
On approximation and estimation of Schrödinger potentials without the curse of dimensionality

Anonymous Authors¹

Abstract

We examine generative modelling approaches based on the construction of Schrödinger bridges between Gaussian noise and a target distribution. It is known that the solution of the dynamic Schrödinger problem is a diffusion process with a drift associated with Doob’s h-transform of a Schrödinger potential. Although its accurate restoration from finite samples is crucial for reliable, high-quality data generation, the existing literature lacks theoretical guarantees regarding this question. In our work, we establish theoretical upper bounds on the complexity of Schrödinger potential approximation and estimation via neural networks. These bounds are determined by the effective dimension of the target distribution. To our knowledge, this is the first result demonstrating that generative modelling methods based on Schrödinger bridges and stochastic optimal control can escape the curse of dimensionality.

1. Introduction

Denoising diffusion (Ho et al., 2020; Song et al., 2021) and flow-based (Liu et al., 2023; Lipman et al., 2023) models showed themselves as a powerful tool for modelling complex multimodal distributions. The idea of using diffusion processes for data generation and translation was further developed in approaches based on Schrödinger bridges (De Bortoli et al., 2021; Shi et al., 2023; Korotin et al., 2024; Gushchin et al., 2024) and stochastic optimal control (SOC) (Domingo-Enrich et al., 2025; Rapakoulias et al., 2025; Puchkin et al., 2025), which have recently gained much popularity. In general, the Schrödinger bridge between two probability distributions P_0 and P_T on \mathbb{R}^D is a stochastic process $\{X_t^* : 0 \leq t \leq T\}$ such that $X_0^* \sim P_0$,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

$X_T^* \sim P_T$, and X_t deviates from a given base (or reference) process $\{X_t^0 : 0 \leq t \leq T\}$ as less as possible (in terms of the Kullback-Leibler divergence between the corresponding path measures). In what follows, we will assume that the source and the target distributions are absolutely continuous with the densities p_0 and p_T , respectively, and $\{X_t^0 : 0 \leq t \leq T\}$ is the Ornstein-Uhlenbeck (OU) process: $X_0 \sim p_0$,

$$dX_t^0 = -bX_t^0 dt + \sqrt{\gamma} dW_t, \quad 0 < t < T. \quad (1)$$

Here $b > 0$ and $\gamma > 0$ are some parameters and W_t is the standard Brownian motion in \mathbb{R}^D . Such a choice of the reference process is motivated by exponentially fast mixing and a recent paper (Puchkin et al., 2026) demonstrating advantages of the OU process over the standard scaled Wiener process $\{\sqrt{\gamma} dW_t : 0 \leq t \leq T\}$ in unconditional generation and data-to-data translation tasks. It is known that (see, for instance, Theorem 2.4 in (Léonard, 2013)) the joint density of X_0^* and X_T^* is given by

$$\pi^*(x, y) = \nu_0^*(x) q_T(y | x) \nu_T^*(y), \quad (2)$$

where

$$q_t(y | x) = (\sqrt{2\pi\gamma\sigma_t})^{-D} \exp \left\{ -\frac{\|y - e^{-bt}x\|^2}{2\gamma\sigma_t^2} \right\},$$
$$\sigma_t^2 = \frac{1 - e^{-2bt}}{2b}, \quad (3)$$

is the transition density of the base process and (ν_0^*, ν_T^*) is a uniquely determined (up to positive multiplicative constants) pair of σ -finite measures on \mathbb{R}^D referred to as *Schrödinger potentials*. The right potential ν_T^* plays a crucial role in generative modelling, since the solution of the Schrödinger bridge problem (as well as the optimally controlled process in the SOC approach) admits a representation

$$dX_t^* = -bX_t^* dt + \gamma \nabla \log \mathcal{T}_{T-t}[\nu_T^*](X_t^*) dt + \sqrt{\gamma} dW_t,$$

where \mathcal{T}_t stands for the Ornstein-Uhlenbeck operator. For any measurable function f on \mathbb{R}^D , it is defined as

$$\mathcal{T}_t[f](x) = \int_{\mathbb{R}^D} f(y) q_t(y | x) dy, \quad x \in \mathbb{R}^D. \quad (4)$$

Unfortunately, as it will become clear a bit later, theoretical study of ν_T^* remains quite limited.

Related work. The factorization (2) is the core of the Sinkhorn algorithm, which was extensively studied in the literature. Some significant results on its rates of convergence appeared in the last years (Conforti et al., 2023; Chizat et al., 2024; Chiarini et al., 2024). However, in generative modelling, the target density p_T is accessible only through i.i.d. samples Y_1, \dots, Y_n , while (Conforti et al., 2023; Chizat et al., 2024; Chiarini et al., 2024) assume that the source and the target distributions are known. One of the ways to overcome this issue is to apply the Sinkhorn algorithm to the empirical measure associated with the training sample, as, for instance, in (Pooladian and Niles-Weed, 2025). This approach has an evident drawback that the estimated potential (or the corresponding drift) requires additional smoothing. In (Pooladian and Niles-Weed, 2025), the authors use nonparametric kernel-type estimates for this purpose, which usually suffer from the curse of dimensionality.

Several recent papers (Korotin et al., 2024; Puchkin et al., 2025; 2026; Belomestny et al., 2026) suggested procedures for direct estimation of the continuous potential ν_T^* or a related quantity p_T/ν_T^* based on empirical risk minimization. The framework of (Puchkin et al., 2025; Belomestny et al., 2026) naturally incorporates neural networks, which proved their ability to adapt to the intrinsic data structure and to succeed in high-dimensional practical tasks. In addition, these methods directly optimize the corresponding objectives, avoiding splitting into smaller subproblems, which is usual for iterative proportional (De Bortoli et al., 2021; Vargas et al., 2021; Chen et al., 2022) or Markovian (Chen et al., 2023; Peluchetti, 2023) fitting.

Rates of convergence and the curse of dimensionality.

The ERM-based approach (Korotin et al., 2024; Puchkin et al., 2025; Belomestny et al., 2026) is supported by rigorous theoretical study. In (Korotin et al., 2024), (Puchkin et al., 2025), and (Puchkin et al., 2026), the authors elaborate on the statistical error and report their results in terms of the excess risk and complexity of the learnable class of potentials, leaving analysis of the approximation error out of the scope. This gap was recently addressed by (Belomestny et al., 2026), who studied an approximation of the ratio p_T/ν_T^* in the $L^2(p_T)$ norm. However, our work focuses on approximating $\log \nu_T^*$ with respect to KL divergence rather than the ratio itself, addressing a distinct problem. To our knowledge, (Belomestny et al., 2026) were the first ones who derived non-asymptotic $\mathcal{O}(\log^D(n)/\sqrt{n})$ rates of convergence. Unfortunately, the unfavourable dependence on the ambient dimension D is inevitable, unless we impose additional assumptions on p_T . While the ability of deep generative models to adapt to the data structure was extensively studied in the context of denoising diffusion models (Chen et al., 2023; Tang and Yang, 2024; Azan-

gulov et al., 2024; Yakovlev and Puchkin, 2025; Yakovlev et al., 2025), *there are no results on approximation or estimation of Schrödinger potentials without the curse of dimensionality*. The present paper aims to make the first step towards filling this gap.

Contributions. We consider an affine-subspace model, which has recently got considerable attention in the context of score estimation and generative diffusion models (Sasaki et al., 2016; Chen et al., 2023; Oko et al., 2023). Our main contributions are following.

- We prove that the right log-potential $\log \nu_T^*$ can be approximated with any prescribed accuracy $\varepsilon > 0$ in terms of KL divergence via a feedforward neural network with polylogarithmic in ε^{-1} number of weights (Theorem 3.4). This represents the first such result in the literature that avoids the curse of dimensionality.
- We prove an $\mathcal{O}(\text{polylog}(n)/n)$ rate of convergence in KL divergence under the affine-subspace model (Theorem 3.6), thereby escaping the curse of dimensionality (Theorem 3.4). We substantially improve over the $\mathcal{O}(n^{-2/(d+5)})$ bound of (Chen et al., 2023) and the $\mathcal{O}(n^{-1/2})$ rate of (Pooladian and Niles-Weed, 2025). In addition, our bound does not depend on the stopping time provided that the underlying data distribution possesses intrinsic variability.
- The proofs of Theorems 3.4 and 3.6 rely on sharp bounds on the $\|\cdot\|_\infty$ -norms of higher-order derivatives of $\log p_T$ and $\log \nu_T^*$ (see Lemmata 5.1 and 5.2), which could be of independent interest.

Notation. We use the following notation throughout the paper. For a vector v and a tensor A with real entries, $\|v\|_\infty$ and $\|A\|_\infty$ denote the maximum absolute value of their entries. Similarly, for a vector v and a matrix A , the number of non-zero entries is denoted by $\|v\|_0$ and $\|A\|_0$, respectively. We often replace $\max\{a, b\}$ and $\min\{a, b\}$ by shorter $a \vee b$ and $a \wedge b$, respectively. The operation \otimes stands for the outer product. For a function $f : \Omega \rightarrow \mathbb{R}$, $\|f\|_{L^\infty(\Omega)} := \sup_{x \in \Omega} |f(x)|$. For $s \geq 1$, the $L_s(p)$ -norm of a measurable function f is $\|f\|_{L_s(p)} := (\int |f(x)|^s p(x) dx)^{1/s}$. For two probability densities p and q on \mathbb{R}^d , such that $p \ll q$, the Kullback–Leibler (KL) divergence is $\text{KL}(p, q) := \int p(x) \log(p(x)/q(x)) dx$. For any $s \geq 1$, the Orlicz ψ_s -norm of a random variable is defined as

$$\|\xi\|_{\psi_s} = \inf \left\{ u > 0 : \mathbb{E} e^{|\xi|^{s/u^s}} \leq 2 \right\}.$$

For a random vector $X \in \mathbb{R}^d$, we define $\|X\|_{\psi_2} := \sup_{v \in S^{d-1}} \|\langle v, X \rangle\|_{\psi_2}$, where $S^{d-1} := \{v \in \mathbb{R}^d : \|v\| = 1\}$.

1} is the Euclidean unit sphere. For some $\Omega \subseteq \mathbb{R}^r$, we define the class of β -Hölder functions with parameter H as

$$\mathcal{H}^\beta(\Omega, H) = \left\{ f : \Omega \rightarrow \mathbb{R} : \sum_{|\mathbf{k}| < \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)} + \sum_{|\mathbf{k}| = \lfloor \beta \rfloor} \sup_{x \neq y} \frac{|\partial^{\mathbf{k}} f(x) - \partial^{\mathbf{k}} f(y)|}{\|x - y\|_\infty^{\beta - \lfloor \beta \rfloor}} \leq H \right\}.$$

Finally, $f \lesssim g$ and $g \gtrsim f$ are equivalent to $f = \mathcal{O}(g)$, and $f \asymp g$ means that $f \lesssim g \lesssim f$.

Paper structure. The rest of the paper is organised as follows. Section 2 introduces the Schrödinger-bridge formulation, our notation, and the function classes we will use. Section 3 states the main approximation result and its statistical corollary, including a comparison with (Chen et al., 2023). Section 4 discusses limitations and an extension to factorisable density models. Section 5 provides high-level insights of our proof strategy. The proofs of all auxiliary results are deferred to the appendix.

