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_\sigma)\cdot\frac{2dL+4d+1}{d}} \cdot (\log N)^{-(d_x+\beta)\cdot\frac{2dL+4d+1}{d}}\right).$$

4288 · **For $\epsilon_{\text{mult},1}$:** By [Lemma I.8](#) with $m = 3$, we have

$$4290 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}(N^{7\beta}).$$

4293 · **For $\epsilon_{\text{mult},2}$:** By [Lemma I.8](#) with $m = 2$, we have

$$4295 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}(N^{5\beta}).$$

4297 · **For $\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2}$:** By [Lemma I.9](#), we have

$$4299 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)}(\log N)^{-3(d_x+\beta)}\right).$$

4302 · **For $\epsilon_{\text{rec},3}$:** By [Lemma I.9](#), we have

$$4304 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+6C_\sigma)}\right).$$

4306 · **For ϵ_{σ_1} :** By [Lemma I.11](#), we have

$$4308 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(27\beta+18C_\sigma)}(\log N)^{-9(d_x+\beta)}\right).$$

4311 · **For ϵ_{σ_2} :** By [Lemma I.11](#), we have

$$4313 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(21\beta+15C_\sigma)}(\log N)^{-6(d_x+\beta)}\right).$$

4317 We select the largest parameter bound from $\epsilon_{\text{mult},1}$ and $\epsilon_{\text{rec},3}$ that remains valid across all other
 4318 approximations. That is, we take $N^{(7\beta+6C_\sigma)}$ as the upper-bound.

4319 * **Parameter Bound on W_O and W_V .**

4320
4321
4322
4323
4324
4325
4326
4327
4328
4329
4330
4331
4332
4333
4334
4335
4336
4337
4338
4339
4340
4341
4342
4343
4344
4345
4346
4347
4348
4349
4350
4351
4352
4353
4354
4355
4356
4357
4358
4359
4360
4361
4362
4363
4364
4365
4366
4367
4368
4369
4370
4371
4372
4373

Given error ϵ , the bound on each operation follows:

• **For ϵ_{f_1} :** By [Lemma I.5](#), we have

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(9\beta+6C_\sigma)}{d}} (\log N)^{\frac{3(d_x+\beta)}{d}}\right).$$

• **For ϵ_{f_2} :** By [Lemma I.6](#), we have

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

• **For $\epsilon_{\text{mult},1}$:** By [Lemma I.8](#) with $m = 3$, we have

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

• **For $\epsilon_{\text{mult},2}$:** By [Lemma I.8](#) with $m = 2$, we have

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-2\beta}\right).$$

• **For $\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2}$:** By [Lemma I.9](#), we have

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

• **For $\epsilon_{\text{rec},3}$:** By [Lemma I.9](#), we have

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

• **For ϵ_{σ_1} :** By [Lemma I.11](#), we have

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(9\beta+6C_\sigma)} (\log N)^{3(d_x+\beta)}\right).$$

• **For ϵ_{σ_2} :** By [Lemma I.11](#), we have

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(7\beta+5C_\sigma)} (\log N)^{2(d_x+\beta)}\right).$$

Note that only ϵ_{f_1} and ϵ_{f_2} involve the reshape operation. From [Lemma H.5](#), we take $\mathcal{O}(\sqrt{d})$ and $\mathcal{O}(d) \|W_V\|_2$ and $\|W_V\|_{2,\infty}$. Moreover, We select the largest parameter bound from $\epsilon_{\text{rec},1}$ and ϵ_{σ_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.

* **Parameter Bound on W_1 .**

Given error ϵ , the bound on each operation follows:

• **For ϵ_{f_1} :** By [Lemma I.5](#), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_\sigma)}{d}} (\log N)^{-\frac{3(d_x+\beta)}{d}} \cdot (\log N)\right).$$

• **For ϵ_{f_2} :** By [Lemma I.6](#), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+2C_\sigma)}{d}} (\log N)^{-\frac{(d_x+\beta)}{d}} \cdot (\log N)\right).$$

4374
4375
4376
4377
4378
4379
4380
4381
4382
4383
4384
4385
4386
4387
4388
4389
4390
4391
4392
4393
4394
4395
4396
4397
4398
4399
4400
4401
4402
4403
4404
4405
4406
4407
4408
4409
4410
4411
4412
4413
4414
4415
4416
4417
4418
4419
4420
4421
4422
4423
4424
4425
4426
4427

- **For $\epsilon_{\text{mult},1}$:** By Lemma I.8 with $m = 3$ and $C = K_2$ in (I.25), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}(K_2 \cdot N^{3\beta}) = \mathcal{O}\left(N^{(4\beta+C_\sigma)}(\log N)^{-\frac{1}{2}(d_x+\beta)}\right).$$

- **For $\epsilon_{\text{mult},2}$:** By Lemma I.8 with $m = 2$ and $C = K_1$ in (I.26), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}(K_1 \cdot N^{2\beta}) = \mathcal{O}\left(N^{(2\beta+C_\sigma)}\right).$$

- **For $\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2}$:** By Lemma I.9, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(6\beta+4C_\sigma)}(\log N)^{-2(d_x+\beta)}\right).$$

- **For $\epsilon_{\text{rec},3}$:** By Lemma I.9, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(2\beta+4C_\sigma)}\right).$$

- **For ϵ_{σ_1} :** By Lemma I.11, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)}(\log N)^{-3(d_x+\beta)} \cdot \log N\right).$$

- **For ϵ_{σ_2} :** By Lemma I.11, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+5C_\sigma)}(\log N)^{-2(d_x+\beta)} \cdot \log N\right).$$

We select the largest parameter bound from $\epsilon_{\text{rec},3}$ that remains valid across all other approximations. That is, we take $N^{(2\beta+4C_\sigma)}$ as the upper-bound.

* **Parameter Bound for W_2 .**

Given error ϵ , the bound on each operation follows:

- **For ϵ_{f_1} :** By Lemma I.5, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_\sigma)}{d}}(\log N)^{-3\frac{(d_x+\beta)}{d}}\right).$$

- **For ϵ_{f_2} :** By Lemma I.6, we have **For ϵ_{f_1} :** By Lemma I.5, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+2C_\sigma)}{d}}(\log N)^{-\frac{(d_x+\beta)}{d}}\right).$$

- **For $\epsilon_{\text{mult},1}$:** By Lemma I.8 with $m = 3$, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}(N^{3\beta}).$$

- **For $\epsilon_{\text{mult},2}$:** By Lemma I.8 with $m = 2$, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}(N^{2\beta}).$$

4428 · **For $\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2}$:** By [Lemma I.9](#), we have

$$4429 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_\sigma)}(\log N)^{-(d_x+\beta)}\right).$$

4432 · **For $\epsilon_{\text{rec},3}$:** By [Lemma I.9](#), we have

$$4433 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(\beta+2C_\sigma)}\right).$$

4437 · **For ϵ_{σ_1} :** By [Lemma I.11](#), we have

$$4438 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)}(\log N)^{-3(d_x+\beta)}\right).$$

4442 · **For ϵ_{σ_2} :** By [Lemma I.11](#), we have

$$4443 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+5C_\sigma)}(\log N)^{-2(d_x+\beta)}\right).$$

4447 We select the largest parameter bound from $\epsilon_{\text{mult},1}$ and $\epsilon_{\text{rec},3}$ that remains valid across all other approximations. That is, we take $N^{(3\beta+2C_\sigma)}$ as the upper-bound.

4448 * **Parameter Bound for E .**

4449 Since only ϵ_{f_1} and ϵ_{f_2} involve the reshape operation. From [Lemma H.5](#), we take $\mathcal{O}(d^{\frac{1}{2}}L^{\frac{3}{2}})$ for $\|E^\top\|_{2,\infty}$.

4450 By integrating results above, we derive the following parameter bounds for the transformer network, ensuring valid approximation across all nine approximations.

$$4451 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+6C_\sigma)}\right);$$

$$4452 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+6C_\sigma)}(\log N)^{3(d_x+\beta)}\right);$$

$$4453 \quad \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);$$

$$4454 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(2\beta+4C_\sigma)}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N/\sigma_t^2}\right);$$

$$4455 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_\sigma)}\right).$$

4456 The last network output bound $C_{\mathcal{T}} = \mathcal{O}(\sqrt{d_x \log N/\sigma_t^2})$ follows the entry-wise minimum bounds $K(C_x \sqrt{d \log N} + 1)/\sigma_t^2$ in ℓ_∞ distance by [Lemma I.12](#).

4457 This completes the proof. □

4474 I.2 MAIN PROOF OF [THEOREM 3.1](#)

4475 In [Lemma I.13](#), we establish the score approximation with transformer that incorporates every 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 $p_t < \epsilon_{\text{low}}$ where the problem of score explosion remains unaddressed.

4480 To combat this, we introduce two additional lemmas. The first lemma gives us the error caused by the truncation of \mathbb{R}^{d_x} within a radius R_1 in ℓ_2 distance.

Lemma I.14 (Truncate x for Score Function, Lemma A.1 of (Fu et al., 2024b)). Assume **Assumption 3.1**. For any $R_1 > 1, y, t > 0$ we have

$$\int_{\|x\|_\infty \geq R_1} p_t(x|y) dx \leq R_1 \exp(-C'_2 R_1^2),$$

$$\int_{\|x\|_\infty \geq R_1} \|\nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \leq \frac{R_1^3}{\sigma_t^4} \exp(-C'_2 R_1^2),$$

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

Remark I.9. Because we only impose assumption on the light tail property of the conditional distribution in **Assumption 3.1**, the unboundedness of x necessitates a truncation for integrals regarding x , or else the result would diverge.

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

Lemma I.15 (Lemma A.2 of (Fu et al., 2024b)). Assume **Assumption 3.1**. For any $R_2, y, \epsilon_{\text{low}} > 0$ we have

$$\int_{\|x\|_\infty \leq R_2} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \cdot p_t(x|y) dx \leq R_2^{d_x} \epsilon_{\text{low}},$$

$$\int_{\|x\|_\infty \leq R_2} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \cdot \|\nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \leq \frac{1}{\sigma_t^4} R_2^{d_x+2} \epsilon_{\text{low}}.$$

Remark I.10. 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 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**.

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 ℓ_2 distance. We approximate the error with **Lemma I.14**

- **(A.2): The approximation for $\mathbb{1}\{p_t(x|y) < \epsilon_{\text{low}}\}$ and $\{\|x\|_\infty \leq R_1\}$.**

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 I.15**.

- **(A.3) The approximation for $\mathbb{1}\{p_t(x|y) \geq \epsilon_{\text{low}}\}$ and $\{\|x\|_\infty \leq R_1\}$.**

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

□

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

To begin with, we set $R_1 = R_2 = \sqrt{2\beta \log N/C'_2}$ in **Lemma I.14** and **Lemma I.15**, and we expand the target into three parts (A_1) , (A_2) , and (A_3) :

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

$$\begin{aligned}
&= \underbrace{\int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx}_{(A_1)} \\
&+ \underbrace{\int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx}_{(A_2)} \\
&+ \underbrace{\int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx}_{(A_3)}.
\end{aligned}$$

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

- **Bounding (A_1) .** We apply [Lemma I.14](#). Note that we have $\|s(x, y, t)\|_\infty \lesssim \sqrt{\log N}/\sigma_t^2$ from the construction of the score estimator in [Lemma I.13](#).

$$\begin{aligned}
&\int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \quad (\text{By expanding the } \ell_2 \text{ norm}) \\
&\leq 2 \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|s(x, y, t)\|_2^2 \cdot p_t(x|y) dx + 2 \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|\nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \\
&\quad (\text{By } \|\cdot\|_2^2 \leq d_x \|\cdot\|_\infty^2) \\
&\leq 2d_x \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|s(x, y, t)\|_\infty^2 \cdot p_t(x|y) dx + 2 \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|\nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \\
&\quad (\text{By the } \ell_\infty \text{ bound on the score function}) \\
&\lesssim 2d_x \left(\frac{\sqrt{\log N}}{\sigma_t^2} \right)^2 \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} p_t(x|y) dx + 2 \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \|\nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \\
&\quad (\text{By } \text{Lemma I.14} \text{ and dropping constant}) \\
&\lesssim 2d_x \left(\frac{\sqrt{\log N}}{\sigma_t^2} \right)^2 \left(\sqrt{\frac{2\beta}{C'_2} \log N} N^{-2\beta} \right) + \frac{2}{\sigma_t^4} \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{3}{2}} N^{-2\beta} \\
&\quad (\text{By dropping constant and lower order term}) \\
&\lesssim \frac{1}{\sigma_t^4} N^{-2\beta} (\log N)^{\frac{3}{2}}.
\end{aligned}$$

- **Bounding (A_2) .** We apply [Lemma I.15](#). Note that we set $\epsilon_{\text{low}} = C_3 N^{-\beta} (\log N)^{(d_x + k_1)/2}$ in [Lemma I.13](#).

$$\begin{aligned}
&\int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \\
&\quad (\text{By expanding the } \ell_2 \text{ norm}) \\
&\leq \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} 2 \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \left(\|s(x, y, t)\|_2^2 + \|\nabla \log p_t(x|y)\|_2^2 \right) \cdot p_t(x|y) dx \\
&\quad (\text{By } \|\cdot\|_2^2 \leq d_x \|\cdot\|_\infty^2) \\
&\leq \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \left(d_x \|s(x, y, t)\|_\infty^2 + \|\nabla \log p_t(x|y)\|_2^2 \right) \cdot p_t(x|y) dx \\
&\quad (\text{By the } \ell_\infty \text{ bound on the score function})
\end{aligned}$$

$$\begin{aligned}
& \lesssim \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| < \epsilon_{\text{low}}\} \left(d_x \left(\frac{\sqrt{\log N}}{\sigma_t^2} \right)^2 + \|\nabla \log p_t(x|y)\|_2^2 \right) \cdot p_t(x|y) dx \\
& \hspace{15em} \text{(By Lemma I.15 and dropping constant)} \\
& \lesssim d_x \left(\frac{\sqrt{\log N}}{\sigma_t^2} \right)^2 \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{d_x}{2}} \epsilon_{\text{low}} + \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{d_x+2}{2}} \frac{\epsilon_{\text{low}}}{\sigma_t^4} \\
& \hspace{15em} \text{(By dropping constant and lower order term)} \\
& \lesssim \frac{1}{\sigma_t^4} (\log N)^{\frac{d_x+2}{2}} \epsilon_{\text{low}}.
\end{aligned}$$

• **Bounding (A_3) .** We apply Lemma I.13.

$$\begin{aligned}
& \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx \\
& \hspace{15em} \text{(By } \|\cdot\|_2^2 \leq d_x \|\cdot\|_\infty^2 \text{)} \\
& \leq \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\} d_x \|s(x, y, t) - \nabla \log p_t(x|y)\|_\infty^2 \cdot p_t(x|y) dx \\
& \hspace{15em} \text{(Multiply with } p_t/p_t \text{)} \\
& = \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \frac{\mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\}}{p_t(x|y)} d_x \|s(x, y, t) - \nabla \log p_t(x|y)\|_\infty^2 \cdot p_t^2(x|y) dx \\
& \hspace{15em} \text{(By Lemma I.13)} \\
& \lesssim \frac{B^2 d_x}{\sigma_t^2} N^{-2\beta} (\log N)^{d_x+k_1+1} \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\} p_t(x|y) dx \\
& \hspace{15em} \text{(Multiply with } \epsilon_{\text{low}}/\epsilon_{\text{low}} \text{)} \\
& = \frac{B^2 d_x}{\sigma_t^2 \epsilon_{\text{low}}} N^{-2\beta} (\log N)^{d_x+k_1+1} \int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \mathbb{1}\{|p_t(x|y)| \geq \epsilon_{\text{low}}\} \frac{\epsilon_{\text{low}}}{p_t(x|y)} dx \\
& \hspace{15em} \text{(By Lemma I.15)} \\
& \lesssim \frac{B^2 d_x}{\sigma_t^2 \epsilon_{\text{low}}} N^{-2\beta} (\log N)^{d_x+k_1+1} \cdot \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{d_x}{2}} \\
& \hspace{15em} \text{(By the choice of } \epsilon_{\text{low}} \text{ and dropping lower order term)} \\
& \lesssim \frac{B^2 d_x}{\sigma_t^4 \epsilon_{\text{low}}} N^{-2\beta} (\log N)^{\frac{3d_x}{2}+k_1+1}.
\end{aligned}$$

• **Combining the Results.**

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

$$\begin{aligned}
& \int_{\mathbb{R}^d} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \\
& \lesssim \underbrace{\frac{N^{-2\beta} (\log N)^{\frac{3}{2}}}{\sigma_t^4}}_{(A_1)} + \underbrace{\frac{\epsilon_{\text{low}} (\log N)^{\frac{d_x+2}{2}}}{\sigma_t^4}}_{(A_2)} + \underbrace{\frac{B^2 d}{\sigma_t^4 \epsilon_{\text{low}}} N^{-2\beta} (\log N)^{\frac{3d_x}{2}+k_1+1}}_{(A_3)}.
\end{aligned}$$

By replacing ϵ_{low} with $C_3 N^{-\beta} (\log N)^{d_x + k_1/2}$ and using the relation $k_1 \leq \beta$,¹⁰ we obtain

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

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

The transformer parameter norm bounds follow [Lemma I.13](#), with the replacement of N with $N^{1/(d_x + d_y)}$ 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 interpretation of the cutoff and early stopping time parameter, we reset C_α as $(d_x + d_y)C_\alpha$ and C_σ as $(d_x + d_y)C_\sigma$ such that $t \in [N^{-C_\alpha}, C_\sigma \log N]$.