2. Problem setup

We consider a problem of unconditional data generation from Gaussian noise. That is, given i.i.d. samples Y_1, \dots, Y_n drawn according to an unknown probability density \mathfrak{p}_T over \mathbb{R}^D , our goal is to produce a new sample from approximately the same distribution. We use the Schrödinger bridge framework for this purpose. Let \mathfrak{p}_0 be the density of the standard Gaussian distribution $\mathcal{N}(0, I_D)$ in \mathbb{R}^D . Given a reference stochastic process $\{X_t^0 : 0 \leq t \leq T\}$ obeying the dynamics (1), we aim to construct $\{X_t^* : 0 \leq t \leq T\}$ of the form

$$dX_t^* = -bX_t^* dt + u(t, X_t^*) dt + \sqrt{\gamma} dW_t, \quad 0 < t < T,$$

such that $X_0^* \sim \mathfrak{p}_0$, $X_T^* \sim \mathfrak{p}_T$, and the Kullback-Leibler divergence between the path measures of $\{X_t^* : 0 \leq t \leq T\}$ and $\{X_t^0 : 0 \leq t \leq T\}$ is as small as possible. It is known that the optimal $u(x, t)$ is given by

$$u(x, t) = \gamma \nabla \log \mathcal{T}_{T-t}[\nu_T^*](x), \quad (5)$$

where ν_T^* and \mathcal{T}_{T-t} stand for the right Schrödinger potential and the Ornstein-Uhlenbeck operator defined in (2) and (4), respectively. Let us note that $u(x, t)$ from (5) also solves a stochastic optimal control problem (Dai Pra, 1991).

In our paper, we are interested in quantitative bounds on complexity of approximation and estimation of ν_T^* in the presence of low-dimensional data structure via neural networks, in particular, via feedforward neural networks with ReLU activation functions. Let us recall that the ReLU activation is given by $\text{ReLU}(x) = 0 \vee x$. For a vector $b \in \mathbb{R}^r$,

we define the shifted activation function $\text{ReLU}_b(x) = (\text{ReLU}(x_1 - b_1), \dots, \text{ReLU}(x_r - b_r))^\top$, where $x \in \mathbb{R}^r$. Given the number of layers $L \in \mathbb{N}$ and an architecture vector $W = (W_0, \dots, W_L) \in \mathbb{R}^{L+1}$, a feedforward neural network $f : \mathbb{R}^{W_0} \rightarrow \mathbb{R}^{W_L}$ is a composition of linear and nonlinear transforms of the form

$$f(x) = -b_L + A_L \circ \text{ReLU}_{b_{L-1}} \circ A_{L-1} \circ \dots \circ A_2 \circ \text{ReLU}_{b_1} \circ A_1 \circ x, \quad x \in \mathbb{R}^{W_0}. \quad (6)$$

Here, $A_j \in \mathbb{R}^{W_{j-1} \times W_j}$ is a weight matrix and $b_j \in \mathbb{R}^{W_j}$ is a bias vector for each $j \in \{1, \dots, L\}$. Finally, we introduce the class of feedforward neural networks with at most $S \in \mathbb{N}$ non-zero weights, maximum width $\|W\|_\infty \in \mathbb{N}$, and weight magnitudes bounded by $B > 0$ as follows:

$$\text{NN}(L, W, S, B) = \left\{ f \text{ as in (6)} : \sum_{j=1}^L (\|A_j\|_0 \vee \|b_j\|_0) \leq S, \max_{1 \leq j \leq L} \{\|A_j\|_\infty \vee \|b_j\|_\infty\} \leq B \right\}.$$

Once we constructed $\tilde{\nu}_T$ approximating ν_T^* , we can consider the corresponding diffusion

$$d\tilde{X}_t = -b\tilde{X}_t dt + \gamma \nabla \log \mathcal{T}_{T-t}[\tilde{\nu}_T](\tilde{X}_t) dt + \sqrt{\gamma} dW_t,$$

where $\tilde{X}_0 \sim \mathfrak{p}_0$. We measure the quality of $\tilde{\nu}_T$ with the Kullback-Leibler divergence between the corresponding endpoint marginals of X_T^* and \tilde{X}_T . Let us note that, in view of (2), it holds that

$$\begin{aligned} \mathfrak{p}_0(x) &= \int_{\mathbb{R}^D} \pi^*(x, y) dy \\ &= \nu_0^*(x) \int_{\mathbb{R}^D} \mathfrak{q}_T(y | x) \nu_T^*(y) dy \\ &= \nu_0^*(x) \mathcal{T}_T[\nu_T^*](x) \end{aligned}$$

and then

$$\begin{aligned} \mathfrak{p}_T(x) &= \int_{\mathbb{R}^D} \pi^*(x, y) dx \\ &= \nu_T^*(y) \int_{\mathbb{R}^D} \mathfrak{q}_T(y | x) \nu_0^*(x) dx \\ &= \nu_T^*(y) \int_{\mathbb{R}^D} \frac{\mathfrak{q}_T(y | x)}{\mathcal{T}_T[\nu_T^*](x)} \mathfrak{p}_0(x) dx. \end{aligned}$$

Similarly, the marginal density of \tilde{X}_T equals to

$$\tilde{\mathfrak{p}}_T(y) = \tilde{\nu}_T(y) \int_{\mathbb{R}^D} \frac{\mathfrak{q}_T(y | x)}{\mathcal{T}_T[\tilde{\nu}_T](x)} \mathfrak{p}_0(x) dx.$$

Hence, the performance of $\tilde{\nu}_T$ is measured with

$$\begin{aligned} \text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) &= \int_{\mathbb{R}^D} \log \left(\frac{\nu_T^*(y)}{\tilde{\nu}_T(y)} \right) \mathfrak{p}_T(y) dy \\ &+ \int_{\mathbb{R}^D} \log \left(\int_{\mathbb{R}^D} \frac{\mathfrak{q}_T(y|x)}{\mathcal{T}_T[\nu_T^*](x)} \mathfrak{p}_0(x) dx \right) \\ &- \int_{\mathbb{R}^D} \log \left(\int_{\mathbb{R}^D} \frac{\mathfrak{q}_T(y|x)}{\mathcal{T}_T[\tilde{\nu}_T](x)} \mathfrak{p}_0(x) dx \right). \end{aligned}$$

In the next section, we address the following question.

Q1: Given a precision parameter $\varepsilon > 0$, how large should be the parameters L, W, S, B of an approximating neural network to ensure that $\text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) \leq \varepsilon$?

Besides approximation guarantees, we are interested in non-asymptotic rates of convergence for Schrödinger potential estimation. In (Puchkin et al., 2025), the authors studied theoretical properties an empirical risk minimizer

$$\begin{aligned} \log \hat{\nu}_T \in \operatorname{argmin}_{\psi \in \Psi} \left\{ -\frac{1}{n} \sum_{i=1}^n \psi(Y_i) \right. \\ \left. - \frac{1}{n} \sum_{i=1}^n \log \left(\int_{\mathbb{R}^D} \frac{\mathfrak{q}_T(Y_i|x)}{\mathcal{T}_T[e^\psi](x)} \mathfrak{p}_0(x) dx \right) \right\}, \quad (7) \end{aligned}$$

where $Y_1, \dots, Y_n \sim \mathfrak{p}_T$ are i.i.d. samples and Ψ is a reference class of potentials. In particular, they proved an oracle inequality for the excess risk of $\hat{\nu}_T$ and specified the estimation error. However, they left the question of approximation out of the scope of their paper. This brings us to the following question.

Q2: Given a training sample of size n , what are the rates of convergence of the estimate (7)?

3. Main results

In this section, we present our main theoretical results on approximation and estimation of the right Schrödinger log-potential. To avoid the curse of dimensionality, we introduce the following assumption on the data distribution.

without suffering from the curse of dimensionality. We begin with introducing the key assumptions on the , which exhibits an underlying low-dimensional structure.

Assumption 3.1. The data distribution $Y \sim \mathfrak{p}_T$ is generated by the following model:

$$Y = GZ + \sigma_{\text{data}}\xi,$$

where $\sigma_{\text{data}} \in (0, 1)$ and $G \in \mathbb{R}^{D \times d}$ has orthonormal columns $G^\top G = I_d$. The random elements $Z \sim \mu$ and $\xi \sim \mathcal{N}(0, I_D)$ are independent, with $Z \in \mathbb{R}^d$ being centered and sub-Gaussian with $\|Z\|_{\psi_2} \leq \tau < +\infty$.

Assumption 3.1 means that the data distribution is concentrated near d -dimensional unknown linear subspace, where the *effective dimension* d can be significantly smaller than the ambient dimension D . Similar models were considered in recent papers on score estimation and denoising diffusion models (Sasaki et al., 2016; Chen et al., 2023; Oko et al., 2023). The parameter $\sigma_{\text{data}} > 0$ measures intrinsic data variability. In practice, it is unlikely that data points lie exactly on a low-dimensional subspace or a manifold but they may concentrate around a low-dimensional set. In many practical settings access to noise-free measurements is either expensive or impossible, see, for instance, (Daras et al., 2024). Our model therefore includes an irreducible noise and still preserves the intrinsic low-dimensional structure of the data. In addition, the requirement $\sigma_{\text{data}} > 0$ is crucial for the stochastic optimal control approach to generative modelling, because, according to (Dai Pra, 1991), the marginal endpoint distribution of the base process $\{X_t^0 : 0 \leq t \leq T\}$ must dominate the target density \mathfrak{p}_T .

In view of Assumption 3.1, we have that the target density admits a decomposition

$$\begin{aligned} \log \mathfrak{p}_T(y) &= \log \mathfrak{p}_{||}(G^\top y) - \frac{D-d}{2} \log(2\pi\sigma_{\text{data}}^2) \\ &- \frac{\|(I_D - GG^\top)y\|^2}{2\sigma_{\text{data}}^2}, \quad (8) \end{aligned}$$

where the on-support density is given by

$$\begin{aligned} \mathfrak{p}_{||}(u) &= (2\pi\sigma_{\text{data}}^2)^{-d/2} \\ &\cdot \int_{\mathbb{R}^d} \exp \left\{ -\frac{\|u - z\|^2}{\sigma_{\text{data}}^2} \right\} d\mu(z), \quad u \in \mathbb{R}^d. \quad (9) \end{aligned}$$

It is important to note that we impose mild conditions on the mixing measure μ . For instance, (Chen et al., 2023) requires that $\mathfrak{p}_{||}(u)$ must be globally Lipschitz. This yields that the log-density Hessian $\nabla^2 \log \mathfrak{p}_{||}(u)$ must have a uniformly bounded operator norm. Direct computations show that

$$\begin{aligned} \nabla^2 \log \mathfrak{p}_{||}(u) &= \frac{1}{2Z^2(u)} \int \int_{\mathbb{R}^d \mathbb{R}^d} \left[(z_1 - z_2)(z_1 - z_2)^\top \right. \\ &\cdot \left. \prod_{j=1}^2 \exp \left\{ \frac{u^\top z_j}{\sigma_{\text{data}}^2} - \frac{\|z_j\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z_1) d\mu(z_2) \right], \quad (10) \end{aligned}$$

where

$$Z(u) = \int \int_{\mathbb{R}^d \mathbb{R}^d} \prod_{j=1}^2 \exp \left\{ \frac{u^\top z_j}{\sigma_{\text{data}}^2} - \frac{\|z_j\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z_1) d\mu(z_2),$$

While the right-hand side of (10) is obviously bounded on \mathbb{R}^d in the case when μ is Gaussian or has a bounded support, it is not clear, whether this claim holds for a general sub-Gaussian measure μ . For this reason, Assumption 3.1 allows us to consider a more general framework compared to (Chen et al., 2023).

We also impose a mild assumption on the potential ν_T^* assuming that it is bounded from above.

Assumption 3.2. *There exists a positive constant M such that $\log \nu_T^*(y) \leq M$ for all $y \in \mathbb{R}^D$. In addition, it holds that $\mathcal{T}_\infty[\log \nu_T^*] = 0$.*

Note that Assumption 3.2 is consistent with the analysis of Schrödinger potential estimation in (Puchkin et al., 2025) and (Puchkin et al., 2026). The requirement $\mathcal{T}_\infty[\log \nu_T^*] = 0$ is not restrictive at all. Let us remind a reader that the potentials ν_0^* and ν_T^* are defined up to multiplicative constant. Hence, the condition $\mathcal{T}_\infty[\log \nu_T^*] = 0$ should be considered as a normalization.

We now establish that, under the assumptions formulated above, the right Schrödinger log-potential can be approximated without the curse of dimensionality. The key observation is that the log-potential admits the following decomposition.

Lemma 3.3 (log-potential decomposition). *Assume p_T satisfies Assumption 3.1, and let p_0 be the standard Gaussian density on \mathbb{R}^D . Let $\rho_{||}(u)$ and $q_{||}(u|x)$ be the densities of the Gaussian distributions $\mathcal{N}(0, I_d)$ and $\mathcal{N}(e^{-bT}x, \gamma\sigma_T^2 I_d)$, where σ_T^2 is defined in (3), respectively. For any measurable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, introduce*

$$\mathcal{T}_{||}[f](x) = \int_{\mathbb{R}^d} f(u) q_{||}(u|x) du, \quad (11)$$

and let $\nu_{||}^*: \mathbb{R}^d \rightarrow \mathbb{R}$ be a solution of the equation

$$\begin{aligned} \log \nu_{||}^*(u) &= \log p_{||}(u) \\ &- \log \int_{\mathbb{R}^d} \frac{q_{||}(u|x) \rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}^*](x)}, \quad u \in \mathbb{R}^d, \end{aligned} \quad (12)$$

such that $\int_{\mathbb{R}^d} \log \nu_{||}^*(u) e^{-\|u\|^2/(2\gamma\sigma_\infty^2)} du = 0$. Then, the log-potential $\log \nu_T^*$ admits a decomposition

$$\begin{aligned} \log \nu_T^*(y) &= \log \nu_{||}^*(G^\top y) \\ &- \frac{\|(I_D - GG^\top)y\|^2 - \gamma\sigma_\infty^2(D-d)}{2} \\ &\cdot \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma\sigma_T^2} \right), \end{aligned}$$

where $\sigma_\infty^2 = (2b)^{-1}$ and σ is the positive root of the quadratic equation $\gamma\sigma_T^2\sigma = e^{-bT}(\sigma_{\text{data}}^2 - \sigma^2)$. Furthermore, $\mathcal{T}_\infty[\log \nu_T^*] = 0$.

The proof of Lemma 3.3 is postponed to Appendix A.1. We would like to note that $\nu_{||}^*$ is well-defined, since it is the right Schrödinger potential for the marginals $\rho_{||}(x)$, $p_{||}(u)$ and the Markov kernel $q_{||}(u|x)$. Importantly, Lemma 3.3 shows that the approximation of the right Schrödinger log-potential $\log \nu_T^*$ can be achieved by approximating $\log \nu_{||}^*$, a d -dimensional function, thus avoiding the curse of dimensionality associated with the ambient D -dimensional space. Nevertheless, existing approximation results from the literature do not apply to $\log \nu_{||}^*$. Overcoming this obstacle is the main technical challenge of the present paper. It is also worth mentioning that the restriction of $\log \nu_T^*(y)$ onto the orthogonal complement of $\text{Im}(G)$ is bounded from above if and only if $\gamma\sigma_T^2 \geq \sigma_{\text{data}}^2$. Thus, the inequality $\gamma\sigma_T^2 \geq \sigma_{\text{data}}^2$ is a necessary condition for Assumption 3.2 to hold. It poses no problems, since a learner is free to choose the parameters γ , b , and T .

Theorem 3.4. *Grant Assumptions 3.1 and 3.2. Assume further that bT is sufficiently large in the sense that it satisfies*

$$bT \gtrsim \log \log(1/\varepsilon) + \log(D\sigma_{\text{data}}^{-2}) + \log \log \sigma_{\text{data}}^{-2}. \quad (13)$$

Then, for every $0 < \varepsilon < A \wedge M_{||} \wedge 1$, where $A \asymp \sigma_{\text{data}}^{-2}$ and $M_{||} \asymp M + D\sigma_{\text{data}}^{-2}$, there exists an approximation $\tilde{\nu}_T(y)$ of the right Schrödinger potential ν_T^* of the form

$$\begin{aligned} \log \tilde{\nu}_T(y) &= \log \tilde{\nu}_{||}(G^\top y) \\ &- \frac{\|(I_D - GG^\top)y\|^2 - \gamma\sigma_\infty^2(D-d)}{2} \\ &\cdot \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma\sigma_T^2} \right), \end{aligned}$$

such that $\log \tilde{\nu}_{||} \in \text{NN}(L, W, S, B)$ fulfils $-A(1 + \|z\|^2) \leq \log \tilde{\nu}_{||}^*(z) \leq M_{||}$ for all $z \in \mathbb{R}^d$ and the terminal density \tilde{p}_T , corresponding to $\tilde{\nu}_T$, satisfies $\text{KL}(p_T, \tilde{p}_T) \leq \varepsilon$. Here $\sigma_\infty^2 = (2b)^{-1}$ and σ is the positive root of the quadratic equation $\gamma\sigma_T^2\sigma = e^{-bT}(\sigma_{\text{data}}^2 - \sigma^2)$. Furthermore, the neural network architecture satisfies

$$\begin{aligned} L &\lesssim \log^3(1/\varepsilon) (\log \sigma_{\text{data}}^{-2} + \log \log(D\sigma_{\text{data}}^{-2})), \\ B &\lesssim \sqrt{\log(\varepsilon^{-1}(D\sigma_{\text{data}}^{-2}))}, \\ \|W\|_\infty \vee S &\lesssim \frac{(\log(1/\varepsilon))^{3d+6}}{\sigma_{\text{data}}^d} \left(\log \frac{D}{\sigma_{\text{data}}^2} \right)^{d/2+2}. \end{aligned}$$

The hidden constants depend solely on the parameters γ , b , M , and d .

The proof of Theorem 3.4 is postponed to Appendix A.2. Theorem 3.4 establishes that the right Schrödinger log-potential $\log \nu_T^*$ can be approximated within the accuracy ε in terms of the KL divergence via a quite small neural network with only $\tilde{\mathcal{O}}(\log(1/\varepsilon))$ weights. It is a consequence of the favourable bounds on the higher-order derivatives of

the log-potential (see Lemmata 5.1 and 5.2). To the best of our knowledge, Theorem 3.4 provides the first approximation result of such kind. A very recent paper (Belomestny et al., 2026) considered a related problem but focused on the approximation of \mathfrak{p}_T/ν_T^* in the $L^2(\mathfrak{p}_T)$ norm. Our goal is different. We also would like to note that the condition (13) should be considered as a condition on the product bT , rather than on the time horizon T . We still can take $T = 1$ but it will require larger values of b .

We proceed with theoretical guarantees on accuracy of estimation of the potential ν_T^* from n i.i.d. samples Y_1, \dots, Y_n drawn according to \mathfrak{p}_T . Lemma 3.3 and Theorem 3.4 suggest us to use a class of neural networks of the following form.

Definition 3.5. We denote by $\mathcal{F}_{\text{NN}}(L, W, S, B, \sigma_{\min})$ the class of neural networks of the form

$$\psi(y) = \varphi(V^\top y) - a (\|(I_D - VV^\top)y\|^2 - \gamma\sigma_\infty^2(D - d)),$$

where $\varphi \in \text{NN}(L, W, S, B)$, the matrix $V \in \mathbb{R}^{D \times d}$ has orthonormal columns, $0 < \sigma_{\min} \leq \sigma < 1$ for some $\sigma_{\min} \leq \sigma_{\text{data}}$, and $|a| \lesssim \sigma_{\min}^{-2} \vee (\gamma\sigma_T^2)^{-1}$. We further assume that

$$-A(1 + \|z\|^2) \leq \varphi(z) \leq M_{\parallel}, \quad \text{for all } z \in \mathbb{R}^d,$$

where $A \lesssim \sigma_{\min}^{-2}$ and $M_{\parallel} \lesssim M + D\sigma_{\min}^{-2}$. Moreover, we assume that $\mathcal{T}_\infty[\varphi] = 0$.

In Definition 3.5, since σ_{data} is unknown, we make it learnable and use the decomposition given by Lemma 3.3. Crucially, we assume that we have an access to a lower bound to σ_{data} , which we denote by $\sigma_{\min} > 0$. Moreover, since the linear subspace given by G from Assumption 3.1 is also unknown, we introduce a learnable matrix V with orthonormal columns. Finally, we emphasize that the log-potential ψ satisfies the normalization constraint $\mathcal{T}_\infty[\psi] = 0$. We now formulate our main result on the estimation of the Schrödinger potential.

Theorem 3.6. Assume conditions of Theorem 3.4 hold. Let $\log \hat{\nu}_T$ be an empirical risk minimizer (7) over the class $\Psi = \mathcal{F}_{\text{NN}}$ (see Definition 3.5) with the parameters

$$\begin{aligned} L &\lesssim \log^3 n (\log \sigma_{\min}^{-2} + \log \log(D\sigma_{\min}^{-2})), \\ B &\lesssim \sqrt{\log(n(D\sigma_{\min}^{-2}))}, \\ \|W\|_\infty \vee S &\lesssim (\log n)^{3d+6} \sigma_{\min}^{-d} (\log(D\sigma_{\min}^{-2}))^{d/2+2}. \end{aligned}$$

Assume further that $bT \gtrsim \log \log n + \log(D\sigma_{\min}^{-2})$. Let $\hat{\mathfrak{p}}_T$ be the density associated with the estimate $\hat{\nu}_T$. Then, for every $\delta \in (0, 1/2)$, with probability at least $(1 - 2\delta)$, it follows that

$$\text{KL}(\mathfrak{p}_T, \hat{\mathfrak{p}}_T) \lesssim \frac{D^2 \mathcal{L}(n, D, \sigma_{\min}) \log(1/\delta)}{n\sigma_{\min}^d},$$

where

$$\mathcal{L}(n, D, \sigma_{\min}) = (\log n)^{3d+16} (\log(D\sigma_{\min}^{-2}))^{d/2+5}.$$

Here, the hidden constants depend on γ, b, M , and d only.

The proof of Theorem 3.6 is moved to Appendix A.3. Let us elaborate on comparison of Theorem 3.6 with relevant results in the existing literature. By a straightforward application of the Girsanov theorem, Theorem 2 in (Chen et al., 2023) yields that

$$\begin{aligned} \mathbb{E}\text{KL}(\mathfrak{p}_{T-t_0}, \hat{\mathfrak{p}}_{T-t_0}) \\ = \tilde{\mathcal{O}}\left(\frac{1}{t_0} \left(n^{-\frac{2}{d+5}} + Dn^{-\frac{d+3}{d+5}}\right)\right), \end{aligned} \quad (14)$$

where t_0 is a stopping time. By contrast, Theorem 3.6 significantly improves the dependence on the sample size. The underlying reason is the approximation theory developed in Theorem 3.4. Our construction guarantees that the number of non-zero parameters scales polylogarithmically with ε^{-1} , thereby enabling the faster convergence rate in the sample size. While the results of (Yakovlev and Puchkin, 2025) extend beyond linear subspace assumption of (Chen et al., 2023), their bound with the smoothness parameter $\beta \asymp \log n$ does not recover the $\tilde{\mathcal{O}}(n^{-1})$ convergence rate. This is due to the fact that the dependence on the ambient dimension D is worse than in Theorem 3.6. Moreover, the bounds (14) and those of (Yakovlev and Puchkin, 2025) explode as $t_0 \rightarrow 0$. The same issue appears in the analysis of diffusion and flow-based generative models (Oko et al., 2023; Yakovlev and Puchkin, 2025; Fukumizu et al., 2025), even when the target distribution is absolutely continuous. On the other hand, the bound established in Theorem 3.6 circumvents this limitation and eliminates the need for stopping time t_0 . The dependence on the noise level appears only through σ_{\min}^2 , which plays a role analogous to t_0 .

Notably, (Pooladian and Niles-Weed, 2025) obtained that expected squared total variation distance between the path measures, corresponding to $\{X_t^* : 0 \leq t \leq T\}$ and the process with the estimated drift $\{\hat{X}_t : 0 \leq t \leq T\}$, scales as $\mathcal{O}(n^{-1/2} + n^{-1}(1 - \tau)^{-k-2})$, where $1 - \tau$ denotes the stopping time and k represents the dimension of the manifold supporting the target distribution. Our Theorem 3.6 correctly captures the dependence on σ_{\min} through the effective dimension d . Notably, unlike (Pooladian and Niles-Weed, 2025), our result avoids the slow $n^{-1/2}$ convergence rate, achieving improved statistical efficiency.