This completes the proof. □

¹⁰Recall the definition of the Hölder smoothness from [Definition 3.1](#).

J PROOF OF THEOREM 3.2

We provide the formal version of Theorem 3.2 at the end of Appendix J.2.

- **Step 0.** We decompose the density function and the score function under Assumption 3.2. In Lemma J.1, we provide details regarding the decomposed form of the score function presented in (3.2). We specify the upper and lower bound on h and ∇h in Lemma J.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 J.3.
- **Step 2.** We construct polynomial approximation based on Taylor expansion of h and ∇h in Lemmas J.4 and J.5. The approximation result captures the local Hölder smoothness, with improved precision relative to the analogous step in Lemma I.3 and Lemma I.4.
- **Step 3.** We approximate h and ∇h with transformer in Lemmas J.6 and J.7. In order to construct the score approximator with transformer, we approximate several additional algebraic operators with transformer in Lemma J.8, Lemma J.9 and Lemma J.10. We incorporate these results into a unified transformer architecture in Lemma J.11.

Organization. Appendix J.1 includes the four steps and auxiliary lemmas supporting our proof. Appendix J.2 includes the formal version and main proof of Theorem 3.2.

J.1 AUXILIARY LEMMAS

Step 0: Decompose the Score with Stronger Hölder Smoothness Assumption. We utilize the condition assumed in Assumption 3.2 to achieve the decomposition.

Lemma J.1 (Lemma B.1 of Fu et al. (2024b)). Assume Assumption 3.2. The conditional distribution at time t has the following expression:

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

Moreover, the score function has the following expression:

$$\nabla \log p_t(x|y) = \frac{-C_2x}{\alpha_t^2 + C_2\sigma_t^2} + \frac{\nabla h(x, y, t)}{h(x, y, t)},$$

where $h(x, y, t) = \int \frac{f(x_0, y)}{\hat{\sigma}_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|x_0 - \hat{\alpha}_t x\|^2}{2\hat{\sigma}_t^2}\right) dx_0$, $\hat{\sigma}_t = \frac{\sigma_t}{(\alpha_t^2 + C_2\sigma_t^2)^{1/2}}$, and $\hat{\alpha}_t = \frac{\alpha_t}{\alpha_t^2 + C_2\sigma_t^2}$.

Proof. From Assumption 3.2, we have the initial conditional density with the form: $p(z|y) = \exp\left(-C_2\|z\|_2^2/2\right) \cdot f(z, y)$.

This allows the decomposition:

$$p_t(x|y) = \int \frac{1}{\sigma_t^d (2\pi)^{d/2}} p(z|y) \exp\left(-\frac{\|x - \alpha_t z\|^2}{2\sigma_t^2}\right) dz, \quad (\text{J.1})$$

$$= \frac{1}{\sigma_t^d (2\pi)^{d/2}} \int \exp\left(-\frac{C_2\|z\|_2^2}{2}\right) f(z, y) \exp\left(-\frac{\|x - \alpha_t z\|^2}{2\sigma_t^2}\right) dz. \quad (\text{J.2})$$

We rearrange the two exponential terms in (J.2) into

$$\exp\left(-\frac{C_2\|z\|_2^2}{2}\right) \exp\left(-\frac{\|x - \alpha_t z\|^2}{2\sigma_t^2}\right) = \exp\left(-\frac{1}{2\sigma_t^2} \sum_{i=1}^d (x[i]^2 - 2\alpha_t x[i]z[i] + \alpha_t^2 z[i]^2 + C_2\sigma_t^2 z[i]^2)\right).$$

Note that, we replace the summation in the exponents by first focusing on one coordinate and then do the product for all d components.

Without loss of generality, we derive the first coordinate of the function:

$$\begin{aligned}
& \exp\left(-\frac{1}{2\sigma_t^2}(x[1]^2 - 2\alpha_t x[1]z[1] + \alpha_t^2 z[1]^2 + C_2\sigma_t^2 z[1]^2)\right), \\
&= \exp\left(-\frac{1}{2\sigma_t^2}(\alpha_t^2 + C_2\sigma_t^2)\left(z[1]^2 - \frac{2\alpha_t}{\alpha_t^2 + C_2\sigma_t^2}x[1]z[1] + \frac{x[1]^2}{\alpha_t^2 + C_2\sigma_t^2}\right)\right), \\
&= \exp\left(-\frac{1}{2\sigma_t^2}(\alpha_t^2 + C_2\sigma_t^2)\left(z[1] - \frac{\alpha_t x[1]}{\alpha_t^2 + C_2\sigma_t^2}\right)^2 - \frac{1}{2\sigma_t^2}\left(\frac{-\alpha_t^2}{\alpha_t^2 + C_2\sigma_t^2} + 1\right)x[1]^2\right), \\
&= \exp\left(-\frac{1}{2\sigma_t^2}(\alpha_t^2 + C_2\sigma_t^2)\left(z[1] - \frac{\alpha_t x[1]}{\alpha_t^2 + C_2\sigma_t^2}\right)^2\right) \exp\left(-\frac{C_2 x[1]^2}{2(\alpha_t^2 + C_2\sigma_t^2)}\right).
\end{aligned}$$

The other $d_x - 1$ coordinates abide by the same derivation. Consider the product of them, we have:

$$\exp\left(-\frac{C_2\|z\|_2^2}{2}\right) \exp\left(-\frac{\|x - \alpha_t z\|_2^2}{2\sigma_t^2}\right), \quad (\text{J.3})$$

$$= \exp\left(-\frac{1}{2\sigma_t^2}(\alpha_t^2 + C_2\sigma_t^2)\left\|z - \frac{\alpha_t x}{\alpha_t^2 + C_2\sigma_t^2}\right\|^2\right) \exp\left(-\frac{C_2}{2(\alpha_t^2 + C_2\sigma_t^2)}\|x\|_2^2\right). \quad (\text{J.4})$$

Following (Fu et al., 2024b), we plug (J.3) into (J.1) and set $\hat{\alpha}_t = \frac{\alpha_t}{\alpha_t^2 + C_2\sigma_t^2}$ and $\hat{\sigma}_t^2 = \frac{\sigma_t^2}{\alpha_t^2 + C_2\sigma_t^2}$ for simplicity. Then we get:

$$\begin{aligned}
& p_t(x|y) \\
&= \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left(-\frac{C_2\|x\|_2^2}{2(\alpha_t^2 + C_2\sigma_t^2)}\right) \int f(z, y) \exp\left(-\frac{1}{2\sigma_t^2}(\alpha_t^2 + C_2\sigma_t^2)\left\|z - \frac{\alpha_t x}{\alpha_t^2 + C_2\sigma_t^2}\right\|^2\right) dz, \\
&= \frac{1}{\sigma_t^d (2\pi)^{d/2}} \exp\left(-\frac{C_2\|x\|_2^2}{2(\alpha_t^2 + C_2\sigma_t^2)}\right) \int f(z, y) \exp\left(-\frac{\|z - \hat{\alpha}_t x\|_2^2}{2\hat{\sigma}_t^2}\right) dz.
\end{aligned}$$

Finally, we define $h(x, y, t) = \int \frac{1}{\hat{\sigma}_t^d (2\pi)^{d/2}} f(z, y) \exp\left(-\frac{\|z - \hat{\alpha}_t x\|_2^2}{2\hat{\sigma}_t^2}\right) dz$ and plug it back to the equation above.

The form of the score function is proved by simply implementing the logarithm and the gradient to the result of $p_t(x|y)$

This completes the proof. \square

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

Lemma J.2 (Lemma B.8 of (Fu et al., 2024b)). Under Assumption 3.2, we have the following bounds for $h(x, y, t)$ and $\frac{\partial h}{\partial \alpha_t} \nabla h(x, y, t)$

$$C_1 \leq h(x, y, t) \leq B, \quad \left\| \frac{\partial h}{\partial \alpha_t} \nabla h(x, y, t) \right\|_\infty \leq \sqrt{\frac{2}{\pi}} B,$$

where C_1 and B are the hyperparameters of $\mathcal{H}^\beta(\mathbb{R}^{d_x} \times [0, 1]^{d_y}, B)$ in Assumption 3.2.

Remark J.1 (Bound on h and ∇h). We reiterate that Lemma J.2 drives the key distinction between the analyses in Theorem 3.1 and Theorem 3.2. Specifically, in Appendix I.2, the decomposed term containing the threshold ϵ_{low} results in lower approximation rate, while bounds on h and ∇h eliminate the need of the threshold with h 's lower bound C_1 , rendering faster approximation rate.

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

Lemma J.3 (Clipping Integral, Lemma B.10 of [Fu et al. \(2024b\)](#)). Assume [Assumption 3.2](#). 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

$$\int_{\mathbb{R}^{d_x} \setminus B_x} \left| \left(\frac{\hat{\alpha}_t x_0 - x}{\hat{\sigma}_t} \right)^\kappa \right| \cdot p(x_0|y) \cdot \frac{1}{\hat{\sigma}_t^d (2\pi)^{d/2}} \exp\left(-\frac{\|\hat{\alpha}_t x_0 - x\|^2}{2\hat{\sigma}_t^2}\right) dx_0 \leq \epsilon, \quad (\text{J.5})$$

where $\left(\frac{\hat{\alpha}_t x_0 - x}{\hat{\sigma}_t} \right)^\kappa := \left(\left(\frac{\hat{\alpha}_t x_0[1] - x[1]}{\hat{\sigma}_t} \right)^{\kappa[1]}, \left(\frac{\hat{\alpha}_t x_0[2] - x[2]}{\hat{\sigma}_t} \right)^{\kappa[2]}, \dots, \left(\frac{\hat{\alpha}_t x_0[d_x] - x[d_x]}{\hat{\sigma}_t} \right)^{\kappa[d_x]} \right)$ and

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

Step 2: Approximate h and ∇h with Polynomials. Similar to the construction of the diffused local polynomials in [Lemma I.5](#) and [Lemma I.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 J.4 (Approximation with Diffused Local Polynomials, Lemma B.4 of ([Fu et al., 2024b](#))). Assume [Assumption 3.2](#). For sufficiently larger $N > 0$ and constant C_2 , 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

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

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 J.5 (Counterpart of [Lemma J.4](#), Lemma B.6 of ([Fu et al., 2024b](#))). Assume [Assumption 3.2](#). 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

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

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

Step 3: Approximate Diffused Local Polynomials and Algebraic Operators with Transformers.

First, we apply the universal approximation theory of transformers to f_1 and f_2 . Second, we adopt a comparable approach to approximate the algebraic operators essential for the final score computation. Last, we introduce [Lemma J.11](#) that outlines how these components fit into a single transformer 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 H.2](#). We utilize network consisting of one transformer block and one feed-forward layer (see [Figure 1](#) and [Definition 2.2](#)).

Lemma J.6 (Approximate Scalar Polynomials with Transformers). Assume [Assumption 3.1](#). Consider the diffused local polynomial f_1 in [Lemma J.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

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

The parameter bounds in the transformer network class follows [Lemma I.5](#).

Proof of Lemma J.6. The proof closely follows [Lemma I.5](#) □

Lemma J.7 (Approximate Vector-Valued Polynomials with Transformers). Assume [Assumption 3.1](#) and consider $f_2(x, y, t) \in \mathbb{R}^{d_x}$ in [Lemma J.5](#). For any $\epsilon > 0$, there exists a transformer $\mathcal{T}_{f_2} \in \mathcal{T}_R^{h,s,r}$ such that

$$\|f_2(x, y, t) - \mathcal{T}_{f_2}\|_\infty \leq \epsilon,$$

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 parameter bounds in the transformer network class follows [Lemma I.5](#).

Proof of Lemma J.7. The proof closely follows [Lemma I.6](#) \square

• **Step 3.2: Approximate Algebraic Operators with Transformers.**

Next, we introduce lemmas regarding the function of time. These are also key components to the proof of [Theorem J.1](#).

Lemma J.8 (Approximation of α^2 with Transformer). For $t \in [t_0, T]$ with $t_0 < 1$, there exists Transformer $\mathcal{T}_{\alpha^2}(t) \in \mathcal{T}_R^{h,s,r}$ such that

$$|\mathcal{T}_{\alpha^2} - \alpha^2| \leq \epsilon_{\hat{\alpha}}.$$

The parameter bounds in the Transformer network class follow [Lemma I.11](#).

Proof. The proof closely follows [Lemma I.11](#). \square

Also, we approximate $\hat{\alpha}$ and $\hat{\sigma}_t$ as well.

Lemma J.9 (Approximation of $\hat{\alpha}$ with Transformer). Consider $\hat{\alpha}_t = \frac{\alpha_t}{\alpha_t^2 + C_2\sigma_t^2}$, for $t \in [t_0, T]$ with $t_0 < 1$, there exists Transformer $\mathcal{T}_{\hat{\alpha}}(t) \in \mathcal{T}_R^{h,s,r}$ such that

$$|\mathcal{T}_{\hat{\alpha}} - \hat{\alpha}| \leq \epsilon_{\hat{\alpha}}.$$

The parameter bounds in the transformer network class follow [Lemma I.11](#).

Proof. The proof closely follows [Lemma I.11](#). \square

Lemma J.10 (Approximation of $\hat{\sigma}$ with Transformer). Consider $\hat{\sigma}_t = \frac{\sigma_t}{(\alpha_t^2 + C_2\sigma_t^2)^{1/2}}$, for $t \in [t_0, T]$ with $t_0 < 1$, there exists Transformer $\mathcal{T}_{\hat{\sigma}}(t) \in \mathcal{T}_R^{h,s,r}$ such that

$$|\mathcal{T}_{\hat{\sigma}} - \hat{\sigma}| \leq \epsilon_{\hat{\sigma}}.$$

The parameter bounds in the transformer network class follow [Lemma I.11](#).

Proof. The proof closely follows [Lemma I.11](#). \square

We have finished establishing the approximation with transformer for every key component for the proof of [Theorem 3.2](#).

• **Step 3.3: Unified Transformer-Based Score Function Approximation.**

We introduce the counterpart of [Lemma I.13](#). It is the core of the proof for [Theorem J.1](#).

Lemma J.11 (Score Approximation with Transformer). Assume [Assumption 3.2](#). For sufficiently large integer N , there exists a mapping from transformer $\mathcal{T}_{\text{score}} \in \mathcal{T}_R^{h,s,r}$ such that

$$\left\| \mathcal{T}_{\text{score}} - \nabla \log h(x, y, t) + \frac{C_2 x}{\alpha_t^2 + C_2 \sigma_t^2} \right\|_\infty \leq \frac{B}{\sigma_t} N^{-\beta} (\log N)^{\frac{k_1+1}{2}},$$

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 parameter bounds in the transformer network class satisfy

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

Proof. Our proof follows the proof structure of (Fu et al., 2024b, Proposition B.3).

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

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

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 J.4** and **Lemma J.5** to approximate $h(x, y, t)$ and $\nabla h(x, y, t)$ respectively.

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

Next, by **Lemma J.4** and **Lemma J.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 ∇h , denoted as $\nabla h[1]$:

$$\begin{aligned} \left| \frac{\nabla h[1]}{h} - \frac{\hat{\alpha}_t f_2[1]}{\hat{\sigma}_t f_1} \right| & \leq \left| \frac{\nabla h[1]}{h} - \frac{\nabla h[1]}{f_1} \right| + \left| \frac{\nabla h[1]}{f_1} - \frac{\hat{\alpha}_t f_2[1]}{\hat{\sigma}_t f_1} \right|, \\ & \leq \left| \frac{\nabla h[1]}{h \cdot f_1} \right| |f_1 - h| + \frac{\hat{\alpha}_t}{\hat{\sigma}_t} \left| \frac{1}{f_1} \right| \left| f_2 - \frac{\hat{\sigma}_t}{\hat{\alpha}_t} \nabla h[1] \right|, \\ & \lesssim \frac{\hat{\alpha}_t}{\hat{\sigma}_t} \left(|f_1 - h| + \left| f_2 - \frac{\hat{\sigma}_t}{\hat{\alpha}_t} \nabla h[1] \right| \right), \quad (\text{By bounds on } h, \nabla h, f_1, f_2) \\ & \lesssim \frac{\hat{\alpha}_t}{\hat{\sigma}_t} \left(BN^{-\beta} (\log N)^{\frac{k_1}{2}} + BN^{-\beta} (\log N)^{\frac{k_1+1}{2}} \right), \\ & \hspace{15em} (\text{By Lemma J.4 and Lemma J.5}) \\ & \lesssim \frac{1}{\sigma_t} \left(BN^{-\beta} (\log N)^{\frac{k_1+1}{2}} \right). \end{aligned}$$

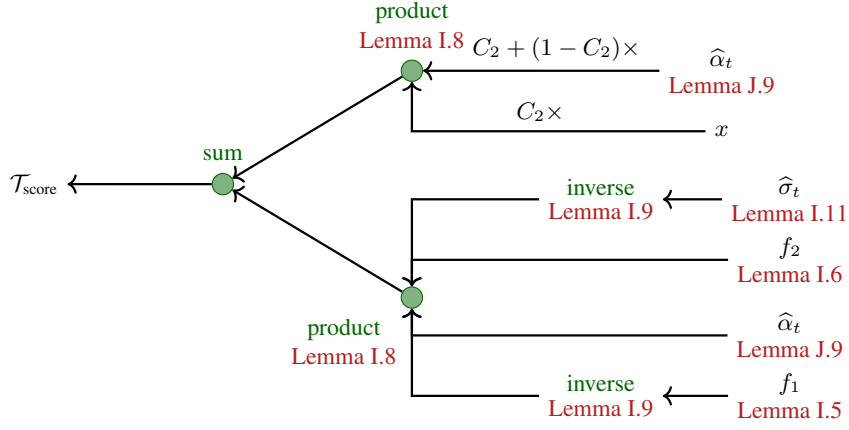


Figure 5: **Approximate Score Function under Assumption 3.2 with Transformer $\mathcal{T}_{\text{score}}$.** The construction of the final score function consists of the approximation of diffused local polynomials f_1 and f_2 with transformer and transformer-approximate operators. We highlight the overall pipeline and related lemmas to ensemble the Transformer network.

Note that in the last line, we utilize

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

By the symmetry of each coordinate in ∇h , we obtain the ℓ_∞ bounds:

$$\left\| \frac{\nabla h(x, y, t)}{h(x, y, t)} - \frac{\hat{\alpha}_t f_2(x, y, t)}{\hat{\sigma}_t f_1(x, y, t)} \right\|_\infty \lesssim \frac{B}{\sigma_t} N^{-\beta} (\log N)^{\frac{k_1+1}{2}}. \quad (\text{J.6})$$

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

Next, we prove that there exist Transformer networks $\mathcal{T}_{\text{score}} \in \mathcal{T}_R^{h,s,r}$ that approximates $f_3(x, y, t)$ with error of order $N^{-\beta}$. We illustrate the overall approximation of f_3 in Figure 5.

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 I.5, Lemma I.6, Lemma J.9 and Lemma J.10 to approximate each one of the component. This gives error ϵ_{f_1} , ϵ_{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 I.9 for the approximation of the inverse of f_1 and $\hat{\sigma}_t$. This gives error

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

and

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

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

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

and K_3 is a positive constant.

From [Lemma I.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 $\hat{\alpha}_{t_0}^{-1}$:

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

and

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

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}_{α^2} in [Lemma J.8](#) for the approximation of α_t^2 . This gives error ϵ_{α^2} .

Next, we utilize $\mathcal{T}_{\text{rec},1}$ in [Lemma I.8](#) for the approximation of the inverse of α_t^2 .

This gives error

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

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

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

and from [Lemma I.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.**

5076 Combining **Step B.1** and **Step B.2**, we obtain the total network with error bounded by
5077

$$5078 \epsilon_{\text{score}} \leq \epsilon_{\text{mult},2} + 2K_3 \left(\epsilon_{\text{rec},1} + \frac{\epsilon_{\alpha^2}}{2\epsilon_{\text{rec},1}} \right) + \epsilon_{\text{mult},1} + 4K_4^3 \epsilon_2.$$

5081
5082 Next, we specify on the choice of ϵ in each approximation to attain a final approximation
5083 error of order $N^{-\beta}$.

5084
5085 · **For the Error of the First Inverse Operator:**

$$5086 \epsilon_{\text{rec},1} = \mathcal{O} \left(N^{-(\beta + \frac{1}{2}C_\alpha)} \right).$$

5087
5088
5089 · **For the Error of the Second and Third Inverse Operator:**

$$5090 \epsilon_{\text{rec},2}, \epsilon_{\text{rec},3} = \mathcal{O} \left(N^{-(\beta + 3C_\sigma)} \right).$$

5091
5092
5093
5094 · **For the Error of f_1 :**

$$5095 \epsilon_{f_1} = \mathcal{O} \left(N^{-(3\beta + 9C_\sigma)} \right).$$

5096
5097
5098 · **For the Error of f_2 :**

$$5099 \epsilon_{f_2} = \mathcal{O} \left(N^{-(\beta + 3C_\sigma)} \right).$$

5100
5101
5102
5103 · **For the Error of $\hat{\sigma}$:**

$$5104 \epsilon_{\hat{\sigma}} = \mathcal{O} \left(N^{-(3\beta + 9C_\sigma)} \right).$$

5105
5106
5107
5108 · **For the Error of $\hat{\alpha}$:**

$$5109 \epsilon_{\hat{\alpha}} = \mathcal{O} \left(N^{-(\beta + 3C_\sigma)} \right).$$

5110
5111
5112 · **For the Error of α^2 :**

$$5113 \epsilon_{\alpha^2} = \mathcal{O} \left(N^{-(3\beta + \frac{3}{2}C_\alpha)} \right).$$

5114
5115
5116
5117 · **For the Error of the Two Product Operators:**

$$5118 \epsilon_{\text{mult},1}, \epsilon_{\text{mult},2} = \mathcal{O}(N^{-\beta}).$$

5119
5120
5121
5122 With above error choice, we have

$$5123 |\mathcal{T}_{\text{score}}(x, y, t) - f_3(x, y, t)| \leq N^{-\beta}. \quad (\text{J.7})$$

5124
5125 Combining (J.6), (J.7) and dropping lower order term, we obtain

$$5126 \|\mathcal{T}_{\text{score}} - \nabla \log p_t(x|y)\|_\infty \lesssim \frac{B}{\sigma_t} N^{-\beta} (\log N)^{\frac{k_1+1}{2}}.$$

5127
5128
5129

We have completed the first part of the proof. Next, we select the parameter bounds based on all the above approximations.

Step C: Transformer Parameter Bound.

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 $\log N$.

– Parameter Bound on W_Q and W_K .

Given error ϵ , the bound on each operation follows:

* For ϵ_{f_1} : By Lemma I.5, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_\sigma)\frac{2dL+4d+1}{d}}\right).$$

* For ϵ_{f_2} : By Lemma I.6, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(\beta+3C_\sigma)\frac{2dL+4d+1}{d}}\right).$$

* For $\epsilon_{\text{mult},1}$: By Lemma I.8 with $m = 4$, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{9\beta}\right).$$

* For $\epsilon_{\text{mult},2}$: By Lemma I.8 with $m = 2$, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{5\beta}\right).$$

* For $\epsilon_{\text{rec},1}$: By Lemma I.9, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{3\beta+\frac{3C_\alpha}{2}}\right).$$

* For $\epsilon_{\text{rec},2}$ and $\epsilon_{\text{rec},3}$: By Lemma I.9, we have

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

* For $\epsilon_{\hat{\alpha}}$: By Lemma I.11, we have

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

* For ϵ_{α^2} : By Lemma I.11, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{9\beta+\frac{9C_\alpha}{2}}\right).$$

* For $\epsilon_{\hat{\sigma}}$: By Lemma I.11, we have

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{9\beta+27C_\sigma}\right).$$

We select the largest parameter bound from ϵ_{f_1} that remains valid across all other approximations.

– Parameter Bound on W_O and W_V .

5184 Given error ϵ , the bound on each operation follows:

5185 * For ϵ_{f_1} : By Lemma I.5, we have

$$5186 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(3\beta+9C_\sigma)}{d}}\right).$$

5187

5188

5189

5190 * For ϵ_{f_2} : By Lemma I.6, we have

$$5191 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(\beta+3C_\sigma)}{d}}\right).$$

5192

5193

5194

5195 * For $\epsilon_{\text{mult},1}$: By Lemma I.8 with $m = 4$, we have

$$5196 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-4\beta}\right).$$

5197

5198

5199 * For $\epsilon_{\text{mult},2}$: By Lemma I.8 with $m = 2$, we have

$$5200 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-2\beta}\right).$$

5201

5202

5203

5204 * For $\epsilon_{\text{rec},1}$: By Lemma I.9, we have

$$5205 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(\beta+\frac{C_\alpha}{2})}\right).$$

5206

5207

5208

5209 * For $\epsilon_{\text{rec},2}$ and $\epsilon_{\text{rec},3}$: By Lemma I.9, we have

$$5210 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(\beta+3C_\sigma)}\right).$$

5211

5212

5213 * For $\epsilon_{\hat{\alpha}}$: By Lemma I.11, we have

$$5214 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(\beta+3C_\sigma)}\right).$$

5215

5216

5217

5218 * For ϵ_{α^2} : By Lemma I.11, we have

$$5219 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+\frac{3C_\alpha}{2})}\right).$$

5220

5221

5222

5223 * For $\epsilon_{\hat{\sigma}}$: By Lemma I.11, we have

$$5224 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+9C_\sigma)}\right).$$

5225

5226

5227

5228 Since we do not impose any relation on C_σ , C_α and β , we simply take looser bound

5229 $\|W_O\|_2, \|W_O\|_{2,\infty} = N^{-\beta}$. Moreover, since only ϵ_{f_1} and ϵ_{f_2} involve the reshape operation.

5230 From Lemma H.5, we take $\mathcal{O}(\sqrt{d})$ and $\mathcal{O}(d) \|W_V\|_2$ and $\|W_V\|_{2,\infty}$.

5231

5232 – **Parameter Bound for W_1 .**

5233 Given error ϵ , the bound on each operation follows:

5234

5235 * For ϵ_{f_1} : By Lemma I.5, we have

$$5236 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+9C_\sigma)}{d}} \cdot \log N\right).$$

5237

5238
5239
5240
5241
5242
5243
5244
5245
5246
5247
5248
5249
5250
5251
5252
5253
5254
5255
5256
5257
5258
5259
5260
5261
5262
5263
5264
5265
5266
5267
5268
5269
5270
5271
5272
5273
5274
5275
5276
5277
5278
5279
5280
5281
5282
5283
5284
5285
5286
5287
5288
5289
5290
5291

* **For ϵ_{f_2} :** By Lemma I.6, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(\beta+3C_\sigma)}{d}} \cdot \log N\right).$$

* **For $\epsilon_{\text{mult},1}$:** By Lemma I.8 with $m = 4$ and $C = K_4$ in (I.25), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}(K_4 \cdot N^{4\beta}) = \mathcal{O}\left(N^{(4\beta+C_\sigma)}\right).$$

* **For $\epsilon_{\text{mult},2}$:** By Lemma I.8 with $m = 2$ and $C = K_3$ in (I.26), we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}(K_3 \cdot N^{2\beta}) = \mathcal{O}\left(N^{(2\beta+\frac{C_\sigma}{2})}\right).$$

* **For $\epsilon_{\text{rec},1}$:** By Lemma I.9, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{2\beta+C_\alpha}\right).$$

* **For $\epsilon_{\text{rec},2}$ and $\epsilon_{\text{rec},3}$:** By Lemma I.9, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(2\beta+6C_\sigma)}\right).$$

* **For $\epsilon_{\hat{\alpha}}$:** By Lemma I.11, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(\beta+3C_\sigma)} \cdot \log N\right).$$

* **For ϵ_{α^2} :** By Lemma I.11, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+\frac{3C_\sigma}{2})} \cdot \log N\right).$$

* **For $\epsilon_{\hat{\sigma}}$:** By Lemma I.11, we have

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+9C_\sigma)} \cdot \log N\right).$$

We select the largest parameter bound from ϵ_{f_1} that remains valid across all other approximations.

– Parameter Bound for W_2 .

Given error ϵ , the bound on each operation follows:

* **For ϵ_{f_1} :** By Lemma I.5, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+9C_\sigma)}{d}}\right).$$

* **For ϵ_{f_2} :** By Lemma I.6, we have

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(\beta+3C_\sigma)}{d}}\right).$$

5292 * **For $\epsilon_{\text{mult},1}$:** By Lemma I.8 with $m = 4$, we have

$$5293 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}(N^{4\beta}).$$

5296 * **For $\epsilon_{\text{mult},2}$:** By Lemma I.8 with $m = 2$, we have

$$5297 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}(N^{2\beta}).$$

5301 * **For $\epsilon_{\text{rec},1}$:** By Lemma I.9, we have

$$5302 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(\beta + \frac{C_\alpha}{2})}\right).$$

5305 * **For $\epsilon_{\text{rec},2}$ and $\epsilon_{\text{rec},3}$:** By Lemma I.9, we have

$$5306 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(\beta + 3C_\sigma)}\right).$$

5310 * **For $\epsilon_{\hat{\alpha}}$:** By Lemma I.11, we have

$$5311 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(\beta + 3C_\sigma)}\right).$$

5315 * **For ϵ_{α^2} :** By Lemma I.11, we have

$$5316 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta + \frac{3C_\alpha}{2})}\right).$$

5319 * **For $\epsilon_{\hat{\sigma}}$:** By Lemma I.11, we have

$$5320 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta + 9C_\sigma)}\right).$$

5324 We select the largest parameter bound from ϵ_{f_1} that remains valid across all other approximations.

5326 – Parameter Bound for E .