The above findings indicate possible gaps in theoretical guarantees for generative diffusion models and flow-based models. We suppose that these problems are primarily related to the properties of the corresponding loss functions (in denoising score matching and flow matching), which

blow up as t_0 approaches zero, rather than to the properties of the score function and the velocity field, which may remain bounded. We emphasize that this comparison does not imply that the approaches, relying on construction of Schrödinger bridges or stochastic optimal control, inherently outperform diffusion- and flow-based methods. Our analysis focuses solely on upper bounds, while lower bounds remain an open question in this context.

4. Extensions and future directions

In view of (8), the target density p_T can be represented as a product of two probability densities on orthogonal spaces. Such a model is quite popular in statistics in the context of non-Gaussian component analysis and linear dimension reduction (see, for instance, (Blanchard et al., 2006) and the references therein). It has a natural extension when p_T is represented as a product of $J \geq 2$ components:

$$p_T(y) = \prod_{j=1}^J p_T^{(j)}(S_j y), \quad y \in \mathbb{R}^D, \quad (15)$$

where the matrices $S_j \in \mathbb{R}^{d_j \times D}$, $j \in \{1, \dots, J\}$, $d_1 + \dots + d_J = D$, are such that $S_j S_j^\top = I_{d_j}$, for every $1 \leq j \leq J$, and $S_j S_k^\top = 0$ for all $j \neq k$, and, for each $j \in \{1, \dots, J\}$, $p_T^{(j)}$ is a probability density on \mathbb{R}^{d_j} . In this case, it is also possible to approximate and estimate the Schrödinger potential ν_T^* , because it is also factorizable. We provide the corresponding result below.

Lemma 4.1. *Assume that p_T admits a decomposition (15), and let p_0 be the density of $\mathcal{N}(0, I_D)$. For each $j \in \{1, \dots, J\}$, let $\rho^{(j)}(\tilde{y}_j)$ and $q_T^{(j)}(\tilde{y}_j | \tilde{x}_j)$ denote the densities of the Gaussian distributions $\mathcal{N}(0, I_{d_j})$ and $\mathcal{N}(e^{-bT} \tilde{x}_j, \gamma \sigma_T^2 I_{d_j})$, respectively. Let $(\nu_0^{(j)}, \nu_T^{(j)})$ be the Schrödinger potentials of the corresponding sub-problem on \mathbb{R}^{d_j} with initial density $\rho^{(j)}$, target density $p_T^{(j)}$, and transition kernel $q_T^{(j)}$. Then the Schrödinger log-potentials admit the decompositions*

$$\begin{aligned} \log \nu_0^*(x) &= \sum_{j=1}^J \log \nu_0^{(j)}(S_j x), \quad x \in \mathbb{R}^D, \\ \log \nu_T^*(y) &= \sum_{j=1}^J \log \nu_T^{(j)}(S_j y), \quad y \in \mathbb{R}^D. \end{aligned}$$

The proof of Lemma 4.1 is postponed to Appendix B. Further extensions of Assumption 3.1 can include situation when the support of the mixing measure μ is a union of affine spaces (Chakraborty and Nguyen, 2026) or nonlinear (Yakovlev and Puchkin, 2025). The nonlinear model with low intrinsic dimension has a practical impact, because real-world data sets often exhibit nonlinear underlying structures (Brown et al., 2023).

5. Proof Insights

In this section, we provide some intuition behind the proofs of our main results, Theorems 3.4 and 3.6. Lemma 3.3 essentially reduces the problem to approximation and estimation of the on-support log-potential $\log \nu_{||}^*$ defined in (12). We rely on a well-known result of (Schmidt-Hieber, 2020) (Lemma A.5) to approximate the summands in the right-hand side of (12) via feedforward neural networks with ReLU activations. However, the most challenging and the most technical part of the proof is to provide sharp bounds on $\|\cdot\|_\infty$ -norms of higher-order derivatives of

$$\log p_{||}(u) \quad \text{and} \quad \log \int_{\mathbb{R}^d} \frac{q_{||}(u|x)\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}^*](x)}.$$

We provide the corresponding results below and believe that they could be of independent interest.

Lemma 5.1. *Assume that the target distribution p_T satisfies Assumption 3.1, and let $p_{||}$ be the on-support density defined in (9). Then, for every integer $k \geq 2$ and every $u \in \mathbb{R}^d$,*

$$\begin{aligned} &\left\| \nabla^k \left(\log p_{||}(u) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2} \right) \right\|_\infty \\ &\leq \left(\frac{36(\|u\| + 2\tau\sqrt{d} + \sigma_{\text{data}})\sqrt{C}}{\sigma_{\text{data}}^2} \right)^k k!, \end{aligned}$$

where $C \geq 1$ is an absolute constant.

Lemma 5.2. *Let $\rho_{||}$ and $q_{||}(\cdot|x)$ be the densities of $\mathcal{N}(0, I_d)$ and $\mathcal{N}(e^{-bT}x, \gamma\sigma_T^2 I_d)$, respectively. Let $\nu_{||} : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy $\mathcal{T}_\infty[\log \nu_{||}] = 0$ and $\log \nu_{||} \leq M_{||}$, and assume that bT obeys*

$$8e^{2-bT} \leq \gamma\sigma_T^2\sqrt{d} \quad \text{and} \quad 14e^2\sqrt{d}e^{-bT/2} \leq \log 2. \quad (16)$$

Then, for every integer $k \geq 3$ and every $y \in \mathbb{R}^d$,

$$\begin{aligned} &\left\| \nabla^k \log \int_{\mathbb{R}^d} \frac{q_{||}(y|x)\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right\|_\infty \\ &\leq \left(\frac{3 \cdot 2^{2+d/2} e^{-bT} \sqrt{C}}{\gamma\sigma_T^2} \right)^k \\ &\quad \cdot \exp \left\{ k(2M_{||} + 3) + \frac{k e^{-2bT} \|y\|^2}{2(\gamma\sigma_T^2)^2} \right\} k!, \end{aligned}$$

where $\mathcal{T}_{||}$ is defined in (11) and $C \geq 1$ is an absolute constant.

The proofs of Lemmata 5.1 and 5.2 are moved to Appendices C.1 and C.2, respectively. The main obstacle is that the density $\rho_{||}(x)$ as well as the mixing measure μ have

unbounded supports. This makes the proof significantly more complicated compared to, for instance, (Yakovlev and Puchkin, 2025) (Lemma 4.1). In the proofs of Lemmata 5.1 and 5.2, we derive explicit expressions for the k -th order derivatives and reduce the problem to upper bounds on the expectations of homogeneous shift-invariant polynomials of sub-Gaussian random variables of degree k (see Lemma D.2). The proof of Lemma 5.1 significantly exploits the following fact from probability theory, which, to our knowledge, was not presented in the literature before.

Lemma 5.3. *Let X and Y be independent random vectors such that $X \sim \mathcal{N}(0, \sigma^2 I_d)$, $\sigma > 0$, and $\mathbb{E}Y = 0$ and $\mathbb{E}\|Y\|^2 \leq d\tau^2$. Set $S = X + Y$. Then, for any $s \in \mathbb{R}^d$, the conditional distribution of X given $S = s$ is sub-Gaussian. Moreover, its conditional ψ_2 -norm and the conditional expectation $\mu(s) := \mathbb{E}[X \mid S = s]$ satisfy the inequalities*

$$\|X - \mu(s)\|_{\psi_2(\cdot|S=s)} \leq 6R + 6\sigma$$

and

$$\|\mu(s)\| \leq 3R + 3\sigma,$$

where $R = \|s\| + 2\tau\sqrt{d}$.

We provide the proof of Lemma 5.3 in Appendix 5.3. Notably, the random vector Y in Lemma 5.3 is not necessarily sub-Gaussian, it should just have a finite second moment.

Lemmata 5.1, 5.2, and 5.3 are the main ingredients of the proof of Theorem 3.4. Once we quantified the complexity of Schrödinger potential approximation, the proof of Theorem 3.6 is not very hard. It just requires to carefully check that all the conditions of Theorem 1 from (Puchkin et al., 2025) are satisfied and combine this results with Theorem 3.4.

References

- I. Azangulov, G. Deligiannidis, and J. Rousseau. Convergence of diffusion models under the manifold hypothesis in high-dimensions. Preprint. ArXiv:2409.18804, 2024.
- D. Belomestny, A. Naumov, N. Puchkin, and D. Suchkov. Schrödinger bridge problem via empirical risk minimization. Preprint. ArXiv:2602.08374, 2026.
- G. Blanchard, M. Kawanabe, M. Sugiyama, V. Spokoiny, and K.-R. Müller. In search of non-gaussian components of a high-dimensional distribution. *Journal of Machine Learning Research*, 7(9):247–282, 2006.
- B. C. Brown, A. L. Caterini, B. L. Ross, J. C. Cresswell, and G. Loaiza-Ganem. Verifying the union of manifolds hypothesis for image data. In *The Eleventh International Conference on Learning Representations*, 2023.
- S. Chakraborty and X. Nguyen. Learning mixtures of nonparametric and convolutional measures on effectively low-dimensional affine spaces. Preprint. ArXiv:2604.17236, 2026.
- M. Chen, K. Huang, T. Zhao, and M. Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 4672–4712. PMLR, 2023.
- T. Chen, G.-H. Liu, and E. Theodorou. Likelihood training of schrödinger bridge using forward-backward SDEs theory. In *International Conference on Learning Representations*, 2022.
- A. Chiarini, G. Conforti, G. Greco, and L. Tamanini. A semiconcavity approach to stability of entropic plans and exponential convergence of Sinkhorn’s algorithm. Preprint. ArXiv:2412.09235, 2024.
- L. Chizat, A. Delalande, and T. Vaškevičius. Sharper exponential convergence rates for Sinkhorn’s algorithm in continuous settings. Preprint. ArXiv:2407.01202, 2024.
- G. Conforti, A. Durmus, and G. Greco. Quantitative contraction rates for Sinkhorn’s algorithm: beyond bounded costs and compact marginals. Preprint. ArXiv:2304.04451, 2023.
- P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. *Applied mathematics and Optimization*, 23(1):313–329, 1991.
- G. Daras, A. Dimakis, and C. C. Daskalakis. Consistent diffusion meets Tweedie: Training exact ambient diffusion models with noisy data. In *Forty-first International Conference on Machine Learning*, 2024.
- V. De Bortoli, J. Thornton, J. Heng, and A. Doucet. Diffusion Schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- C. Domingo-Enrich, M. Drozdal, B. Karrer, and R. T. Q. Chen. Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control. In *The Thirteenth International Conference on Learning Representations*, 2025.
- K. Fukumizu, T. Suzuki, N. Isobe, K. Oko, and M. Koyama. Flow matching achieves almost minimax optimal convergence. In *The Thirteenth International Conference on Learning Representations*, 2025.