5327 Since only ϵ_{f_1} and ϵ_{f_2} involve the reshape operation. From Lemma H.5, we take $\mathcal{O}(d^{1/2}L^{3/2})$.

5330 By integrating results above, we derive the following parameter bounds for the transformer network, ensuring valid approximation across all ten approximations.

$$5331 \quad \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta + 9C_\sigma)\frac{2dL+4d+1}{d}}\right);$$

$$5332 \quad \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}(N^{-\beta});$$

$$5333 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{4\beta + 9C_\sigma + \frac{3C_\alpha}{2}} \cdot \log N\right); \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}}L^{\frac{3}{2}}\right);$$

$$5334 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{4\beta + 9C_\sigma + \frac{3C_\alpha}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t\right).$$

5341 This completes the proof. □

5342
5343
5344
5345

5346 J.2 MAIN PROOF OF **THEOREM 3.2**

5347 We state the formal version of **Theorem 3.2**.

5349 Next, similar to the proof of **Theorem 3.1**, we need the truncation of x due to the unboundedness as
5350 well.

5351 **Lemma J.12** (Truncate x , Lemma B.2 of (Fu et al., 2024b)). Assume **Assumption 3.2**. For any
5352 $R_3 > 1$, we have:

$$5354 \int_{\|x\|_\infty \geq R_3} p_t(x|y) dx \lesssim R_3 \exp(-C'_2 R_3^2).$$

$$5355 \int_{\|x\|_\infty \geq R_3} \|\nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \lesssim R_3 \exp(-C'_2 R_3^2) \lesssim \frac{1}{\sigma_t^2} R_3^3 \exp(-C'_2 R_3^2),$$

5356 where $C'_2 = C_2 / (2 \max(1, C_2))$.

5357 Again, unlike result under **Assumption 3.1**, the explicit form of $p_t(x|y)$ in (J.1) and the upper and the
5358 lower bound of the joint distribution **Lemma J.2** automatically allow us to skip the threshold ϵ_{low} as
5359 in **Lemma I.15**.

5360 **Theorem J.1** (Approximation Score Function with Transformer under Stronger Hölder Assumption
5361 (Formal Version of **Theorem 3.2**)). Assume **Assumption 3.2** and $d_x = \Omega(\frac{\log N}{\log \log N})$. For any precision
5362 parameter $0 < \epsilon < 1$ and smoothness parameter $\beta > 0$, let $\epsilon \leq \mathcal{O}(N^{-\beta})$ for some $N \in \mathbb{N}$. For
5363 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
5364 $\mathcal{T}_{\text{score}}(x, y, t) \in \mathcal{T}_R^{h,s,r}$ such that the conditional score approximation satisfies

$$5365 \int_{\mathbb{R}^{d_x}} \|\mathcal{T}_{\text{score}}(x, y, t) - \nabla \log p_t(x|y)\|_2^2 \cdot p_t(x|y) dx = \mathcal{O}\left(\frac{B^2}{\sigma_t^2} \cdot N^{-\frac{2\beta}{d_x+d_y}} \cdot (\log N)^{\beta+1}\right).$$

5366 Notably, for $\epsilon = \mathcal{O}(N^{-\beta})$, the approximation error has the upper bound $\tilde{\mathcal{O}}(\epsilon^{2/(d_x+d_y)} / \sigma_t^2)$.
5367 The parameter bounds in the transformer network class satisfy

$$5368 \begin{aligned} 5369 & \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{\frac{3\beta(2d_x+4d+1)}{d(d_x+d_y)} + \frac{9C_\alpha(2d_x+4d+1)}{d}}\right); \\ 5370 & \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{\beta}{d_x+d_y}}\right); \\ 5371 & \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{4\beta}{d_x+d_y} + 9C_\sigma + \frac{3C_\alpha}{2}} \cdot \log N\right); \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}} L^{\frac{3}{2}}\right); \\ 5372 & \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{4\beta}{d_x+d_y} + 9C_\sigma + \frac{3C_\alpha}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N} / \sigma_t\right). \end{aligned}$$

5373 *Proof Sketch of Theorem J.1.* We decompose the integral into two terms based on **Lemma J.12**.

5374 • **(A.1): The approximation for region outside of the truncation $\|x\| > R_3$:**

5375 We give the error bound via **Lemma J.12**.

5376 • **(A.2): The approximation for region within the truncation $\|x\|_\infty \leq R_3$:**

5377 We give the error bound via **Lemma J.11**.

5378 □

5379 *Proof of Theorem 3.2.* For simplicity, we change the variable N to $N^{\frac{1}{d_x+d_y}}$ in the following subsection.
5380 We put the original form back at the end of the proof.

We take $C_x = \sqrt{\frac{2\beta}{C'_2}}$ in [Lemma J.11](#) and $R_3 = C_x \sqrt{\log N}$ in [Lemma J.12](#).

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

$$\begin{aligned}
& \int_{\mathbb{R}^{d_x}} \|\mathcal{T}_{\text{score}} - \nabla \log p_t\|_2^2 p_t dx \\
& \leq \underbrace{\int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \left(2\|\mathcal{T}_{\text{score}}\|_2^2 + 2\|\nabla \log p_t\|_2^2 \right) p_t dx}_{(A.1)} + \underbrace{\int_{\|x\|_\infty \leq \sqrt{\frac{2\beta}{C'_2} \log N}} \left(\|\mathcal{T}_{\text{score}} - \nabla \log p_t\|_2^2 \right) p_t dx}_{A.2} \\
& \lesssim \int_{\|x\|_\infty > \sqrt{\frac{2\beta}{C'_2} \log N}} \left(2 \left(\frac{\sqrt{\log N}}{\sigma_t} \right)^2 + 2\|\nabla \log p_t\|_2^2 \right) p_t dx + \frac{B^2}{\sigma_t^2} N^{-2\beta} (\log N)^{k_1+1} \\
& \hspace{15em} (\text{By } \ell_2 \text{ bound on } \mathcal{T}_{\text{score}} \text{ and } \a href="#">\text{Lemma J.11}) \\
& \lesssim 2d_x \frac{\sqrt{\log N}}{\sigma_t^2} \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{1}{2}} N^{-2\beta} + \frac{2}{\sigma_t^2} \left(\frac{2\beta}{C'_2} \log N \right)^{\frac{3}{2}} N^{-2\beta} + \frac{B^2}{\sigma_t^2} N^{-2\beta} (\log N)^{k_1+1} \\
& \hspace{15em} (\text{By } \a href="#">\text{Lemma J.12}) \\
& \lesssim \frac{B^2}{\sigma_t^2} N^{-2\beta} (\log N)^{\beta+1}. \hspace{10em} (\text{By dropping lower order term})
\end{aligned}$$

The transformer parameter norm bounds follow [Lemma J.11](#), with the replacement of N with 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 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 that $t \in [N^{-C_\alpha}, C_\sigma \log N]$.

This completes the proof. \square

K PROOF OF THE ESTIMATION RESULTS FOR CONDITIONAL DITS

Overview of Our Proof Strategy of [Theorem 3.3](#).

Step 0. Preliminaries. We introduce the mixed risk that accounts for risk with the distribution of the mask signal in [Definition K.1](#). We restate the loss function and the score matching technique in [Definition K.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 K.4](#). We bound the error from the truncation from the assumed light tail condition in [Lemma K.1](#).

Step 2. Derive the Covering Number of Transformer Network. We introduce the covering number of a given function class in [Definition K.5](#). We provide lemma detailing the calculation of the covering number for transformer architecture in [Lemma K.2](#). We derive the covering numbers under the respective parameter configurations for our two previous main results in [Lemma K.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 K.4](#).

Overview of Our Proof Strategy of [Theorem 3.4](#). We decompose the total variation into three components and we bound the separately.

Step 1. We bound the total variation distance between the true distributions evaluated at $t = 0$ and early-stopping time $t = t_0$.

Step 2. We bound the total variation between the true distribution at t_0 and the reverse process distribution using the true score function.

Step 3. We bound the total variation between the reverse process distributions using the true and estimated score functions at t_0 .

Organization. [Appendix K.1](#) includes auxiliary lemmas for supporting our proof of [Theorem 3.3](#). [Appendix K.2](#) includes the main proof of [Theorem 3.3](#). [Appendix K.5](#) includes auxiliary lemmas for supporting our proof of [Theorem 3.4](#). [Appendix K.6](#) includes the main proof of [Theorem 3.4](#).

K.1 AUXILIARY LEMMAS FOR [THEOREM 3.3](#)

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

$$\mathcal{R}(s_W) := \int_{t_0}^T \frac{1}{T-t_0} \mathbb{E}_{x_t, y} \|s_W(x_t, y, t) - \nabla \log p_t(x_t|y)\|_2^2 dt. \quad (\text{K.1})$$

Definition K.1 (Mixed Risk). The risk [\(K.1\)](#) 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 $\tau = \{\emptyset, \text{id}\}$ with $P(\tau = \emptyset) = P(\tau = \text{id}) = 0.5$:

$$\mathcal{R}_m(s_W) := \int_{t_0}^T \frac{1}{T-t_0} \mathbb{E}_{(x_t, y, \tau)} \left[\|s_W(x_t, \tau y, t) - \nabla \log p_t(x_t|\tau y)\|_2^2 \right] dt, \quad (\text{K.2})$$

Remark K.1. Given the score estimator \hat{s} trained from the empirical loss [\(G.8\)](#), 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:

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

Definition K.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:

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

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

Remark K.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 K.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

Definition K.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 the empirical risk:

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

Remark K.4. The key distinction between \mathcal{R}_m and \mathcal{L} lies in their formulations. Specifically, \mathcal{R}_m takes input x_t and compares s_W to the ground truth $\nabla \log p_t(x|y)$. In contrast, the score matching loss \mathcal{L} provides an explicit calculation based on the sample. It averages the squared difference between s_W and $\nabla \log p_t(x|x_0)$ over the sample and time interval.

Remark K.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 K.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 K.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 K.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:

$$\ell^{\text{trunc}}(x, y; s_W) := \ell(x, y; s_W) \mathbb{1}\{\|x\|_\infty \leq R_{\mathcal{T}}\}. \quad (\text{K.3})$$

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

$$\mathcal{S}(R_{\mathcal{T}}) := \{\ell(\cdot, \cdot; s_W) : \mathcal{D} \rightarrow \mathbb{R} \mid s_W \in \mathcal{T}_R^{h,s,r}\}. \quad (\text{K.4})$$

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

Lemma K.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 Assumption 3.1, we have $|\ell(x, y; s_W)| \lesssim 1/t_0$.

5562 Consider the parameter configuration in [Theorem 3.1](#), it holds:

$$5563 \mathbb{E}_{x,y} [|\ell(x, y, t) - \ell^{\text{trunc}}(x, y, s)|] \lesssim \exp(-C_2 R_{\mathcal{T}}^2) R_{\mathcal{T}} \left(\frac{1}{t_0}\right).$$

5564 Moreover, under [Assumption 3.2](#), we have $|\ell(x, y; s_W)| \lesssim \log(1/t_0)$. Consider the parameter configuration in [Theorem J.1](#), it holds:

$$5565 \mathbb{E}_{x,y} [|\ell(x, y, t) - \ell^{\text{trunc}}(x, y, s)|] \lesssim \exp(-C_2 R_{\mathcal{T}}^2) R_{\mathcal{T}} \log\left(\frac{1}{t_0}\right).$$

5572 **Step 2: Covering Number of Transformer Network Class.** We begin with the definition.

5573 **Definition K.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) $\mathcal{C} \in \mathcal{F}$ such that for any $f \in \mathcal{F}$, there exist $\hat{f} \in \mathcal{C}$ satisfying

$$5574 \max_i \|f(X_i) - \hat{f}(X_i)\| \leq \epsilon.$$

5575 Further, we define the covering number with respect to the data distribution as

$$5576 \mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|) = \sup_{\{X_i\}_{i=1}^n \sim P} \mathcal{N}(\epsilon, \mathcal{F}, \{X_i\}_{i=1}^n, \|\cdot\|).$$

5577 Next, we introduce the following lemma that provides results for the calculation of the covering number for transformer networks.

5578 **Lemma K.2** (Modified from Theorem A.17 of [Edelman et al. \(2022\)](#)).

5579 Let $\mathcal{T}_R^{h,s,r}(C_{\mathcal{T}}, C_Q^{2,\infty}, C_Q, C_K^{2,\infty}, C_K, C_V^{2,\infty}, C_V, C_O^{2,\infty}, C_O, C_E, C_{f_1}^{2,\infty}, C_{f_1}, C_{f_2}^{2,\infty}, C_{f_2}, L_{\mathcal{T}})$

5580 represent the class of functions of one transformer block satisfying the norm bound for matrix and Lipschitz property for feed-forward layers. Then for all data point $\|X\|_{2,\infty} \leq R_{\mathcal{T}}$ we have

$$5581 \log \mathcal{N}(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2)$$

$$5582 \leq \frac{\log(nL_{\mathcal{T}})}{\epsilon_c^2} \cdot \left(\alpha^{\frac{2}{3}} \left(d^{\frac{2}{3}} \left(C_F^{2,\infty} \right)^{\frac{4}{3}} + d^{\frac{2}{3}} \left(2(C_F)^2 C_{OV} C_{KQ}^{2,\infty} \right)^{\frac{2}{3}} + 2 \left((C_F)^2 C_{OV}^{2,\infty} \right)^{\frac{2}{3}} \right) \right)^3,$$

5583 where $\alpha := (C_F)^2 C_{OV} (1 + 4C_{KQ})(R_{\mathcal{T}} + C_E)$.

5584 **Remark K.7.** We modify ([Edelman et al., 2022](#), Theorem A.17) in seven aspects:

- 5585 1. We do not consider the last linear layer in the model: converting each column vector of the transformer output to a scalar. Therefore, we ignore the item related to the last linear layer in [Edelman et al. \(2022, Theorem A.17\)](#).
- 5586 2. We do not consider the normalization layer in our model. Because the normalization layer in the original proof only applies $\|\prod_{\text{norm}}(X_1) - \prod_{\text{norm}}(X_2)\|_{2,\infty} \leq \|X_1 - X_2\|_{2,\infty}$, ignoring this layer does not change the result.
- 5587 3. Our activation function is ReLU, we replace the Lipschitz upper bound of the activate function by 1.
- 5588 4. We consider the positional encoding in our work, we need to replace the upper bound $R_{\mathcal{T}}$ for the inputs with the upper bound $R_{\mathcal{T}} + C_E$. Besides, for multi-layer transformer, the original conclusion in [Edelman et al. \(2022, Theorem A.17\)](#) considers the upper bound for the $2, \infty$ -norm of inputs is 1, we add the upper bound for the inputs in [Lemma K.2](#).

- 5616 5. We use the feed-forward layer, including two linear layers and a residual layer. Thus, in
 5617 **Lemma K.2**, we replace the original upper bound for the norm of the weight matrix with the
 5618 upper bound for the norm of $I_d + W_2 W_1$. In the following, we use \mathcal{O} to estimate the log-covering
 5619 number, thus we ignore the item for I_d here for convenience. This is the same for the self-attention
 5620 layer.
 5621
 5622 6. We use multi-head attention, and we add the number of heads τ in our result, similar to (Edelman
 5623 et al., 2022, Theorem A.12).
 5624
 5625 7. In our work, we use transformer $\mathcal{T}_R^{1,4,1}$, i.e., with $h = 1$ head, $r = 4$ MLP dimension, and $s = 1$
 5626 hidden dimension, following the configuration for transformers' universality in **Theorem H.2**
 5627 and **Corollary H.2.1**. We remark that this configuration is minimally sufficient to achieve DiTs'
 5628 score approximation result **Theorem 3.1** but not necessary. More complex configurations can also
 5629 achieve transformer universality, as reported in (Hu et al., 2024; Kajitsuka and Sato, 2024; Yun
 5630 et al., 2020).