- 440 N. Gushchin, S. Kholkin, E. Burnaev, and A. Korotin. Light and optimal Schrödinger bridge matching. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 17100–17122. PMLR, 2024.
- 441
442
443
444
445 J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- 446
447
448
449
450 A. Korotin, N. Gushchin, and E. Burnaev. Light Schrödinger bridge. In *The Twelfth International Conference on Learning Representations*, 2024.
- 451
452
453
454 C. Léonard. A survey of the Schrödinger problem and some of its connections with optimal transport. Preprint. ArXiv:1308.0215, 2013.
- 455
456
457
458 Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- 459
460
461
462 X. Liu, C. Gong, and Q. Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023.
- 463
464
465
466
467 R. Nakada and M. Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020.
- 468
469
470
471
472 K. Oko, S. Akiyama, and T. Suzuki. Diffusion models are minimax optimal distribution estimators. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 26517–26582. PMLR, 2023.
- 473
474
475
476
477 S. Peluchetti. Diffusion bridge mixture transports, Schrödinger bridge problems and generative modeling. *Journal of Machine Learning Research*, 24(374):1–51, 2023.
- 478
479
480
481
482 A.-A. Pooladian and J. Niles-Weed. Plug-in estimation of Schrödinger bridges. *SIAM Journal on Mathematics of Data Science*, 7(3):1315–1336, 2025.
- 483
484
485
486 N. Puchkin, I. Pustovalov, Y. Sapronov, D. Suchkov, A. Naumov, and D. Belomestny. Sample complexity of Schrödinger potential estimation. Preprint. ArXiv:2506.03043, 2025.
- 487
488
489
490
491 N. Puchkin, D. Suchkov, A. Naumov, and D. Belomestny. Tight bounds for Schrodinger potential estimation in unpaired data translation. In *The Fourteenth International Conference on Learning Representations*, 2026.
- 492
493
494
495 G. Rapakoulias, A. R. Pedram, F. Liu, L. Zhu, and P. Tsionas. Go with the flow: Fast diffusion for gaussian mixture models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- 496
497 P. Rigollet and J.-C. Hütter. High-dimensional statistics, 2023.
- 498
499
500 H. Sasaki, G. Niu, and M. Sugiyama. Non-gaussian component analysis with log-density gradient estimation. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 1177–1185. PMLR, 2016.
- 501
502
503
504 J. Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- 505
506
507
508 Y. Shi, V. De Bortoli, A. Campbell, and A. Doucet. Diffusion Schrödinger bridge matching. In *Advances in Neural Information Processing Systems*, volume 36, pages 62183–62223. Curran Associates, Inc., 2023.
- 509
510
511
512 Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- 513
514
515
516 R. Tang and Y. Yang. Adaptivity of diffusion models to manifold structures. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 1648–1656. PMLR, 2024.
- 517
518
519
520 F. Vargas, P. Thodoroff, A. Lamacraft, and N. Lawrence. Solving Schrödinger bridges via maximum likelihood. *Entropy*, 23(9), 2021.
- 521
522
523
524 R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- 525
526
527
528 K. Yakovlev and N. Puchkin. Generalization error bound for denoising score matching under relaxed manifold assumption. In *Proceedings of Thirty Eighth Conference on Learning Theory*, volume 291 of *Proceedings of Machine Learning Research*, pages 5824–5891. PMLR, 2025.
- 529
530
531
532 K. Yakovlev, A. Markovich, and N. Puchkin. Implicit score matching meets denoising score matching: improved rates of convergence and log-density Hessian estimation. Preprint. ArXiv:2512.24378, 2025.

495	Contents	
496		
497	1 Introduction	1
498		
499	2 Problem setup	3
500		
501	3 Main results	4
502		
503	4 Extensions and future directions	7
504		
505	5 Proof Insights	7
506		
507		
508		
509	A Proofs of the results from Section 3	11
510	A.1 Proof of Lemma 3.3	11
511	A.2 Proof of Theorem 3.4	12
512	A.3 Proof of Theorem 3.6	14
513	A.4 Proof of Lemma A.1	16
514	A.5 Proof of Lemma A.2	17
515		
516		
517		
518	B Proof of Lemma 4.1	20
519		
520		
521	C Proofs of the results from Section 5	21
522	C.1 Proof of Lemma 5.1	21
523	C.2 Proof of Lemma 5.2	23
524	C.3 Proof of Lemma 5.3	29
525		
526		
527		
528	D Auxiliary results	31
529	D.1 Proof of Lemma D.1	31
530	D.2 Proof of Lemma D.2	33
531		
532		
533		
534		
535		
536		
537		
538		
539		
540		
541		
542		
543		
544		
545		
546		
547		
548		
549		

A. Proofs of the results from Section 3

A.1. Proof of Lemma 3.3

First, observe that the density of X_T satisfies

$$p_T(y) + \frac{D}{2} \log(2\pi\sigma_{\text{data}}^2) = -\frac{\|(I_D - GG^\top)y\|^2}{2\sigma_{\text{data}}^2} + \log \int_{\mathbb{R}^d} \exp\left(-\frac{\|G^\top y - u\|^2}{2\sigma_{\text{data}}^2}\right) d\mu(u).$$

We now show that the Ornstein-Uhlenbeck operator also admits a similar decomposition. Let $G_\perp \in \mathbb{R}^{D \times (D-d)}$ be the matrix that completes G to an orthonormal basis of \mathbb{R}^D . Let us find the log-potential in the form $\log \nu_T^*(y) = \log \nu_{\parallel}^*(G^\top y) + \log \nu_{\perp}^*(G_\perp^\top y)$. Then, it holds that

$$\mathcal{T}_T[\nu_T^*](x) = \mathbb{E}_{\xi \sim \mathcal{N}(e^{-bT}x, \sigma_T^2 \gamma I_D)}[\nu_{\parallel}^*(G_{\parallel}\xi) \nu_{\perp}^*(G_{\perp}y)] = \mathbb{E}_{\xi}[\nu_{\parallel}^*(G_{\parallel}\xi)] \cdot \mathbb{E}_{\xi}[\nu_{\perp}^*(G_{\perp}\xi)],$$

where we used the fact that $G^\top \xi$ and $G_\perp^\top \xi$ are independent. Consequently, defining

$$\begin{aligned} \mathcal{T}_{\parallel}[\nu_{\parallel}](u) &= \mathbb{E}_{\xi \sim \mathcal{N}(e^{-bT}u, \sigma_T^2 \gamma I_d)}[\nu_{\parallel}(\xi)], & u \in \mathbb{R}^d, \\ \mathcal{T}_{\perp}[\nu_{\perp}](v) &= \mathbb{E}_{\xi \sim \mathcal{N}(e^{-bT}v, \sigma_T^2 \gamma I_{D-d})}[\nu_{\perp}(\xi)], & v \in \mathbb{R}^{D-d}, \end{aligned}$$

we conclude that

$$\mathcal{T}_T[\nu_T^*](x) = \mathcal{T}_{\parallel}[\nu_{\parallel}^*](G^\top x) \cdot \mathcal{T}_{\perp}[\nu_{\perp}^*](G_\perp^\top x).$$

We now define transition densities

$$\begin{aligned} \log q_{\parallel}(u' | u) &= -\frac{d}{2} \log(2\pi\gamma\sigma_T^2) - \frac{\|u' - e^{-bT}u\|^2}{2\gamma\sigma_T^2}, & u', u \in \mathbb{R}^d, \\ \log q_{\perp}(v' | v) &= -\frac{D-d}{2} \log(2\pi\gamma\sigma_T^2) - \frac{\|v' - e^{-bT}v\|^2}{2\gamma\sigma_T^2}, & v', v \in \mathbb{R}^{D-d}. \end{aligned}$$

Let us also define the standard Normal density in \mathbb{R}^d and \mathbb{R}^{D-d} as p_{\parallel} and p_{\perp} respectively. Thus, it follows that

$$\log \int_{\mathbb{R}^D} \frac{q_T(y|x)p_0(x) dx}{\mathcal{T}_T[\nu_T^*](x)} = \log \int_{\mathbb{R}^d} \frac{q_{\parallel}(G^\top y|x)p_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}^*](x)} + \log \int_{\mathbb{R}^{D-d}} \frac{q_{\perp}(G_\perp^\top y|x)p_{\perp}(x) dx}{\mathcal{T}_{\perp}[\nu_{\perp}^*](x)}.$$

Therefore, we conclude that the potentials ν_{\parallel}^* and ν_{\perp}^* can be found from the corresponding Schrödinger systems:

$$\begin{aligned} \log \nu_{\parallel}^*(z) &= -\frac{d}{2} \log(2\pi\sigma_{\text{data}}^2) + \log \int_{\mathbb{R}^d} \exp\left(-\frac{\|z - u\|^2}{2\sigma_{\text{data}}^2}\right) d\mu(u) - \log \int_{\mathbb{R}^d} \frac{q_{\parallel}(z|x)p_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}^*](x)}, \\ \log \nu_{\perp}^*(v) &= -\frac{D-d}{2} \log(2\pi\sigma_{\text{data}}^2) - \frac{\|v\|^2}{2\sigma_{\text{data}}^2} - \log \int_{\mathbb{R}^{D-d}} \frac{q_{\perp}(v|x)p_{\perp}(x) dx}{\mathcal{T}_{\perp}[\nu_{\perp}^*](x)}. \end{aligned}$$

Here, $z \in \mathbb{R}^d$ and $v \in \mathbb{R}^{D-d}$. Note that ν_{\perp}^* solves the Schrödinger problem between two Gaussian distributions $\mathcal{N}(0, I_{D-d})$ and $\mathcal{N}(0, \sigma_{\text{data}}^2 I_{D-d})$. This problem has a closed-form solution, which we formalize in the following lemma.

Lemma A.1. *Let the base process be given by (1) Let also $p_0 = \mathcal{N}(0, I_D)$ and $p_T = \mathcal{N}(\mu_T, \sigma_{\text{data}}^2 I_D)$ for some $\mu_T \in \mathbb{R}^D$ and $\sigma_{\text{data}} \in [0, 1)$. Then, the coupling between X_0^* and X_T^* in the optimally controlled process has the form*

$$\begin{pmatrix} X_0^* \\ X_T^* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ \mu_T \end{pmatrix}, \begin{pmatrix} I_D & \sigma I_D \\ \sigma I_D & \sigma_{\text{data}}^2 I_D \end{pmatrix} \right), \quad \frac{\sigma}{\sigma_{\text{data}}} = \frac{-1 + \sqrt{1 + 4C^2}}{2C}, \quad C = \frac{\sigma_{\text{data}} e^{-bT}}{\gamma \sigma_T^2}.$$

In addition, for all $x, y \in \mathbb{R}^D$, the Schrödinger potentials are given by

$$\log \nu_T^*(y) = -\frac{\|y\|^2}{2} \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma \sigma_T^2} \right) + \frac{y^\top \mu_T}{\sigma_{\text{data}}^2 - \sigma^2}$$

and

$$\begin{aligned} \log \nu_0^*(x) = & -\frac{\|x\|^2}{2} \left(\frac{\sigma_{\text{data}}^2 - e^{-bT} \sigma}{\sigma_{\text{data}}^2 - \sigma^2} \right) - \frac{x^\top \mu_T \sigma}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{\|\mu_T\|^2}{2(\sigma_{\text{data}}^2 - \sigma^2)} \\ & - \frac{D}{2} \log \left(\frac{2\pi(\sigma_{\text{data}}^2 - \sigma^2)}{\gamma\sigma_T^2} \right). \end{aligned}$$

The proof of Lemma A.1 is postponed to Appendix A.4. Therefore, applying Lemma A.1, we obtain the explicit form of $\log \nu_\perp^*$:

$$\log \nu_\perp^*(v) = -\frac{\|v\|^2 - \gamma\sigma_\infty^2(D-d)}{2} \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma\sigma_T^2} \right).$$

This log-potential satisfies the normalization condition $\mathcal{T}_\infty[\log \nu_\perp^*] = 0$. Finally, we observe that $\|G_\perp^\top y\|^2 = \|(I_D - GG^\top)y\|^2$, for all $y \in \mathbb{R}^D$. The proof is now complete. \square

A.2. Proof of Theorem 3.4

Lemma D.1 together with condition (13) implies that, for all $z \in \mathbb{R}^d$,

$$-A(1 + \|z\|^2) \leq \log \nu_\parallel^*(z) \leq M_\parallel, \quad A \lesssim \sigma_{\text{data}}^{-2}, \quad M_\parallel \lesssim M + D\sigma_{\text{data}}^{-2}. \quad (17)$$

Furthermore, Lemma D.1 claims that $\mathcal{T}_\infty[\log \nu_\parallel^*] = 0$. Next, we approximate the right log-potential on a compact. This is established by the following lemma.

Lemma A.2. *Grant the assumptions of Theorem 3.4. Fix an arbitrary $R \geq e$ and assume that*

$$bT \gtrsim M_\parallel + \log R, \quad (18)$$

where M_\parallel is given by (17). Then, for any $\varepsilon \in (0, 1)$, there exists a ReLU-neural network $\tilde{f} \in \text{NN}(L, W, S, B)$ such that

$$\begin{aligned} (i) \quad & \|\tilde{f} - \log \nu_\parallel^*\|_{L^\infty([-R, R]^d)} \leq \varepsilon, \\ (ii) \quad & \min_{u \in [-R, R]^d} \log \nu_\parallel^*(u) - \varepsilon \leq \tilde{f}(x) \leq M_\parallel + \varepsilon, \quad \text{for all } x \in \mathbb{R}^d. \end{aligned} \quad (19)$$

Furthermore, we have that

$$\begin{aligned} L & \lesssim \log^2(1/\varepsilon) \log(\sigma_{\text{data}}^{-2} R), \quad B \lesssim R, \\ \|W\|_\infty \vee S & \lesssim (\log(1/\varepsilon))^{5d/2+5} (R\sigma_{\text{data}}^{-1})^d \log(\sigma_{\text{data}}^{-2} R). \end{aligned}$$

The proof of Lemma A.2 is deferred to Appendix A.5. In what follows, we assume that $\varepsilon \leq M_\parallel \wedge A$, which implies that the approximation $\log \tilde{\nu}_\parallel$ from Lemma A.2 satisfies

$$-A \cdot d\|x\|^2 - 2A \leq \log \tilde{\nu}_\parallel(x) \leq 2M_\parallel, \quad \text{for all } x \in \mathbb{R}^d.$$

Next, observe that

$$\begin{aligned} \text{KL}(\mathbf{p}_T, \tilde{\mathbf{p}}_T) & \leq \int_{\mathbb{R}^D} \mathbf{p}_T(y) |\log \nu_T^*(y) - \log \tilde{\nu}_T(y)| \, dy \\ & + \int_{\mathbb{R}^D} \mathbf{p}_T(y) \left| \log \int_{\mathbb{R}^D} \frac{\mathbf{q}_T(y|x)\mathbf{p}_0(x) \, dx}{\mathcal{T}_\parallel[\nu_T^*](x)} - \log \int_{\mathbb{R}^D} \frac{\mathbf{q}_T(y|x)\mathbf{p}_0(x) \, dx}{\mathcal{T}_\parallel[\tilde{\nu}_T](x)} \right| \, dy. \end{aligned}$$

In view of Lemma 3.3, we have that

$$\begin{aligned} \text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) &\leq \int_{\mathbb{R}^d} \mathfrak{p}_{||}(z) \left| \log \nu_{||}^*(z) - \log \tilde{\nu}_{||}(z) \right| dz \\ &\quad + \int_{\mathbb{R}^d} \mathfrak{p}_{||}(z) \left| \log \int_{\mathbb{R}^d} \frac{\mathfrak{q}_{||}(z|u)\mathfrak{p}_{||}(u) du}{\mathcal{T}_{||}[\nu_{||}^*](u)} - \log \int_{\mathbb{R}^d} \frac{\mathfrak{q}_{||}(z|u)\mathfrak{p}_{||}(u) du}{\mathcal{T}_{||}[\tilde{\nu}_{||}](u)} \right| dz. \end{aligned}$$

To bound the second term in the above bound, we need to formulate the following helper lemma.

Lemma A.3 ((Puchkin et al., 2025), Lemma B.2, adapted). *Consider arbitrary functions $f_0, f_1 : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathcal{T}_\infty[f_0] = \mathcal{T}_\infty[f_1] = 0$. Assume that there exist constants $M \in \mathbb{R}$, $A \geq 0$, and $B \geq M \vee 0$ such that, for all $x \in \mathbb{R}^d$ and $i \in \{0, 1\}$,*

$$-A\|x\|^2 - B \leq f_i(x) \leq M.$$

Assume further that $bT \gtrsim 1$. Then, for all $y \in \mathbb{R}^d$, it follows that

$$\begin{aligned} &\left| \log \int_{\mathbb{R}^d} \frac{\mathfrak{q}_T(y|x)\mathfrak{p}_0(x) dx}{\mathcal{T}_T[e^{f_0}](x)} - \log \int_{\mathbb{R}^d} \frac{\mathfrak{q}_T(y|x)\mathfrak{p}_0(x) dx}{\mathcal{T}_T[e^{f_1}](x)} \right| \lesssim (\mathcal{T}_\infty[|f_1 - f_0|])^{1/\mathcal{K}(T)} \\ &\cdot (d^2 + \|y\|^2)^{1-1/\mathcal{K}(T)} \exp \left\{ \mathcal{O}(d + (M \log(A \vee B))\sqrt{d}e^{-bT} + e^{-bT}(1 + \|y\|^2)) \right\}. \end{aligned}$$

Here, the hidden constant behind $\mathcal{O}(\cdot)$ depends on γ and b only, and $\mathcal{K}(T)$ is given by (36).

Applying Lemma A.3 and using the bound (43), we obtain

$$\begin{aligned} \text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) &\lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^1(\mathfrak{p}_{||})} + \left(\mathcal{T}_\infty[|\log \nu_{||}^* - \log \tilde{\nu}_{||}|] \right)^{1/\mathcal{K}(T)} \\ &\cdot \exp \left\{ \mathcal{O}(d + (M_{||} \log(A \cdot d \vee B))\sqrt{d}e^{-bT}) \right\} \mathbb{E} \left[(d + \|\xi\|) \exp \{ \mathcal{O}(e^{-bT}(1 + \|\xi\|^2)) \} \right], \end{aligned}$$

where ξ has density $\mathfrak{p}_{||}$. Since the distribution of ξ coincides with the distribution of the convolution $\eta + \sigma_{\text{data}}\zeta$, where $\eta \sim \mu$ and $\zeta \sim \mathcal{N}(0, I_d)$ are independent, we deduce that $\|\|\xi\|^2\|_{\psi_1} \lesssim d$. Therefore, applying the Hölder inequality and (Vershynin, 2018, Proposition 2.7.1), we arrive at

$$\begin{aligned} \text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) &\lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^1(\mathfrak{p}_{||})} + \left(\mathcal{T}_\infty[|\log \nu_{||}^* - \log \tilde{\nu}_{||}|] \right)^{1/\mathcal{K}(T)} \\ &\cdot \exp \left\{ \mathcal{O}(d + (M_{||} \log(A \cdot d \vee B))\sqrt{d}e^{-bT}) \right\} \sqrt{d} \exp \{ \mathcal{O}(d^2 e^{-2bT}) \}. \end{aligned}$$

Thus, using condition (13), we deduce that

$$\text{KL}(\mathfrak{p}_T, \tilde{\mathfrak{p}}_T) \lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^1(\mathfrak{p}_{||})} + \left(\mathcal{T}_\infty[|\log \nu_{||}^* - \log \tilde{\nu}_{||}|] \right)^{1/\mathcal{K}(T)} \exp \{ \mathcal{O}(d) \}.$$

Using condition (17), we have that, for any $R > 0$,

$$\|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^1(\mathfrak{p}_{||})} \lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^\infty(\mathcal{B}(0, R))} + \sqrt{(M_{||} + A \cdot d + B)\mathbb{P}(\|\xi\| \geq R)}.$$

Applying the Gaussian concentration bound (Rigollet and Hütter, 2023, Theorem 1.19), we deduce that setting

$$R \asymp \sqrt{d \log(\varepsilon^{-1}(M_{||} + A \cdot d + B))}$$

yields

$$\|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^1(\mathfrak{p}_{||})} \lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^\infty(\mathcal{B}(0, R))} + \varepsilon \lesssim \varepsilon,$$

where the last inequality uses Lemma A.2. Similarly, we have that

$$\mathcal{T}_\infty[\|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|] \lesssim \|\log \nu_{||}^* - \log \tilde{\nu}_{||}\|_{L^\infty(\mathcal{B}(0,R))} + \varepsilon \lesssim \varepsilon.$$

In view of the bound (43) and condition (13), we have that

$$1 - 1/\mathcal{K}(T) \leq 1 - \exp(-14e^2\sqrt{d}e^{-bT}) \leq 14e^2\sqrt{d}e^{-bT} \lesssim \frac{1}{\log(1/\varepsilon)}.$$

Hence, for any $a \in (0, \varepsilon]$, we have $a^{1/\mathcal{K}(T)} \lesssim a$, and therefore

$$\text{KL}(\mathbf{p}_T, \tilde{\mathbf{p}}_T) \lesssim \varepsilon.$$

Finally, Lemma A.2 suggests that the architecture of $\log \tilde{\nu}_{||} \in \text{NN}(L, W, S, B)$ with

$$\begin{aligned} L &\lesssim \log^3(1/\varepsilon) (\log \sigma_{\text{data}}^{-2} + \log \log(M_{||} + A + B)), \quad B \lesssim \sqrt{\log(\varepsilon^{-1}(M_{||} + A + B))}, \\ \|W\|_\infty \vee S &\lesssim (\log(1/\varepsilon))^{3d+6} \sigma_{\text{data}}^{-d} (\log(M_{||} + A + B))^{d/2+1} \log \sigma_{\text{data}}^{-2}. \end{aligned}$$

By substituting the bounds for A and $M_{||}$ as specified in (17), we derive

$$\begin{aligned} L &\lesssim \log^3(1/\varepsilon) (\log \sigma_{\text{data}}^{-2} + \log \log(M + D\sigma_{\text{data}}^{-2})), \quad B \lesssim \sqrt{\log(\varepsilon^{-1}(M + D\sigma_{\text{data}}^{-2}))}, \\ \|W\|_\infty \vee S &\lesssim (\log(1/\varepsilon))^{3d+6} \sigma_{\text{data}}^{-d} (\log(M + D\sigma_{\text{data}}^{-2}))^{d/2+1} \log \sigma_{\text{data}}^{-2}. \end{aligned}$$

By appropriately rescaling ε with an absolute constant, we complete the proof. \square

A.3. Proof of Theorem 3.6

Step 1: proximity of log-potential estimates. Fix an arbitrary $y \in \mathbb{R}^D$. Consider $\psi^{(1)}, \psi^{(2)} \in \mathcal{F}_{\text{NN}}(L, W, S, B)$ of the form

$$\psi^{(j)}(y) = \varphi^{(j)}(V_{(j)}^\top y) - a_{(j)} \left(\|(I_D - V_{(j)} V_{(j)}^\top) y\|^2 - \gamma \sigma_\infty^2 (D - d) \right), \quad (20)$$

where $j \in \{1, 2\}$. Let $\{A_l^{(j)}, b_l^{(j)}\}_{l=1}^L$ be the corresponding weight matrices and biases of $\varphi^{(j)}$ for each $j \in \{1, 2\}$. Then, it follows that

$$\begin{aligned} |\psi^{(1)}(y) - \psi^{(2)}(y)| &\lesssim |\varphi^{(1)}(V_{(1)}^\top y) - \varphi^{(2)}(V_{(2)}^\top y)| \\ &\quad + |a_{(1)} - a_{(2)}| \left(\|(I_D - V_{(1)} V_{(1)}^\top) y\|^2 + D \right) \\ &\quad + \sigma_{\min}^{-2} \left| \|(I_D - V_{(1)} V_{(1)}^\top) y\|^2 - \|(I_D - V_{(2)} V_{(2)}^\top) y\|^2 \right|. \end{aligned} \quad (21)$$

We now analyze each term individually. Let us denote

$$\varepsilon = \max_{1 \leq l \leq L} (\|A_l^{(1)} - A_l^{(2)}\|_\infty \vee \|b_l^{(1)} - b_l^{(2)}\|).$$

Proceeding as in the proof of Lemma 5 in (Schmidt-Hieber, 2020), we obtain

$$\begin{aligned} &|\varphi^{(1)}(V_{(1)}^\top y) - \varphi^{(2)}(V_{(2)}^\top y)| \\ &= \left| \varphi^{(1)} \left((\|y\|_\infty + 1) V_{(1)}^\top \frac{y}{\|y\|_\infty + 1} \right) - \varphi^{(2)} \left((\|y\|_\infty + 1) V_{(2)}^\top \frac{y}{\|y\|_\infty + 1} \right) \right| \\ &\leq (\|W\|_\infty + 1)^{L+2} (L+1) (\|y\|_\infty + 1) \left(\varepsilon + \|A_1^{(1)} V_{(1)}^\top - A_1^{(2)} V_{(2)}^\top\|_\infty \right) \end{aligned}$$