5631 With **Lemma K.2**, we derive the covering number under transformer weights configuration in
 5632 **Theorem 3.1** and **Theorem J.1**.

5633 **Lemma K.3** (Covering Number for $\mathcal{S}(R_{\mathcal{T}})$). Given $\epsilon_c > 0$ and consider $\|x\|_{\infty} \leq R_{\mathcal{T}}$. With
 5634 sample $\{x_i, y_i\}_{i=1}^n$, the ϵ_c -covering number for $\mathcal{S}(R_{\mathcal{T}})$ with respect to $\|\cdot\|_{L_{\infty}}$ under the network
 5635 configuration in **Theorem 3.1** satisfies

$$5636 \log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_{\infty}) \lesssim \frac{\log n}{\epsilon_c^2} N^{\nu_1} (\log N)^{\nu_2} (R_{\mathcal{T}})^2,$$

5637 where $\nu_1 = 172\beta/(d_x + d_y) + 104C_{\sigma}$ and $\nu_2 = 12d_x + 12\beta + 2$. Moreover, under network
 5638 configuration in **Theorem J.1**, we have

$$5639 \log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_{\infty}) \lesssim \frac{\log n}{\epsilon_c^2} N^{\nu_3} (\log N)^{10} (R_{\mathcal{T}})^2,$$

5640 where $\nu_3 = 48d\beta(L+2)(d_x + 2d + 1)/(d_x + d_y) + 144dC_{\sigma}(L+2) - 8\beta$.

5641 *Proof.* Applying **Lemma K.2**, we have

$$5642 \log \mathcal{N}(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2) \\
 5643 \leq \frac{\log n}{\epsilon_c^2} \cdot \alpha^2 \left(\underbrace{2 \left((C_F)^2 C_{OV}^2 \right)^{\frac{2}{3}}}_{\text{(i)}} + \underbrace{d^{\frac{2}{3}} \left(C_F^2 \right)^{\frac{4}{3}}}_{\text{(ii)}} + \underbrace{d^{\frac{2}{3}} \left(2(C_F)^2 C_{OV} C_{KQ}^2 \right)^{\frac{2}{3}}}_{\text{(iii)}} \right)^3, \quad (\text{K.5})$$

5644 where $\alpha := (C_F)^2 C_{OV} (1 + 4C_{KQ})(R_{\mathcal{T}} + C_E)$.

5645 Note that we drop $L_{\mathcal{T}}$ because it is inconsequential under **Assumptions 3.1** and **3.2**.

5646 • **Step A: Covering Number for Transformer with Network Configuration in **Theorem 3.1****
 5647 (under **Assumption 3.1**).

5648 Recall that from the network configuration in **Theorem 3.1**:

$$5649 \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O} \left(N^{\frac{7\beta}{d_x+d_y} + 6C_{\sigma}} \right); \\
 5650 \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O} \left(N^{-\frac{3\beta}{d_x+d_y} + 6C_{\sigma}} (\log N)^{3(d_x+\beta)} \right); \\
 5651 \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \quad \|W_V\|_{2,\infty} = \mathcal{O}(d); \\
 5652 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O} \left(N^{\frac{2\beta}{d_x+d_y} + 4C_{\sigma}} \right); \quad \|E^{\top}\|_{2,\infty} = \mathcal{O} \left(d^{\frac{1}{2}} L^{\frac{3}{2}} \right);$$

5670
5671
5672
5673
5674
5675
5676
5677
5678
5679
5680
5681
5682
5683
5684
5685
5686
5687
5688
5689
5690
5691
5692
5693
5694
5695
5696
5697
5698
5699
5700
5701
5702
5703
5704
5705
5706
5707
5708
5709
5710
5711
5712
5713
5714
5715
5716
5717
5718
5719
5720
5721
5722
5723

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{3\beta}{d_x+d_y}+2C_\sigma}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t\right).$$

Note that $W_{K,Q} = W_Q W_K^\top$, we take $\|W_Q\|_{2,\infty} \cdot \|W_K\|_{2,\infty}$ as the upper bound for $\|W_{KQ}\|_{2,\infty}$. Since W_Q, W_K share identical upper-bound, we calculate $(C_K^{2,\infty})^4$ for $(C_{K,Q}^{2,\infty})^2$. Similarly we use $\|W_O\|_{2,\infty} \cdot \|W_V\|_{2,\infty}$ as the upper bound for $\|W_{OV}\|_{2,\infty}$. Moreover, we take $C_F = \max\{C_{f_1}, C_{f_2}\}$. Since we do not impose any relation on β and C_σ here, we take $N^{3\beta/(d_x+d_y)+4C_\sigma}$ such that the upper-bound holds for both W_1 and W_2 .

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 $\log N$.

Among three terms, it is obvious that **(III)** dominates the other two. so we begin with:

$$\begin{aligned} \text{(III)} &\lesssim \left((C_F)^4 (C_{OV})^2 (C_{KQ}^{2,\infty})^2 \right)^{\frac{1}{3}} \\ &\lesssim \left(\underbrace{N^{\frac{12\beta}{d_x+d_y}+16C_\sigma}}_{(C_F)^4} \underbrace{N^{-\frac{6\beta}{d_x+d_y}+12C_\sigma} (\log N)^{6(d_x+\beta)}}_{(C_{OV})^2} \underbrace{N^{\frac{28\beta}{d_x+d_y}+24C_\sigma}}_{(C_K^{2,\infty})^4} \right)^{\frac{1}{3}}, \\ &\lesssim \left(N^{\frac{34\beta}{d_x+d_y}+52C_\sigma} (\log N)^{6(d_x+\beta)} \right)^{\frac{1}{3}}. \end{aligned}$$

Recall $\alpha := (C_F)^2 C_{OV} (1 + 4C_{KQ}) (R_{\mathcal{T}} + C_E)$,

$$\begin{aligned} \alpha^2 &\lesssim (C_F)^4 (C_{OV})^2 (C_{KQ})^2 (R_{\mathcal{T}} + C_E)^2, \\ &\lesssim \underbrace{N^{\frac{12\beta}{d_x+d_y}+16C_\sigma}}_{(C_F)^4} \underbrace{N^{-\frac{6\beta}{d_x+d_y}+12C_\sigma} (\log N)^{6(d_x+\beta)}}_{(C_{OV})^2} \underbrace{N^{\frac{28\beta}{d_x+d_y}+24C_\sigma}}_{(C_K^{2,\infty})^4} \underbrace{R_{\mathcal{T}}^2 dL^3}_{(R_{\mathcal{T}}^2 C_E^2)}, \\ &\lesssim \left(\underbrace{N^{\frac{34\beta}{d_x+d_y}+52C_\sigma} (\log N)^{6(d_x+\beta)}}_{\text{(III)}^3} (R_{\mathcal{T}})^2 \right). \end{aligned}$$

Putting all together, we obtain

$$\log \mathcal{N}\left(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2\right) \lesssim \frac{\log n}{\epsilon_c^2} N^{\frac{68\beta}{d_x+d_y}+104C_\sigma} (\log N)^{12d_x+12\beta} (R_{\mathcal{T}})^2. \quad (\text{K.6})$$

• **Step B: Covering Number for Transformer with Network Configuration in Theorem J.1 (under Assumption 3.2).**

Recall that from the network configuration in Theorem J.1

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

We derive the covering number for result under second assumption by the same procedure.

Similar to previous step, we bound **(III)** in (K.5). First, we calculate:

5724 – **Bound on** $(C_F)^4 = (C_{f_1})^4$.

$$5725 (C_{f_1})^4 \lesssim \mathcal{O}\left(N^{\frac{16\beta}{d_x+d_y}+36C_\sigma+6C_\alpha} \cdot (\log N)^4\right)$$

5728 – **Bound on** $(C_K^{2,\infty})^4$.

$$5729 (C_K^{2,\infty})^4 \lesssim N^{\frac{12\beta(2d_x+4d+1)}{d(d_x+d_y)} + \frac{36C_\alpha(2d_x+4d+1)}{d}}$$

5733 The upper-bound on **(III)** follows:

$$5735 \begin{aligned} \text{(III)} &\lesssim \left(d^2(C_{f_1})^4(C_{OV})^2(C_K^{2,\infty})^2\right)^{\frac{1}{3}}, \\ &\lesssim \left(\underbrace{N^{\frac{24\beta d_x+64\beta d+12\beta}{d(d_x+d_y)} + \frac{72C_\alpha d_x+150C_\alpha d+36C_\alpha}{d} + 36C_\sigma (\log N)^4}_{(C_{f_1})^4 \cdot (C_K^{2,\infty})^4} \underbrace{N^{-\frac{2\beta}{d_x+d_y}}}_{(C_{OV})^2}\right)^{\frac{1}{3}} \\ &\left(N^{\frac{24\beta d_x+62\beta d+12\beta}{d(d_x+d_y)} + \frac{72C_\alpha d_x+150C_\alpha d+36C_\alpha}{d} + 36C_\sigma (\log N)^4\right) \end{aligned}$$

5744 Second we bound α in **(K.5)**.

$$5745 \alpha^2 \lesssim (C_{f_1})^4(C_{OV})^2(C_{KQ})^2(R_{\mathcal{T}} + C_E)^2 \lesssim \text{(III)}^3 \cdot (R_{\mathcal{T}})^2.$$

5749 Combining **(III)** and α^2 for network configuration in **Theorem J.1**, we obtain

$$5750 \log \mathcal{N}\left(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2\right) \lesssim \frac{\log n}{\epsilon_c^2} N^{\frac{4(12\beta d_x+31\beta d+6\beta)}{d(d_x+d_y)} + \frac{12(12C_\alpha d_x+25C_\alpha \cdot d+6C_\alpha)}{d} + 72C_\sigma} (\log N)^8 \cdot (R_{\mathcal{T}})^2. \quad (\text{K.7})$$

5757 • **Step C: Covering Number under Domain Truncation.**

5758 We extend the result to the covering number for $\mathcal{S}(R_{\mathcal{T}})$ defined in **(K.4)**.

5759 First note that we obtain the score estimator from \mathcal{T}_2 by virtue of arranging x, y, t into a row vector and treating them as a sequence for execution, so we convert our $\ell_{2,\infty}$ case into ℓ_∞ as stated in **Fu et al. (2024b)** without loss of generality.

5760 For two score estimator $s_1(x, y, t), s_2(x, y, t) \in \mathcal{T}_R^{h,s,r}$ such that $\|s_1 - s_2\|_{L_\infty, \mathcal{D}} \leq \epsilon$, Proof of lemma D.3 in **Fu et al. (2024b)** shows the difference between the loss $\ell(\cdot, \cdot, s_1)$ and $\ell(\cdot, \cdot, s_2)$ in L_∞ is bounded by

$$5761 |\ell(\cdot, \cdot, s_1) - \ell(\cdot, \cdot, s_2)| \lesssim \epsilon \log N. \quad (\text{K.8})$$

5762 Therefore, by replacing ϵ_c with $\epsilon_c / \log N$ in **(K.6)** we obtain the log-covering number for transformer under **Assumption 3.1**

$$5763 \begin{aligned} \log \mathcal{N}\left(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_\infty\right) &\lesssim \frac{\log n}{\epsilon_c^2} N^{\frac{172\beta}{d_x+d_y}+104C_\sigma} (\log N)^{12d_x+12\beta+2} (R_{\mathcal{T}})^2 \\ &:= \frac{\log n}{\epsilon_c^2} N^{\nu_1} (\log N)^{\nu_2} (R_{\mathcal{T}})^2, \end{aligned}$$

5764 where $\nu_1 = 68\beta/(d_x + d_y) + 104C_\sigma$ and $\nu_2 = 12d_x + 12\beta + 2$.

Moreover, by replacing ϵ_c with $\epsilon_c / \log N$ in (K.7) we obtain the log-covering number for transformer under [Assumption 3.2](#)

$$\log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_{\infty}) = \frac{\log n}{\epsilon_c^2} N^{\nu_3} (\log N)^{10} (R_{\mathcal{T}})^2.$$

$$\text{where } \nu_3 = \frac{4(12\beta d_x + 31\beta d + 6\beta)}{d(d_x + d_y)} + \frac{12(12C_{\alpha} d_x + 25C_{\alpha} \cdot d + 6C_{\alpha})}{d} + 72C_{\sigma}.$$

This completes the proof. \square

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

Definition K.6. Let $s^{\circ} := \nabla \log p_t(x|y)$ 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:

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

Remark K.8. 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 K.4](#).

Definition K.7. Given the truncated loss function class $\mathcal{S}(R_{\mathcal{T}})$, we define its ϵ_c -covering with the minimum cardinality in the L^{∞} 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 $\|\ell_J - \ell(x_i, y_i; s_W)\|_{\infty} \leq \epsilon_c$.

Note that [Lemma K.3](#) provides the upper-bound on the ϵ_c -covering number of $\mathcal{S}(R_{\mathcal{T}})$ for score estimator trained from transformer network class. Next, we bound the difference function.

Lemma K.4 (Bound on Difference Function). Consider i.i.d training samples $\{x_{0,i}, y_i\}_{i=1}^n$ and score estimator \widehat{s} from (2.1). Under [Assumption 3.1](#) and parameter configuration in [Theorem 3.1](#), it holds:

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

where $\mathcal{N}(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2)$ is the covering number of transformer network class. Moreover, Under [Assumption 3.2](#) and parameter configuration in [Theorem J.1](#), it holds:

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

Proof. In this proof, we let $z_i := (x_{0,i}, y_i)$, $\widehat{\ell}(z_i) := \ell^{\text{trunc}}(z_i; \widehat{s})$ and $\ell^{\circ}(z_i) := \ell^{\text{trunc}}(z_i; s^{\circ})$. 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 J.1](#).

• **Step A: Rewrite the true risk.**

To derive the upper-bound of the true risk, we introduce a different set of i.i.d samples $\{x'_{0,i}, y'_i\}_{i=1}^n$ independent of the training data drawn from the same distribution.

This allows us to rewrite the true risk as:

$$\mathcal{R}_m(\widehat{s}) - \mathcal{R}_m(s^{\circ}) = \mathcal{L}(\widehat{s}) - \mathcal{L}(s^{\circ}) = \mathbb{E}_{\{x'_i, y'_i\}_{i=1}^n} \left[\frac{1}{n} \sum_{i=1}^n (\ell(x'_i, y'_i, \widehat{s}) - \ell(x'_i, y'_i, s^{\circ})) \right]. \quad (\text{K.9})$$

With (K.9), we rewrite the difference function:

$$\Delta_n(\widehat{s}, s^\circ) = \left| \frac{1}{n} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\sum_{i=1}^n \left(\widehat{\ell}(z_i) - \ell^\circ(z_i) \right) - \left(\widehat{\ell}(z'_i) - \ell^\circ(z'_i) \right) \right] \right|. \quad (\text{K.10})$$

• **Step B: Introduce the ϵ_c -covering.**

Before further decomposing (K.10), we introduce three definitions.

$$- \omega_J(z) := \ell_J(z) - \ell^\circ(z) \text{ and } \widehat{\omega}(z) := \widehat{\ell}(z) - \ell^\circ(z).$$

$$- \Omega := \max_{1 \leq J \leq \mathcal{N}} \left| \sum_{i=1}^n \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right|.$$

$$- h_J := \max\{\mathcal{A}, \sqrt{\mathbb{E}_{z'}[\ell_J(z') - \ell^\circ(z')]} \} \text{ with constant } \mathcal{A} \text{ to be chosen later.}$$

With h_j , ω_j and Ω , we start bounding (K.10) by writing

$$\begin{aligned} \Delta_n(\widehat{s}, s^\circ) &= \left| \frac{1}{n} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\sum_{i=1}^n \left(\widehat{\ell}(z_i) - \ell^\circ(z_i) \right) - \left(\widehat{\ell}(z'_i) - \ell^\circ(z'_i) \right) \right] \right| \\ &\leq \left| \frac{1}{n} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\sum_{i=1}^n (\omega_J(z_i) - \omega_J(z'_i)) \right] \right| + 2\epsilon_c \quad (\text{By Replacing } \widehat{\ell} \text{ with } \ell_J) \\ &\leq \frac{1}{n} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J \Omega] + 2\epsilon_c \quad (\text{By introducing } \Omega \text{ and } h_J) \\ &\leq \frac{1}{n} \sqrt{\mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J^2] \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2]} + 2\epsilon_c \quad (\text{By Cauchy-Schwarz inequality}) \\ &\leq \frac{1}{n} \left(\frac{n}{2} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J^2] + \frac{1}{2n} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2] \right) + 2\epsilon_c \quad (\text{By AM-GM inequality}) \\ &= \underbrace{\frac{1}{2} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J^2]}_{\text{(I)}} + \underbrace{\frac{1}{2n^2} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2]}_{\text{(II)}} + 2\epsilon_c. \end{aligned} \quad (\text{K.11})$$

– **Step B.1: Bounding (I).**

By the definition of h_J ,

$$\begin{aligned} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J^2] &\leq \mathcal{A}^2 + \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\mathbb{E}_{z'}[\omega_J^2(z)]] \\ &\leq \mathcal{A}^2 + \mathbb{E}_{z'}[\widehat{\omega}^2(z')] + 2\epsilon_c \\ &= \mathcal{A}^2 + \mathbb{E}_{\{z_i\}_{i=1}^n} [\mathcal{R}_m^{\text{trunc}}(\widehat{s})] + 2\epsilon_c. \end{aligned} \quad (\text{K.12})$$

– **Step B.2: Bounding (II).**

By Lemma K.1, we have $|\ell(z; s_W)| \lesssim \kappa$, and by the definition of Ω^2 , we write

$$\begin{aligned} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\sum_{i=1}^n \left(\frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right)^2 \right] &\leq \sum_{i=1}^n \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\left(\frac{\omega_J(z_i)}{h_J} \right)^2 + \left(\frac{\omega_J(z'_i)}{h_J} \right)^2 \right] \\ &\quad (\text{By the independence between } z_i \text{ and } z'_i) \\ &\leq \kappa \sum_{i=1}^n \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} \left[\frac{\omega_J^2(z_i)}{h_J} + \frac{\omega_J^2(z'_i)}{h_J} \right] \\ &\leq 2n\kappa. \end{aligned}$$

From the following two facts

$$* (1) \left| \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right| \leq \kappa/\mathcal{A}$$

$$* (2) \sum_{i=1}^n \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \text{ is centered}$$

we further write

$$P \left(\left(\sum_{i=1}^n \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right)^2 \geq \omega \right) = 2P \left(\left(\sum_{i=1}^n \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right) \geq \sqrt{\omega} \right) \leq 2 \exp \left(- \frac{\omega/2}{\kappa \left(2n + \frac{\sqrt{\omega}}{3\mathcal{A}} \right)} \right),$$

(By Bernstein's inequality)

for any J and $\omega \geq 0$. Therefore, we have

$$P(\Omega^2 \geq \omega) \leq \sum_{J=1}^{\mathcal{N}} P \left(\left(\sum_{i=1}^n \frac{\omega_J(z_i) - \omega_J(z'_i)}{h_J} \right)^2 \geq \omega \right) \leq 2\mathcal{N} \exp \left(- \frac{\omega/2}{\kappa \left(2n + \frac{\sqrt{\omega}}{3\mathcal{A}} \right)} \right).$$

For some $\omega_0 > 0$, we bound Ω^2 by

$$\begin{aligned} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2] &= \int_0^{\omega_0} P(\Omega^2 \geq \omega) d\omega + \int_{\omega_0}^{\infty} P(\Omega^2 \geq \omega) d\omega, & (\text{By integral identity}) \\ &\leq \omega_0 + \int_{\omega_0}^{\infty} 2\mathcal{N} \exp \left(- \frac{\omega/2}{\kappa \left(2n + \frac{\sqrt{\omega}}{3\mathcal{A}} \right)} \right) d\omega, \\ &\leq \omega_0 + 2\mathcal{N} \int_{\omega_0}^{\infty} \left\{ \exp \left(- \frac{\omega}{8n\kappa} \right) + \exp \left(- \frac{3\mathcal{A}\sqrt{\omega}}{4\kappa} \right) \right\} d\omega, \\ &\leq \omega_0 + 2\mathcal{N} \left\{ 8n\kappa \exp \left(- \frac{\omega_0}{8n\kappa} \right) + \left(\frac{8\kappa\sqrt{\omega_0}}{3\mathcal{A}} + \frac{32\kappa}{9\mathcal{A}^2} \right) \exp \left(- \frac{3\mathcal{A}\sqrt{\omega_0}}{4\kappa} \right) \right\}. \end{aligned}$$

Taking $\mathcal{A} = \sqrt{\omega_0}/6n$ and $\omega_0 = 8n\kappa \log \mathcal{N}$, we have

$$\mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2] \leq n\kappa \log \mathcal{N}. \quad (\text{K.13})$$

• **Step C: Altogether.**

Combining (K.12) and (K.13), we obtain:

$$\begin{aligned} \Delta_n(\hat{s}, s^\circ) &\leq \frac{1}{2} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [h_J^2] + \frac{1}{2n^2} \mathbb{E}_{\{z_i, z'_i\}_{i=1}^n} [\Omega^2] + 2\epsilon_c \\ &\lesssim \frac{1}{2} \mathbb{E}_{\{z_i\}_{i=1}^n} [\mathcal{R}_m^{\text{trunc}}(\hat{s})] + \frac{\kappa}{2n} \log \mathcal{N} + \frac{7}{2} \epsilon_c. \end{aligned}$$

Recall Definition K.6 and multiply the above inequality with 2, we have

$$\mathbb{E}_{\{z_i\}_{i=1}^n} [\mathcal{R}_m^{\text{trunc}}(\hat{s})] \lesssim 2\mathbb{E}_{\{z_i\}_{i=1}^n} [\widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s})] + \frac{\kappa}{n} \log \mathcal{N} + 7\epsilon_c.$$

Therefore,

$$\begin{aligned} \Delta_n(\hat{s}, s^\circ) &\lesssim \mathbb{E}_{\{z_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s}) \right] + \frac{\kappa}{n} \log \mathcal{N} + 7\epsilon_c && \text{(By Lemma K.1)} \\ &\lesssim \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}_m(\hat{s}) \right] + \kappa \left(R_{\mathcal{T}} \exp(-C_2 R_{\mathcal{T}}^2) + \frac{1}{n} \log \mathcal{N} \right) + 7\epsilon_c, \end{aligned}$$

This completes the proof. \square

K.2 PROOF OF THEOREM 3.3

Proof of Theorem 3.3. 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 J.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 K.3 and Lemma K.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:

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

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

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

$$\begin{aligned} \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\hat{s})] &= \underbrace{\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_{\{x'_i, y'_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}'_m(\hat{s}) - \widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s}) \right] \right]}_{\text{(I)}} \\ &\quad + \underbrace{\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_{\{x'_i, y'_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s}) - \widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s}) \right] \right]}_{\text{(II)}} \\ &\quad + \underbrace{\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}_m^{\text{trunc}}(\hat{s}) - \widehat{\mathcal{R}}_m(\hat{s}) \right]}_{\text{(III)}} + \underbrace{\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\widehat{\mathcal{R}}_m(\hat{s}) \right]}_{\text{(IV)}} \end{aligned}$$

• Step B: Derive the Upper Bound.

– Step B.1: Bound Each Term.

* By Lemma K.1, we have both (I), (III) $\lesssim \kappa \exp(-C_2 R_{\mathcal{T}}^2) R_{\mathcal{T}}$.

* By Lemma K.4, we have (II) \lesssim (IV) $+ \kappa \left(R_{\mathcal{T}} \exp(-C_2 R_{\mathcal{T}}^2) + \frac{1}{n} \log \mathcal{N} \right) + 7\epsilon_c$,

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

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

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

$$\begin{aligned} \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\hat{s})] &\leq 2 \min_{s, w \in \mathcal{T}_R^{h, s, r}} \int_{t_0}^T \frac{1}{T-t_0} \mathbb{E}_{x_t, y, \tau} \left[\|s(x_t, \tau y, t) - \nabla \log p_t(x_t | \tau y)\|_2^2 \right] dt \\ &\quad + \mathcal{O}\left(\frac{\kappa}{n} \log \mathcal{N}\right) + \mathcal{O}(\exp(-C_2 R_{\mathcal{T}}^2) \kappa) + \mathcal{O}(\epsilon_c). \end{aligned} \quad (\text{K.14})$$

By taking $R_{\mathcal{T}} = \sqrt{\frac{(C_\sigma + 2\beta) \log N}{C_2(d_x + d_y)}}$ along with the result in [Lemma K.3](#), we further write

$$\begin{aligned} \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}_m(\hat{s})] &\leq 2 \min_{s \in \mathcal{T}_R^{h, s, r}} \int_{t_0}^T \frac{1}{T-t_0} \mathbb{E}_{\tau, x_t, y} \left[\|s(x, \tau y, t) - \nabla \log p_t(x|y)\|_2^2 \right] dt \\ &\quad \mathcal{O}\left(\frac{\kappa}{n} \log \mathcal{N}\right) + \mathcal{O}\left(N^{-\frac{2\beta}{d_x + d_y}}\right) + \mathcal{O}(\epsilon_c). \end{aligned} \quad (\text{K.15})$$

where we invoke $\kappa \lesssim \frac{1}{t_0} = N^{C_\sigma}$ to obtain the second term on the RHS.

Step C: Altogether.

To apply the previous approximation theorems ([Theorem 3.1](#) and [Theorem J.1](#)) to the first term on the RHS of [\(K.14\)](#), we rewrite the expectation as

$$\begin{aligned} &\mathbb{E}_{x_t, y, \tau} \left[\|s(x_t, \tau y, t) - \nabla \log p_t(x_t | \tau y)\|_2^2 \right] \\ &= \frac{1}{2} \int_{\mathbb{R}^{d_x}} \|s(x, \emptyset, t) - \nabla \log p_t(x|y)\|_2^2 p_t(x) dx + \frac{1}{2} \mathbb{E}_y \left[\int_{\mathbb{R}^{d_x}} \|s(x, y, t) - \nabla \log p_t(x|y)\|_2^2 p_t(x|y) dx \right]. \end{aligned} \quad (\text{K.16})$$

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 [Assumption 3.1](#).

From [Theorem 3.1](#), we rewrite [\(K.15\)](#) as

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

Moreover, from [Lemma K.1](#) we have $\kappa = \mathcal{O}(1/t_0)$ and from [Lemma K.3](#) we have

$$\begin{aligned} \log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_\infty) &\lesssim \frac{\log n}{\epsilon_c^2} N^{\frac{172\beta}{d_x + d_y} + 104C_\sigma} (\log N)^{12d_x + 12\beta + 2} (R_{\mathcal{T}})^2 \\ &:= \frac{\log n}{\epsilon_c^2} N^{\nu_1} (\log N)^{\nu_2} (R_{\mathcal{T}})^2, \end{aligned}$$

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

Taking $N = n^{\frac{1}{\nu_1} \frac{d_x + d_y}{(d_x + d_y) + \beta}}$ and $\epsilon_c = N^{-\frac{1}{4} \frac{\nu_1 \beta}{(d_x + d_y)}}$ renders error

$$\text{* (i)} = \mathcal{O}\left(\frac{1}{t_0} (\log n)^{d_x + \frac{\beta}{2} + 1} n^{-\frac{\beta}{\nu_1(d_x + d_y) + \beta}}\right) \text{ from (K.16) and Theorem 3.1}$$

$$\text{* (ii)} = \mathcal{O}\left(n^{-\frac{2\beta}{\nu_1(d_x + d_y) + \beta}}\right)$$

$$\text{* (iii)} = \mathcal{O}\left(\kappa n^{-1} n^{\frac{1}{2} \frac{\beta}{d_x + d_y + \beta}} (\log n) n^{\frac{d_x + d_y}{d_x + d_y + \beta}} (\log n)^{\nu_2} (\log n)\right)$$

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

$$* \text{ (iv)} = \mathcal{O}\left(n^{-\frac{1}{4} \frac{\beta}{d_x+d_y+\beta}}\right)$$

The total error is bounded by

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

Step C.2: Result under **Assumption 3.2**.

With **Theorem J.1**, we further write **(K.15)** as

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

Moreover, from **Lemma K.1** we have $\kappa = \mathcal{O}(\log \frac{1}{t_0})$ and from **Lemma K.3**

$$\log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_{\infty}) = \frac{\log n}{\epsilon_c^2} N^{\nu_3} (\log N)^{10} (R_{\mathcal{T}})^2.$$

where $\nu_3 = \frac{4(12\beta d_x + 31\beta d + 6\beta)}{d(d_x+d_y)} + \frac{12(12C_{\alpha} d_x + 25C_{\alpha} \cdot d + 6C_{\alpha})}{d} + 72C_{\sigma}$.

Taking $N = n^{\frac{(d_x+d_y)}{\nu_3(d_x+d_y+2\beta)}}$ and $\epsilon_c = N^{-\frac{1}{4} \frac{\nu_3\beta}{(d_x+d_y)}}$ renders error

$$* \text{ (i)} = \mathcal{O}\left(\log \frac{1}{t_0} (\log n)^{\beta+1} n^{-\frac{1}{\nu_3} \frac{2\beta}{(d_x+d_y+2\beta)}}\right) \text{ from (K.16) and Theorem 3.1}$$

$$* \text{ (ii)} = \mathcal{O}\left(n^{-\frac{2\beta}{\nu_3(d_x+d_y+2\beta)}}\right)$$

$$* \text{ (iii)} = \mathcal{O}\left(\frac{\kappa}{n} n^{\frac{1}{2} \frac{\beta}{d_x+d_y+2\beta}} (\log n) n^{\frac{d_x+d_y}{d_x+d_y+2\beta}} (\log n)^{10} (\log n)\right)$$

Rearranging the expression we have **(iii)** = $\mathcal{O}\left(\log \frac{1}{t_0} n^{-\frac{3}{2} \frac{\beta}{d_x+d_y+2\beta}} (\log n)^{12}\right)$

$$* \text{ (iv)} = \mathcal{O}\left(n^{-\frac{1}{4} \frac{\beta}{d_x+d_y+2\beta}}\right)$$

The total error is bounded by

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

This completes the proof. □

K.3 DOMINANCE TRANSITION BETWEEN N AND $\log N$ FOR ALL NORM BOUNDS UNDER ASSUMPTION 3.1

Here we show that there is a sharp transition between the dominance of N and $\log N$ in all norm bounds for using transformers to approximate score function under Assumption 3.1 (in Theorem 3.1).

We remark that this sharp transition necessitates separate analyses for the low-dimensional region ($d_x \ll n$) in Corollaries 3.3.1 and 3.4.1.