Next, we note that $\|V_{(1)}\|_\infty \vee \|V_{(2)}\|_\infty \leq 1$ and, therefore,

$$\begin{aligned} \|A_1^{(1)}V_{(1)}^\top - A_1^{(2)}V_{(2)}^\top\|_\infty &\leq d\|A_1^{(1)} - A_1^{(2)}\|_\infty\|V_{(1)}\|_\infty + d\|A_1^{(2)}\|_\infty\|V_{(1)} - V_{(2)}\|_\infty \\ &\leq d(\varepsilon + B\|V_{(1)} - V_{(2)}\|_\infty). \end{aligned}$$

This immediately yields that

$$\begin{aligned} &|\varphi^{(1)}(V_{(1)}^\top y) - \varphi^{(2)}(V_{(2)}^\top y)| \\ &\leq (d+1)(\|W\|_\infty + 1)^{L+2}(L+1)(B+1)(\|y\| + 1)(\varepsilon + \|V_{(1)} - V_{(2)}\|_\infty). \end{aligned}$$

Next, we deduce that

$$|a_{(1)} - a_{(2)}| \left(\|(I_D - V_{(1)}V_{(1)}^\top)y\|^2 + D \right) \lesssim D(1 + \|y\|^2)|a_{(1)} - a_{(2)}|.$$

Put $V_\perp^{(j)} = I_D - V_{(j)}V_{(j)}^\top$. Then have that

$$\| \|V_\perp^{(1)}y\|^2 - \|V_\perp^{(2)}y\|^2 \| \leq \|V_\perp^{(1)} - V_\perp^{(2)}\| \cdot \|y\|(\|V_\perp^{(1)}y\| + \|V_\perp^{(2)}y\|) \leq 2\|V_\perp^{(1)} - V_\perp^{(2)}\| \cdot \|y\|^2,$$

where the last inequality uses the fact that $\|V_\perp^{(j)}\| \leq 1$ for all $j \in \{1, 2\}$. Therefore, using the fact that $\|V_\perp^{(1)} - V_\perp^{(2)}\| \lesssim D\|V_{(1)} - V_{(2)}\|_\infty$, we conclude that

$$\begin{aligned} |\psi^{(1)}(y) - \psi^{(2)}(y)| &\lesssim (\|W\|_\infty + 1)^{L+2}(L+1)(B+1)(\varepsilon + \|V_{(1)} - V_{(2)}\|_\infty)(1 + \|y\|^2) \\ &\quad + D(1 + \|y\|^2)|a_{(1)} - a_{(2)}| + \sigma_{\min}^{-2}D(1 + \|y\|^2)\|V_{(1)} - V_{(2)}\|_\infty. \end{aligned} \quad (22)$$

Recall from Lemma 3.3 that

$$\frac{\sigma}{\sigma_{\text{data}}} = \frac{-1 + \sqrt{1 + 4C^2}}{2C}, \quad C = \frac{\sigma_{\text{data}}e^{-bT}}{\gamma\sigma_T^2}.$$

Since $bT \gtrsim 1$, we can assume that $C \leq 1$. Consequently, by the mean value theorem, we have that

$$\sigma_{\text{data}}^2 - \sigma^2 = \sigma_{\text{data}}^2 \frac{\sqrt{1 + 4C^2} - 1}{2C^2} = \sigma_{\text{data}}^2 \frac{4C^2}{4C^2\sqrt{1 + 4C^2}\theta} \geq \frac{\sigma_{\text{data}}^2}{\sqrt{5}},$$

where $\theta \in (0, 1)$. Hence, we obtain

$$\left| \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma\sigma_T^2} \right) \right| \lesssim \sigma_{\text{data}}^{-2} \vee (\gamma\sigma_T^2)^{-1} \leq \sigma_{\min}^{-2} \vee (\gamma\sigma_T^2)^{-1}.$$

This observation, combined with Lemma 3.3, demonstrates that for any precision parameter ε , the approximation constructed in Theorem 3.4 belongs to the class of log-potential estimators defined in Definition 3.5 for appropriately chosen neural network parameters L , $\|W\|_\infty$, S , and B .

Deriving the final bound. We now apply the following generalization error bound result.

Lemma A.4 ((Puchkin et al., 2025), Theorem 1, adapted). *Grant assumptions of Theorem 3.6. Let $bT \gtrsim \log \sigma_{\text{data}}^{-2}$. Assume that there exist $J > 0$ such that, for every $\psi^{(1)}, \psi^{(2)} \in \mathcal{F}_{\text{NN}}$ of the form (20) and every $y \in \mathbb{R}^D$, it holds that*

$$\begin{aligned} &|\psi^{(1)}(y) - \psi^{(2)}(y)| \leq J(1 + \|y\|)^2 \\ &\quad \cdot \max_{1 \leq l \leq L} (\|A_l^{(1)} - A_l^{(2)}\|_\infty \vee \|b_l^{(1)} - b_l^{(2)}\|_\infty) \vee \|V_{(1)} - V_{(2)}\|_\infty \vee |\sigma_{(1)} - \sigma_{(2)}|. \end{aligned}$$

Then, for every $\delta \in (0, 1/2)$, with probability at least $(1 - 2\delta)$, it follows that

$$\text{KL}(\mathbf{p}_T, \hat{\mathbf{p}}_T) - \inf_{\psi \in \mathcal{F}_{\text{NN}}(L, W, S, B)} \text{KL}(\mathbf{p}_T, \mathbf{p}_T^\psi) \lesssim \sqrt{\Upsilon(n, \delta)} \inf_{\psi \in \mathcal{F}_{\text{NN}}(L, W, S, B)} \text{KL}(\mathbf{p}_T, \mathbf{p}_T^\psi) + \Upsilon(n, \delta),$$

where

$$\Upsilon(n, \delta) \lesssim D^2 \left(\log \frac{Jn}{\delta} + (M_{\parallel} + D\sigma_{\min}^{-2} + A)\sqrt{De^{-bT}} \right) \frac{SL \log(L(\|W\|_\infty + 1)B) \log n}{n}.$$

The proof of Lemma A.4 is a slight adaptation of the proof of Theorem 1 in (Puchkin et al., 2025). In that work step 6 establishes a uniform large-deviation bound that controls the difference between the empirical KL divergence and its population counterpart uniformly over the parameter class. Since an empirical risk minimizer is selected over the class of neural networks with at most S non-zero weights, we apply a union bound and observe that there are at most

$$\sum_{s=0}^S \binom{L\|W\|_\infty(\|W\|_\infty + 1)}{s} \leq (S+1)L^S(\|W\|_\infty + 1)^{2S}$$

ways to pick a sparsity pattern with at most S non-zero weights. In view of (22), we have that

$$\log J \lesssim L \log(L(\|W\|_\infty + 1)(B+1)D\sigma_{\min}^{-2}).$$

Applying Lemma A.4 together with the Young inequality, we obtain that, with probability at least $(1 - 2\delta)$,

$$\begin{aligned} \text{KL}(\mathfrak{p}_T, \hat{\mathfrak{p}}_T) &\lesssim \inf_{\psi \in \mathcal{F}_{\text{NN}}(L, W, S, B)} \text{KL}(\mathfrak{p}_T, \mathfrak{p}_T^\psi) \\ &\quad + \frac{D^2 S L^2 \log^2(L(\|W\|_\infty + 1)(B+1)) \log^2 n \log(1/\delta)}{n}. \end{aligned}$$

Applying Theorem 3.4 for some $0 < \varepsilon < A \wedge M_{\parallel} \wedge 1$, which will be determined later in the proof, we deduce that

$$\text{KL}(\mathfrak{p}_T, \hat{\mathfrak{p}}_T) \lesssim \varepsilon + \frac{D^2 \sigma_{\min}^{-d} (\log(1/\varepsilon))^{3d+14} (\log(M + D\sigma_{\min}^{-2}))^{d/2+5} \log^2 n \log(1/\delta)}{n},$$

with probability at least $(1 - 2\delta)$. Here, we put σ_{\min} instead of σ_{data} into the architecture given by Theorem 3.4 to ensure that the approximation $\log \tilde{\nu}_T$ belongs to the class $\mathcal{F}_{\text{NN}}(L, W, S, B)$. Therefore, setting $\varepsilon = (A \wedge M_{\parallel} \wedge 1)/n$ implies that, with probability at least $(1 - 2\delta)$,

$$\text{KL}(\mathfrak{p}_T, \hat{\mathfrak{p}}_T) \lesssim \frac{D^2 \sigma_{\min}^{-d} (\log n)^{3d+16} (\log(M + D\sigma_{\min}^{-2}))^{d/2+5} \log(1/\delta)}{n}.$$

Finally, from Theorem 3.4 we find that the architecture satisfies

$$\begin{aligned} L &\lesssim \log^3 n (\log \sigma_{\min}^{-2} + \log \log(M + D\sigma_{\min}^{-2})), \quad B \lesssim \sqrt{\log(n(M + D\sigma_{\min}^{-2}))}, \\ \|W\|_\infty \vee S &\lesssim (\log n)^{3d+6} \sigma_{\min}^{-d} (\log(M + D\sigma_{\min}^{-2}))^{d/2+2}. \end{aligned}$$

Furthermore, we have that $bT \gtrsim \log \log n + \log(M + D\sigma_{\min}^{-2})$. The proof is finished. \square

A.4. Proof of Lemma A.1

Let π^* be the joint density of (X_0^*, X_T^*) . Then, according to (Léonard, 2013) and (3), we have that, for all $x, y \in \mathbb{R}^D$,

$$\begin{aligned} \log \pi^*(x, y) &= \log \nu_0^*(x) + \log \mathfrak{q}_T(y | x) + \log \nu_T^*(y) \\ &= \log \nu_0^*(x) - \frac{D}{2} \log(2\pi\gamma\sigma_T^2) - \frac{\|y - e^{-bT}x\|^2}{2\gamma\sigma_T^2} + \log \nu_T^*(y). \end{aligned} \quad (23)$$

Let us verify that the form of the coupling described in claim of the proposition has the above form. From now on, assume that $0 < \sigma < \sigma_{\text{data}}$. Then, taking into account that

$$\det \begin{pmatrix} I_D & \sigma I_D \\ \sigma I_D & \sigma_{\text{data}}^2 I_D \end{pmatrix} = \det \left(\begin{pmatrix} 1 & \sigma \\ \sigma & \sigma_{\text{data}}^2 \end{pmatrix} \otimes I_D \right) = (\sigma_{\text{data}}^2 - \sigma^2)^D,$$

we deduce that

$$\log \pi^*(x, y) = -D \log(2\pi) - \frac{D}{2} \log(\sigma_{\text{data}}^2 - \sigma^2) - \frac{1}{2} \begin{pmatrix} x \\ y - \mu_T \end{pmatrix}^\top \begin{pmatrix} I_D & \sigma I_D \\ \sigma I_D & \sigma_{\text{data}}^2 I_D \end{pmatrix}^{-1} \begin{pmatrix} x \\ y - \mu_T \end{pmatrix}.$$

Applying the Schur complement argument, we obtain that

$$\begin{aligned} \log \pi^*(x, y) &= -D \log(2\pi) - \frac{D}{2} \log(\sigma_{\text{data}}^2 - \sigma^2) - \frac{1}{2} \|x\|^2 \left(1 - \frac{\sigma^2}{\sigma_{\text{data}}^2}\right)^{-1} \\ &+ x^\top (y - \mu_T) \left(1 - \frac{\sigma^2}{\sigma_{\text{data}}^2}\right)^{-1} \frac{\sigma}{\sigma_{\text{data}}^2} - \frac{1}{2} \|y - \mu_T\|^2 \left(\sigma_{\text{data}}^{-2} + \frac{\sigma^2}{\sigma_{\text{data}}^4} \left(1 - \frac{\sigma^2}{\sigma_{\text{data}}^2}\right)^{-1}\right). \end{aligned} \quad (24)$$