Lemma K.5 (Dominance Transition between N and $\log N$ for All Norm Bounds). Let d_x be the feature dimension of the data. Let N be the discretization resolution of the locally diffused polynomial defined in Lemma I.1 and Remark I.1. Under Assumption 3.1, $d_x = \Theta\left(\frac{\log N}{\log \log N}\right)$ divides the dependence of N and $\log N$ into two regions for the required norm bounds on attention weights W_K, W_Q, W_O, W_1, W_2 in score approximation using transformer networks (Theorem 3.1):

- **High-Dimensional Region:** If $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$, N dominates over $\log N$.
- **Mild and Low-Dimensional Region:** If $d_x = o\left(\frac{\log N}{\log \log N}\right)$, $\log N$ dominates over N .

Proof of Lemma K.5. Recall the required parameter norm bounds for approximating score function with transformer networks from Step C of Lemma I.13. We provide a comprehensive summary of all parameter bounds involving terms dependent on N and $\log N$ from each respective operation.

- **Bound on W_Q and W_K .**

- **For ϵ_{f_1} :**

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)\cdot\frac{2dL+4d+1}{d}} \cdot (\log N)^{-3(d_x+\beta)\cdot\frac{2dL+4d+1}{d}}\right).$$

Since $d_x = dL$, N and $\log N$ balance at

$$N^{\mathcal{O}(d_x)} = (\log N)^{\mathcal{O}(d_x^2)},$$

and hence

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

- **For ϵ_{f_2} :**

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_\sigma)\cdot\frac{2dL+4d+1}{d}} \cdot (\log N)^{-(d_x+\beta)\cdot\frac{2dL+4d+1}{d}}\right).$$

Since $d_x = dL$, N and $\log N$ balance at

$$N^{\mathcal{O}(d_x)} = (\log N)^{\mathcal{O}(d_x)},$$

and hence

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

- **For $\epsilon_{\text{rec},1}$ and $\epsilon_{\text{rec},2}$:**

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)} (\log N)^{-3(d_x+\beta)}\right).$$

6156 N and $\log N$ balance at

$$6157 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6159 and hence

$$6161 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6162
6163
6164
6165 **– For $\epsilon_{\sigma,1}$:**

$$6166 \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_Q\|_{2,\infty} = \mathcal{O}\left(N^{(27\beta+18C_\sigma)}(\log N)^{-9(d_x+\beta)}\right).$$

6169 N and $\log N$ balance at

$$6171 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6172 and hence

$$6174 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6175
6176
6177
6178 **– For $\epsilon_{\sigma,3}$:**

$$6179 \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_Q\|_{2,\infty} = \mathcal{O}\left(N^{(21\beta+15C_\sigma)}(\log N)^{-6(d_x+\beta)}\right).$$

6182 N and $\log N$ balance at

$$6184 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6185 and hence

$$6187 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6188
6189
6190
6191 **• Bound on W_O .**

6192 **– For ϵ_{f_1}**

$$6193 \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(9\beta+6C_\sigma)}{d}}(\log N)^{\frac{3(d_x+\beta)}{d}}\right).$$

6197 N and $\log N$ balance at

$$6198 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6200 and hence

$$6203 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6204
6205
6206
6207 **– For ϵ_{f_2}**

$$6208 \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(3\beta+2C_\sigma)}{d}}(\log N)^{\frac{(d_x+\beta)}{d}}\right).$$

6209

6210 N and $\log N$ balance at

6211

$$6212 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6213

6214 and hence

6215

$$6216 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6217

6218

6219 – For $\epsilon_{\text{rec},1}$ and $\epsilon_{\text{rec},2}$:

6220

$$6221 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-(3\beta+6C_\sigma)}(\log N)^{d_x+\beta}\right).$$

6222

6223 N and $\log N$ balance at

6224

$$6225 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6226

6227 and hence

6228

$$6229 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6230

6231

6232 – For ϵ_{σ_1} :

6233

$$6234 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(dN^{-(9\beta+6C_\sigma)}(\log N)^{3(d_x+\beta)}\right).$$

6235

6236 N and $\log N$ balance at

6237

$$6238 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6239

6240 and hence

6241

$$6242 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6243

6244

6245 – For ϵ_{σ_2} :

6246

$$6247 \quad \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(dN^{-(7\beta+5C_\sigma)}(\log N)^{2(d_x+\beta)}\right).$$

6248

6249 N and $\log N$ balance at

6250

$$6251 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6252

6253 and hence

6254

$$6255 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6256

6257

6258

6259 • **Bound on W_1 .**

6260

6261 – For ϵ_{f_1} :

6262

$$6263 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_\sigma)}{d}}(\log N)^{-\frac{3(d_x+\beta)}{d}} \cdot (\log N)\right).$$

6264 N and $\log N$ balance at

$$6265 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6267
6268 and hence

$$6269 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6272
6273 **– For ϵ_{f_2} :**

$$6274 \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{(3\beta+2C_\sigma)}{d}} (\log N)^{\frac{(d_x+\beta)}{d}}\right).$$

6277 N and $\log N$ balance at

$$6278 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6280
6281 and hence

$$6282 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6285
6286 **– For $\epsilon_{\text{mult},1}$:**

$$6287 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(4\beta+C_\sigma)} (\log N)^{-\frac{1}{2}(d_x+\beta)}\right).$$

6290 N and $\log N$ balance at

$$6291 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6293
6294 and hence

$$6295 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6298
6299 **– For $\epsilon_{\text{rec},1}, \epsilon_{\text{rec},2}$:**

$$6300 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(6\beta+4C_\sigma)} (\log N)^{-2(d_x+\beta)}\right).$$

6303 N and $\log N$ balance at

$$6304 N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6306
6307 and hence

$$6308 d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6311
6312 **– For ϵ_{σ_1} :**

$$6313 \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)} (\log N)^{-3(d_x+\beta)} \cdot \log N\right).$$

6315
6316
6317

6318 N and $\log N$ balance at

6319

6320

$$N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6321

6322 and hence

6323

6324

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6325

6326

6327 **– For ϵ_{σ_2} :**

6328

6329

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+5C_\sigma)}(\log N)^{-2(d_x+\beta)} \cdot \log N\right).$$

6330

6331 N and $\log N$ balance at

6332

6333

$$N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6334

6335 and hence

6336

6337

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6338

6339

6340

6341 **• Bound on W_2 .**

6342

6343

6344

6345

6346 **– For ϵ_{f_1} and ϵ_{f_2} :**

6347

6348

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(9\beta+6C_\sigma)}{d}}(\log N)^{-3\frac{(d_x+\beta)}{d}}\right).$$

6349

6350

6351

N and $\log N$ balance at

6352

6353

$$N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6354

6355

and hence

6356

6357

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6358

6359

6360

6361

6362

6363

$$\|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{(3\beta+2C_\sigma)}{d}}(\log N)^{-\frac{(d_x+\beta)}{d}}\right).$$

6364

6365

6366

N and $\log N$ balance at

6367

6368

$$N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6369

6370

and hence

6371

$$d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6372 – For $\epsilon_{\text{rec},1}$ and $\epsilon_{\text{rec},2}$:

$$6374 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(3\beta+2C_\sigma)}(\log N)^{-(d_x+\beta)}\right).$$

6376 N and $\log N$ balance at

$$6378 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6380 and hence

$$6382 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6385 For ϵ_{σ_1} :

$$6387 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)}(\log N)^{-3(d_x+\beta)}\right).$$

6389 N and $\log N$ balance at

$$6391 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6393 and hence

$$6395 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6398 For ϵ_{σ_2} :

$$6400 \quad \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(7\beta+5C_\sigma)}(\log N)^{-2(d_x+\beta)}\right).$$

6402 N and $\log N$ balance at

$$6404 \quad N^{\mathcal{O}(1)} = (\log N)^{\mathcal{O}(d_x)},$$

6406 and hence

$$6408 \quad d_x = \mathcal{O}\left(\frac{\log N}{\log \log N}\right).$$

6413 This completes the proof. □

6415 K.4 PROOF OF COROLLARY 3.3.1

6417 By brute force, we know $N = \mathcal{O}(n^{d_x^\kappa})$ with¹¹ $\kappa = -2, 1$ under **Assumption 3.1**. This indicates the
6418 positive proportionality between the sample size n and the resolution N .

6419 By **Lemma K.5**, we conclude:

- 6421 • High-Dimension: $d_x = \Omega\left(\frac{\log N}{\log \log N}\right)$, and $\kappa = 1$.
- 6422
- 6423 • Mild and Low-Dimensional Region: $d_x = o\left(\frac{\log N}{\log \log N}\right)$ and $\kappa = -2$.
- 6424

6425 ¹¹The options of κ values are from the hindsight. One must compute all norm bounds to identify the available values

Low-Dimension Approximation Result. For $d_x = o(\log N/(\log \log N))$, since the dominant term in the norm bounds differs ([Lemma K.5](#)), we obtain a distinct score approximation result from [Theorem 3.1](#):

Theorem K.1 (Conditional Score Approximation under [Assumption 3.1](#) and $d_x = o(\log N/(\log \log N))$). Assume [Assumption 3.1](#) and $d_x = o(\log N/(\log \log N))$. 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

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

Notably, for $\epsilon = \mathcal{O}(N^{-\beta})$, the approximation error has the upper bound $\tilde{\mathcal{O}}(\epsilon^{1/(d_x+d_y)}/\sigma_t^2)$. The parameter bounds for the transformer network class are as follows:

$$\begin{aligned} & \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} \\ &= \mathcal{O}\left(N^{\frac{9\beta(2d_x+4d+1)}{d(d_x+d_y)} + \frac{6C_\sigma(2d_x+4d+1)}{d}} \cdot (\log N)^{-3(d_x+\beta) \cdot \frac{2dL+4d+1}{d}}\right); \\ & \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{\beta}{d_x+d_y}}\right); \\ & \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{9\beta}{d_x+d_y} + 6C_\sigma} (\log N)^{-2(d_x+\beta)+1}\right); \\ & \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{9\beta}{d_x+d_y} + 6C_\sigma} (\log N)^{-2(d_x+\beta)}\right); \\ & \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}} L^{\frac{3}{2}}\right); \mathcal{C}_T = \mathcal{O}\left(\sqrt{\log N}/\sigma_t^2\right). \end{aligned}$$

Proof of Theorem K.1. We show the proof by the following two steps.

• **Step A: Upper-Bound Selection.**

For $d_x = o(\log N/(\log \log N))$, N dominates the $\log N$ term. We set the parameter based on the order of N when N and $\log N$ coexist. By [Step C](#) in the proof of [Lemma I.13](#), we have:

– **Bound on W_Q and W_K .**

We set the parameter to the largest upper bound determined by the approximation error ϵ_{f_1} :

$$\|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma) \cdot \frac{2dL+4d+1}{d}} \cdot (\log N)^{-3(d_x+\beta) \cdot \frac{2dL+4d+1}{d}}\right).$$

– **Parameter Bound on W_O and W_V .**

We set the parameter to the largest upper bound determined by the approximation error $\epsilon_{\text{mult},2}$ and $\epsilon_{\text{rec},3}$:

$$\|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}(N^{-\beta}).$$

Note that only ϵ_{f_1} and ϵ_{f_2} involve the reshape operation. That is, approximation other than f_1 and f_2 has $\|W_V\|_2, \|W_V\|_{2,\infty} = \mathcal{O}(1)$. Therefore, we take $\mathcal{O}(\sqrt{d})$ and $\mathcal{O}(d)$ for $\|W_V\|_2$ and $\|W_V\|_{2,\infty}$ by [Lemma H.5](#) respectively.

– **Parameter Bound on W_1 .**

We set the parameter to the largest upper bound determined by the approximation error $\epsilon_{\sigma,1}$ and $\epsilon_{\sigma,2}$. That is, we take $N^{(9\beta+6C_\sigma)}$ from the former and we take $(\log N)^{-2(d_x+\beta)}$ from the latter.

$$\|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)} (\log N)^{-2(d_x+\beta)} \cdot \log N\right).$$

6480 – **Parameter Bound on W_2 .**

6481 Following the same argument for W_1 , we have

$$6482 \quad \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{(9\beta+6C_\sigma)}(\log N)^{-2(d_x+\beta)}\right).$$

6487 • **Step B: Change of Variables.**

6488 Recalling from the last step in the proof of [Theorem 3.1](#) (in [Appendix I](#)), we replace N with
6489 $N^{1/(d_x+d_y)}$ and C_σ with $(d_x+d_y)C_\sigma$ to obtain the final approximation result. Here we perform
6490 the same change of variables.
6491

6492 This completes the proof. □

6493 We compute the covering number for the function class of truncated loss $\mathcal{S}(R_{\mathcal{T}})$ (defined in [Defini-](#)
6494 [tion K.4](#)) under [Assumption 3.1](#) in low-dimensional region $d_x = o(\log N/(\log \log N))$.
6495

6496 **Lemma K.6** (Covering Number for $\mathcal{S}(R_{\mathcal{T}})$). Given $\epsilon_c > 0$ and consider $\|x\|_\infty \leq R_{\mathcal{T}}$. With
6497 sample $\{x_i, y_i\}_{i=1}^n$, the ϵ_c -covering number for $\mathcal{S}(R_{\mathcal{T}})$ with respect to $\|\cdot\|_{L_\infty}$ under the network
6498 configuration in [Theorem 3.1](#) satisfies
6499

$$6500 \quad \log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_2) \lesssim \frac{\log n}{\epsilon_c^2} N^{\nu_4} (\log N)^{\nu_5} (R_{\mathcal{T}})^2,$$

6501 where $\nu_4 = 144d\beta(L+2)(d_x+2d+1)/(d_x+d_y) + 96dC_\sigma(L+2)(d_x+2d+1) - 8\beta$ and
6502 $\nu_5 = -16d(d_x+\beta)(L+2)(3d_x+6d+2) + 2$.
6503

6504 *Proof of [Lemma K.6](#).* The proof closely follows [Lemma K.3](#). Applying [Lemma K.2](#), we calculate
6505

$$6506 \quad \begin{aligned} & \log \mathcal{N}(\epsilon_c, \mathcal{T}_R^{h,s,r}, \|\cdot\|_2) \\ & \leq \frac{\log n}{\epsilon_c^2} \cdot \alpha^2 \left(\underbrace{2 \left((C_F)^2 C_{OV}^{2,\infty} \right)^{\frac{2}{3}}}_{\text{(I)}} + \underbrace{\left(d^{\frac{2}{3}} \left(C_F^{2,\infty} \right)^{\frac{4}{3}} \right)}_{\text{(II)}} + \underbrace{d^{\frac{2}{3}} \left(2(C_F)^2 C_{OV} C_{KQ}^{2,\infty} \right)^{\frac{2}{3}}}_{\text{(III)}} \right)^3, \end{aligned}$$

6507 where [\(III\)](#) dominates [\(I\)](#) and [\(II\)](#).
6508

6509 Plug in the network configuration from [Theorem K.1](#):
6510

$$6511 \quad \begin{aligned} & \|W_Q\|_2, \|W_K\|_2, \|W_Q\|_{2,\infty}, \|W_K\|_{2,\infty} \\ & = \mathcal{O}\left(N^{\frac{9\beta(2d_x+4d+1)}{d(d_x+d_y)} + \frac{6C_\sigma(2d_x+4d+1)}{d}} \cdot (\log N)^{-3(d_x+\beta) \cdot \frac{2dL+4d+1}{d}}\right); \\ & \|W_V\|_2 = \mathcal{O}(\sqrt{d}); \|W_V\|_{2,\infty} = \mathcal{O}(d); \|W_O\|_2, \|W_O\|_{2,\infty} = \mathcal{O}\left(N^{-\frac{\beta}{d_x+d_y}}\right); \\ & \|W_1\|_2, \|W_1\|_{2,\infty} = \mathcal{O}\left(N^{\frac{9\beta}{d_x+d_y} + 6C_\sigma} (\log N)^{-2(d_x+\beta)+1}\right); \\ & \|W_2\|_2, \|W_2\|_{2,\infty} = \mathcal{O}\left(N^{\frac{9\beta}{d_x+d_y} + 6C_\sigma} (\log N)^{-2(d_x+\beta)}\right); \\ & \|E^\top\|_{2,\infty} = \mathcal{O}\left(d^{\frac{1}{2}} L^{\frac{3}{2}}\right); C_{\mathcal{T}} = \mathcal{O}\left(\sqrt{\log N}/\sigma_t^2\right). \end{aligned}$$

6512 Note that $W_{K,Q} = W_Q W_K^\top$, we take $\|W_Q\|_{2,\infty} \cdot \|W_K\|_{2,\infty}$ as the upper bound for $\|W_{KQ}\|_{2,\infty}$.
6513 Since W_Q, W_K share identical upper-bound, we calculate $(C_K^{2,\infty})^4$ for $(C_{K,Q}^{2,\infty})^2$. Similarly we use
6514 $\|W_O\|_{2,\infty} \cdot \|W_V\|_{2,\infty}$ as the upper bound for $\|W_{OV}\|_{2,\infty}$. Moreover, we take $C_F = \max\{C_{f_1}, C_{f_2}\}$.
6515

6534 • **Bound on $C_F^4 = (C_{f_2})^4$:**

6535

6536

6537

6538

6539 • **Bound on $(C_K^{2,\infty})^4$:**

6540

6541

6542

6543

6544

6545

6546

6547

6548

6549

6550

6551

6552

6553

6554

6555

6556

6557

6558

6559

6560

6561

6562

6563

6564

6565

6566

6567

6568

6569

6570

6571

6572

6573

6574

6575

6576

6577

6578

6579

6580

6581

6582

6583

6584

6585

6586

6587

The bound on **(III)** follows:

$$\begin{aligned} \text{(III)} &\lesssim \left((C_{f_2})^4 (C_{OV})^2 (C_{KQ}^{2,\infty})^2 \right)^{\frac{1}{3}} \\ &\lesssim \left(\underbrace{N^{\frac{36\beta(2d_x+5d+1)}{d(d_x+d_y)} + \frac{24C_\sigma(2d_x+5d+1)}{d}}}_{(C_{f_2})^4 \cdot (C_K^{2,\infty})^4} (\log N)^{-\frac{(d_x+\beta)(24dL+56d+12)}{d}} \cdot \underbrace{N^{-2\beta}}_{(C_{OV})^2} \right)^{\frac{1}{3}}. \end{aligned}$$

Moreover, $\alpha := (C_F)^2 C_{OV} (1 + 4C_{KQ})(R_{\mathcal{T}} + C_E)$, we have:

$$\alpha^2 \lesssim (C_{f_1})^4 (C_{OV})^2 (C_{KQ})^2 (R_{\mathcal{T}} + C_E)^2 \lesssim \text{(III)}^3 \cdot R_{\mathcal{T}}^2.$$

By the **Step C** in **Lemma K.3**, we extend the log-covering number of transformer to the truncated loss $\mathcal{S}(R_{\mathcal{T}})$ with $\|x\|_\infty \leq R_{\mathcal{T}}$ by replacing ϵ_c with $\epsilon_c/\log N$.

Combining **(III)** and α^2 for network configuration in **Theorem J.1**, we obtain:

$$\begin{aligned} \log \mathcal{N}(\epsilon_c, \mathcal{S}(R_{\mathcal{T}}), \|\cdot\|_2) &\lesssim N^{\frac{72\beta(2d_x+5d+1)}{d(d_x+d_y)} + \frac{48C_\sigma(2d_x+5d+1)}{d} - 4\beta} (\log N)^{-\frac{8(d_x+\beta)(6dL+14d+3)}{d} + 2} \cdot (R_{\mathcal{T}})^2 \\ &:= \frac{\log n}{\epsilon_c^2} N^{\nu_4} (\log N)^{\nu_5} (R_{\mathcal{T}})^2, \end{aligned}$$

where $\nu_4 = \frac{72\beta(2d_x+5d+1)}{d(d_x+d_y)} + \frac{48C_\sigma(2d_x+5d+1)}{d} - 4\beta$ and $\nu_5 = -\frac{8(d_x+\beta)(6dL+14d+3)}{d} + 2$.

This completes the proof. \square

Proof of Corollary 3.3.1. The proof closely follows the high-dimensional result where $d_x = \Omega(\log N/(\log \log N))$ in **Appendix K.2**. The only distinction lies in the covering number with transformer network (defined in **Definition K.5**), characterized by ν_i with $i \in [5]$. Therefore, we replace ν_1, ν_2 in **Theorem 3.3** with ν_4 and ν_5 .

Specifically, for score estimation under **Assumption 3.1**, by taking $N = n^{\frac{1}{\nu_4} \cdot \frac{d_x+d_y}{\beta+d_x+d_y}}$, $t_0 = N^{-C_\sigma} < 1$ and $T = C_\alpha \log n$, it holds

$$\begin{aligned} \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})] &= \mathcal{O} \left(\frac{1}{t_0} n^{-\frac{1}{\nu_4} \cdot \frac{\beta}{d_x+d_y+\beta}} (\log n)^{\nu_5+2} \right) \\ &= \mathcal{O} \left(\frac{1}{t_0} n^{-\frac{1}{\nu_4} \cdot \frac{\beta}{d_x+d_y+\beta}} \right), \quad (n \text{ term surpasses } \log n \text{ term}) \end{aligned}$$

$\nu_4 = \frac{72\beta(2d_x+5d+1)}{d(d_x+d_y)} + \frac{48C_\sigma(2d_x+5d+1)}{d} - 4\beta$ and $\nu_5 = -\frac{8(d_x+\beta)(6dL+14d+3)}{d} + 2$.

This completes the proof. \square

6588 K.5 AUXILIARY LEMMAS FOR **THEOREM 3.4**.
6589

6590 We give the following two lemmas serving as the key components in the proof of **Theorem 3.4**.
6591

6592 **Lemma K.7** (Proposition D.1 of [Oko et al. \(2023\)](#), Lemma D.4 of [Fu et al. \(2024b\)](#) and also [Chen et al. \(2022\)](#)). Consider probability distribution p_0 and two stochastic processes $h = \{h_t\}_{t \in [0, T]}$ and
6593 $h' = \{h'_t\}_{t \in [0, T]}$ that satisfy the following SDE respectively
6594

$$\begin{aligned} dh_t &= b(h_t, t)dt + dW_t \quad h_0 \sim p_0 \\ dh'_t &= b'(h'_t, t)dt + dW_t \quad h'_0 \sim p_0. \end{aligned}$$

6598 Plus denote the distribution of the two processes at time t as p_t and p'_t . Then suppose
6599

$$\int_x p_t(x) \|(b - b')(x, t)\| dx \leq C \quad (\text{K.17})$$

6603 holds for any $t \in [0, T]$, then we have
6604

$$\text{KL}(p_T \parallel p'_T) = \int_0^T \frac{1}{2} \int_x p_t(x) \|(b - b')(x, t)\| dx dt$$

6608 The bound for KL divergence stems from Girsanov's Theorem, with the extension to the case where
6609 the Novikov's condition is replaced with (K.17) by [Chen et al. \(2022\)](#). Moreover, we need the
6610 following lemma to bound to total variation.

6611 **Lemma K.8** (Lemma D.5 of [Fu et al. \(2024b\)](#)). Assume [Assumption 3.1](#) or [Assumption 3.2](#). For any
6612 $y \in [0, 1]^{d_y}$ we have
6613

$$\text{TV}(P_0(\cdot|y), P_{t_0}(\cdot|y)) = \mathcal{O}\left(\sqrt{t_0} \log^{\frac{d_x+1}{2}}\left(\frac{1}{t_0}\right)\right).$$

6616 With the above lemmas and discussion, we begin the proof of **Theorem 3.4**.
6617

6618 K.6 MAIN PROOF OF **THEOREM 3.4**
6619

6620 *Proof of **Theorem 3.4**.* Given label y , we let $\widehat{P}_{t_0}(\cdot|y)$ denote the data distribution with early-stopped
6621 time t_0 generated by the reverse process with the score estimator from transformer network class.
6622

6623 The decomposition of the total variation between the processes driven by the ground truth and the
6624 score estimator follows

$$\text{TV}\left(P(\cdot|y), \widehat{P}_{t_0}(\cdot|y)\right) \lesssim \underbrace{\text{TV}\left(P(\cdot|y), P_{t_0}(\cdot|y)\right)}_{\text{(I)}} + \underbrace{\text{TV}\left(P_{t_0}(\cdot|y), \widetilde{P}_{t_0}(\cdot|y)\right)}_{\text{(II)}} + \underbrace{\text{TV}\left(\widetilde{P}_{t_0}(\cdot|y), \widehat{P}_{t_0}(\cdot|y)\right)}_{\text{(III)}}$$

6630 • **Step A: Derive the Upper Bound**

6631 – **Step A.1: Bounding (I).**

6633 From [Lemma K.8](#) we have $\text{TV}\left(P(\cdot|y), \widetilde{P}_{t_0}(\cdot|y)\right) = \mathcal{O}\left(\sqrt{t_0} \log^{\frac{d_x+1}{2}}\left(\frac{1}{t_0}\right)\right)$.
6634

6635 – **Step A.2: Bounding (II).**

6636 We use the following process that represents the reverse process starting with standard Gaussian.
6637

$$d\widetilde{X}_t^{\leftarrow} = \left[\frac{1}{2} d\widetilde{X}_t^{\leftarrow} + \nabla \log p_{T-t}(\widetilde{X}_t^{\leftarrow} | y) \right] dt + d\overline{W}_t \quad \widetilde{X}_0^{\leftarrow} \sim N(0, I_{d_x}).$$

6641 The distribution of $\widetilde{X}_t^{\leftarrow}$ conditioned on the label y is denoted by $\widetilde{P}_{T-t}(\cdot|y)$.

Next, by Data Processing Inequality and Pinsker's Inequality (Canonne, 2022, Lemma 2) we have

$$\begin{aligned} \text{TV} \left(P_{t_0}(\cdot|y), \tilde{P}_{t_0}(\cdot|y) \right) &\lesssim \sqrt{\text{KL}(P_{t_0}(\cdot|y) \parallel \tilde{P}_{t_0}(\cdot|y))} \\ &\lesssim \sqrt{\text{KL}(P_T(\cdot|y) \parallel N(0, I_{d_x}))} \\ &\lesssim \sqrt{\text{KL}(P(\cdot|y) \parallel N(0, I_{d_x}))} \exp(-T). \end{aligned} \quad (\text{K.18})$$

Therefore for (II), from (K.18) we have

$$\begin{aligned} \text{TV} \left(P_{t_0}(\cdot|y), \tilde{P}_{t_0}(\cdot|y) \right) &\lesssim \sqrt{\text{KL}(P(\cdot|y) \parallel N(0, I_{d_x}))} \exp(-T) \\ &\lesssim \exp(-T) \end{aligned}$$

– **Step A.3: Bounding (III).**

From (K.18) and Lemma K.7, we have

$$\text{TV} \left(\tilde{P}_{t_0}(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \lesssim \sqrt{\int_{t_0}^T \frac{1}{2} \int_x p_t(x|y) \|\hat{s}(x, y, t) - \nabla \log p_t(x|y)\|^2 dx dt.}$$

• **Step B: Altogether.**

Combining (I) (II) and (III), we take the expectation to the total variation with respect to y

$$\begin{aligned} &\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \\ &\lesssim \sqrt{t_0} \log \frac{d_x+1}{2} \left(\frac{1}{t_0} \right) + \exp(-T) + \sqrt{\int_{t_0}^T \frac{1}{2} \int_x p_t(x|y) \|\hat{s}(x, y, t) - \nabla \log p_t(x|y)\|^2 dx dt} \\ &\hspace{15em} (\text{By Jensen's inequality}) \\ &\lesssim \sqrt{t_0} \log \frac{d_x+1}{2} \left(\frac{1}{t_0} \right) + \exp(-T) + \sqrt{\frac{T}{2} \mathcal{R}(\hat{s})}. \end{aligned}$$

Lastly, take the expectation with respect to the sample $\{x_i, y_i\}_{i=1}^n$ and take $T = C_\alpha \log n$ we have

$$\begin{aligned} &\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \right] \\ &\lesssim \sqrt{t_0} \log \frac{d_x+1}{2} \left(\frac{1}{t_0} \right) + n^{-C_\alpha} + \sqrt{\log n} \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\sqrt{\mathcal{R}(\hat{s})} \right] \quad (\text{By Jensen's Inequality}) \\ &\lesssim \underbrace{\sqrt{t_0} \log \frac{d_x+1}{2} \left(\frac{1}{t_0} \right)}_{\text{(i)}} + n^{-C_\alpha} + \underbrace{\sqrt{\log n} \sqrt{\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} [\mathcal{R}(\hat{s})]}}_{\text{(ii)}} \end{aligned}$$

– **Step B.1: Result under Assumption 3.1.**

We apply Theorem 3.3 and setting $C_\alpha = \frac{2\beta}{d_x+d_y+2\beta}$ and $t_0 = n^{-\beta/(d_x+d_y+\beta)}$, we further write the above expression into

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \right]$$

$$\lesssim \underbrace{n^{-\frac{\beta}{2(d_x+d_y+\beta)}} (\log n)^{\frac{d_x+1}{2}}}_{(i)} + n^{-\frac{2\beta}{d_x+d_y+2\beta}} + \underbrace{(\log n)^{\frac{1}{2}} \left(\frac{1}{t_0} n^{-\frac{\beta}{\nu_1(d_x+d_y+\beta)}} (\log n)^{\nu_2+2} \right)^{\frac{1}{2}}}_{(ii)}$$

Therefore, under [Assumption 3.1](#) we have

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{\beta}{2(\nu_1-1)(d_x+d_y+\beta)}} (\log n)^{\frac{\nu_2}{2} + \frac{3}{2}} \right)$$

– **Step B.2: Result under [Assumption 3.2](#).**

We apply [Theorem 3.3](#) and set $t_0 = n^{-\frac{d\beta}{d_x+d_y+2\beta}-1}$. Note that we have

$$\sqrt{t_0} \left(\log \frac{1}{t_0} \right)^{\frac{d_x+1}{2}} \lesssim n^{-\frac{2\beta}{d_x+d_y+2\beta}}.$$

We further write

$$\begin{aligned} & \mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \right] \\ & \lesssim \underbrace{n^{-\frac{2\beta}{d_x+d_y+2\beta}}}_{(i)} + n^{-\frac{2\beta}{d_x+d_y+2\beta}} + \underbrace{(\log n)^{\frac{1}{2}} \left(\log \frac{1}{t_0} n^{-\frac{1}{\nu_3} \frac{\beta}{d_x+d_y+2\beta}} (\log n)^{\max(10, \beta+1)} \right)^{\frac{1}{2}}}_{(ii)}. \end{aligned}$$

Therefore we have

$$\mathbb{E}_{\{x_i, y_i\}_{i=1}^n} \left[\mathbb{E}_y \left[\text{TV} \left(P(\cdot|y), \hat{P}_{t_0}(\cdot|y) \right) \right] \right] = \mathcal{O} \left(n^{-\frac{1}{2\nu_3} \frac{\beta}{d_x+d_y+2\beta}} (\log n)^{\max(6, (\beta+3)/2)} \right)$$

This completes the proof. \square

K.7 PROOF OF [COROLLARY 3.4.1](#)

Proof of [Corollary 3.4.1](#). The proof closely follows the high-dimensional result where $d_x = \Omega(\log N / (\log \log N))$ in [Appendix K.2](#). The only distinction lies in the covering number with transformer network (defined in [Definition K.5](#)), characterized by ν_i with $i \in [5]$. Therefore, we replace ν_1, ν_2 in [Theorem 3.4](#) with ν_4 and ν_5 . This completes the proof. \square