Therefore, it suffices to take σ such that the terms with $x^\top y$ coincide. Specifically, let

$$\left(1 - \frac{\sigma^2}{\sigma_{\text{data}}^2}\right)^{-1} \frac{\sigma}{\sigma_{\text{data}}^2} = \frac{e^{-bT}}{\gamma \sigma_T^2}. \quad (25)$$

Solving this quadratic equation, we conclude that

$$\frac{\sigma}{\sigma_{\text{data}}} = \frac{-1 + \sqrt{1 + 4C^2}}{2C}.$$

Note that $0 < \sigma < \sigma_{\text{data}}$ and, therefore, the first claim follows. In view of (23) and (24), we have that the right Schrödinger potential is proportional to (up to additive constant)

$$\log \nu_T^*(y) = \log \pi^*(x, y) - \log q_T(y | x) - \log \nu_0^*(x) \propto -\frac{\|y\|^2}{2} \left(\frac{1}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{1}{\gamma \sigma_T^2}\right) + \frac{y^\top \mu_T}{\sigma_{\text{data}}^2 - \sigma^2}.$$

Therefore, we can choose ν_T^* as shown on the right hand side of the above expression, and the second claim follows. Then, the left Schrödinger potential is expressed as

$$\begin{aligned} \log \nu_0^*(x) &= -\frac{\|x\|^2}{2} \left(\frac{\sigma_{\text{data}}^2}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{e^{-2bT}}{\gamma \sigma_T^2}\right) - \frac{x^\top \mu_T \sigma}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{\|\mu_T\|^2}{2(\sigma_{\text{data}}^2 - \sigma^2)} \\ &- \frac{D}{2} \log \left(\frac{2\pi(\sigma_{\text{data}}^2 - \sigma^2)}{\gamma \sigma_T^2}\right). \end{aligned}$$

Consequently, using (25), we obtain

$$\log \nu_0^*(x) = -\frac{\|x\|^2}{2} \left(\frac{\sigma_{\text{data}}^2 - e^{-bT} \sigma}{\sigma_{\text{data}}^2 - \sigma^2}\right) - \frac{x^\top \mu_T \sigma}{\sigma_{\text{data}}^2 - \sigma^2} - \frac{\|\mu_T\|^2}{2(\sigma_{\text{data}}^2 - \sigma^2)} - \frac{D}{2} \log \left(\frac{2\pi(\sigma_{\text{data}}^2 - \sigma^2)}{\gamma \sigma_T^2}\right).$$

This completes the proof. □

A.5. Proof of Lemma A.2

Recall that, by Lemma 3.3, the actual log-potential decomposes as

$$\begin{aligned} \log \nu_{||}^*(u) &= \log p_{||}(u) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2} \\ &- \log \int_{\mathbb{R}^d} \exp \left(\frac{e^{-bT} u^\top x}{\gamma \sigma_T^2} - \frac{e^{-2bT} \|x\|^2}{2\gamma \sigma_T^2} \right) \frac{p_0(x) dx}{\mathcal{T}_T[\nu_{||}^*](x)} - \frac{d}{2} \log(2\pi\gamma\sigma_T^2). \end{aligned}$$

Consequently, it suffices to approximate the log-density term (including the quadratic component) and the log-integral term separately.

Step 1: log-density approximation. Now we aim to approximate

$$f_1(u) = \log p_{||}(u) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2}$$

on a compact set $[-R, R]^d$. Our proof is based on the fundamental result on approximation capabilities of ReLU-neural networks given below.

Lemma A.5 ((Schmidt-Hieber, 2020), Theorem 5, adapted). *For any function $f \in \mathcal{H}^\beta([0, 1]^d, H)$ and any integers m and $N \geq (\beta \vee 1)^d \vee (H + 1)e^d$, there exists $\tilde{f} \in \text{NN}(L, W, S, 1)$ with*

$$L \lesssim m \log(d + \beta), \quad S \lesssim (d + \beta + 1)^{d+3} N m, \quad \|W\|_\infty \lesssim (d + \beta) N,$$

such that

$$\|\tilde{f} - f\|_{L^\infty([0, 1]^d)} \lesssim (H + 1)(d^2 + \beta^2) 6^d N 2^{-m} + H 3^\beta N^{-\beta/d}.$$

Let $\beta \in \mathbb{N}$ will be specified later in the proof. Note that, for any $f \in \mathcal{H}^\beta(\Omega, H)$,

$$\begin{aligned} & \sum_{|\mathbf{k}| < \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)} + \sum_{\mathbf{k} = \lfloor \beta \rfloor} \sup_{x \neq y} \frac{|\partial^{\mathbf{k}} f(x) - \partial^{\mathbf{k}} f(y)|}{\|x - y\|_\infty^{\beta - \lfloor \beta \rfloor}} \\ & \leq \max_{|\mathbf{k}| < \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)} |\{\mathbf{k} : |\mathbf{k}| < \beta\}| + \sqrt{d} \max_{|\mathbf{k}| = \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)} |\{\mathbf{k} : |\mathbf{k}| = \beta - 1\}| \\ & \leq 2\sqrt{d} \max_{|\mathbf{k}| \leq \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)} \binom{d + \beta - 1}{d}. \end{aligned}$$

Furthermore, we observe that, for any $R > 0$ and $f \in \mathcal{H}^\beta([-R, R]^d, H)$, it follows that $x \mapsto f(-R + 2Rx) \in \mathcal{H}^\beta([0, 1]^d, 2^\beta (R \vee 1)^\beta H)$. By Lemma 5.1, we have that

$$\max_{|\mathbf{k}| \leq \beta} \|\partial^{\mathbf{k}} f_1\|_{L^\infty([-R, R]^d)} \leq 2 \left(\frac{24\sqrt{C}\beta(\sqrt{d}R + 2\sqrt{d} + \sigma_{\text{data}})}{\sigma_{\text{data}}^2} \right)^\beta \beta!,$$

where $C \geq 1$ is an absolute constant. Using the bound

$$\binom{d + \beta - 1}{d} \leq (d + \beta - 1)^d \lesssim \beta^d, \quad (26)$$

we deduce that $f_1 \in \mathcal{H}^\beta([-R, R]^d, \mathbb{R}, H^{(1)})$ with

$$H^{(1)} \lesssim \beta^d \left(24\sigma_{\text{data}}^{-2} \sqrt{C\beta d} (R + 3) \right)^\beta \beta!.$$

Applying Lemma A.5, we obtain that there exists $\tilde{f}_1 \in \text{NN}(L^{(1)}, W^{(1)}, S^{(1)}, 2R)$ such that

$$\|\tilde{f}_1 - f_1\|_{L^\infty([-R, R]^d)} \lesssim R^\beta \beta^{d+2} (N 2^{-m} + 3^\beta N^{-\beta/d}) \left(24\sigma_{\text{data}}^{-2} \sqrt{C\beta d} (R + 3) \right)^\beta \beta!.$$

Applying the Stirling approximation and setting m such that $N 2^{-m} = 3^\beta N^{-\beta/d}$, we arrive at

$$\|\tilde{f}_1 - f_1\|_{L^\infty([-R, R]^d)} \lesssim \beta^{d+2} \left(\frac{72\sigma_{\text{data}}^{-2} \sqrt{C\beta d} \beta^{3/2} R (R + 3)}{e \cdot N^{1/d}} \right)^\beta.$$

Consequently, setting $\beta \asymp \log(1/\varepsilon)$ and

$$\begin{aligned} N &= \left(72\sigma_{\text{data}}^{-2} \sqrt{C\beta d} \beta^{3/2} R (R + 3) \right)^d + (\beta \vee 1)^d \vee (H^{(1)} + 1)e^d \\ &\asymp (\sigma_{\text{data}}^{-2} R^2)^d (\log(1/\varepsilon))^{3d/2}, \end{aligned}$$

we deduce that $\|\tilde{f}_1 - f_1\|_{L^\infty([-R, R]^d)} \lesssim \varepsilon$. In addition, we have that

$$m \asymp \log(N^{1+\beta/d} 3^{-\beta}) \lesssim (1 + \beta) \log N \lesssim \log(1/\varepsilon) \log(\sigma_{\text{data}}^{-2} R \log(1/\varepsilon)).$$

Therefore, by Lemma A.5, the architecture satisfies

$$\begin{aligned} L^{(1)} &\lesssim \log^2(1/\varepsilon) \log(\sigma_{\text{data}}^{-2} R), \\ S^{(1)} \vee \|W^{(1)}\|_{\infty} &\lesssim (\log(1/\varepsilon))^{5d/2+5} (\sigma_{\text{data}}^{-2} R^2)^d \log(\sigma_{\text{data}}^{-2} R). \end{aligned} \quad (27)$$

Step 2: log-integral term approximation. We now aim to approximate

$$f_2(u) = \log \int_{\mathbb{R}^d} \exp \left(\frac{e^{-bT} u^\top x}{\gamma \sigma_T^2} - \frac{e^{-2bT} \|x\|^2}{2\gamma \sigma_T^2} \right) \frac{\mathbf{p}_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}^*](x)}$$

on $[-R, R]^d$. Using Lemma 5.2 together with the condition (18), we obtain

$$\max_{|\mathbf{k}| \leq \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty([-R, R]^d)} \leq \left(\frac{2^{3+d/2} e^{-bT} \sqrt{C} \beta}{\gamma \sigma_T^2} \right)^\beta \beta!.$$

Therefore, by Stirling's formula,

$$\max_{|\mathbf{k}| \leq \beta} \|\partial^{\mathbf{k}} f\|_{L^\infty([-R, R]^d)} \lesssim \sqrt{\beta} \left(\frac{2^{3+d/2} \sqrt{C} \beta^{3/2}}{e \gamma \sigma_T^2} \right)^\beta.$$

Hence, applying Lemma A.5 with the Hölder norm parameter

$$H^{(2)} = 2\sqrt{d} \max_{|\mathbf{k}| \leq \beta} \|\partial^{\mathbf{k}} f_2\|_{L^\infty(\Omega)} \binom{d + \beta - 1}{d}$$

and using the bound (26), we obtain that there exists $\tilde{f}_2 \in \text{NN}(L^{(2)}, W^{(2)}, S^{(2)}, R)$ such that

$$\|\tilde{f}_2 - f_2\|_{L^\infty([-R, R]^d)} \lesssim R^\beta \beta^{d+1/2} \left(\frac{2^{3+d/2} \sqrt{C} \beta^{3/2}}{e \gamma \sigma_T^2} \right)^\beta \left(\beta^2 N 2^{-m} + 3^\beta N^{-\beta/d} \right).$$

Consequently, choosing $m \in \mathbb{N}$ such that $N 2^{-m} = 3^\beta N^{-\beta/d}$, we arrive at

$$\|\tilde{f}_2 - f_2\|_{L^\infty([-R, R]^d)} \lesssim \beta^{d+3} \left(\frac{3R \cdot 2^{3+d/2} \sqrt{C} \beta^{3/2}}{e \gamma \sigma_T^2 N^{1/d}} \right)^\beta.$$

Hence, setting $\beta \asymp \log(1/\varepsilon)$ and

$$N = \left(\frac{3R \cdot 2^{3+d/2} \sqrt{C} \beta^{3/2}}{\gamma \sigma_T^2} \right)^d + (\beta \vee 1)^d \vee (H^{(2)} + 1)e^d \asymp R^d (\log(1/\varepsilon))^{3d/2},$$

we conclude that

$$\|\tilde{f}_2 - f_2\|_{L^\infty([-R, R]^d)} \lesssim \varepsilon.$$

In addition, we have that

$$m \asymp (1 + \beta) \log N \lesssim \log(1/\varepsilon) (\log R + \log \log 1/\varepsilon).$$

Furthermore, Lemma A.5 suggests that the architecture satisfies

$$\begin{aligned} L^{(2)} &\lesssim m \log(1 + \beta) \lesssim \log^2(1/\varepsilon) \log R, \\ \|W^{(2)}\|_{\infty} \vee S^{(2)} &\lesssim \beta^{d+3} N m \lesssim (\log(1/\varepsilon))^{5d/2+5} R^d \log R. \end{aligned} \quad (28)$$

Step 3: ensuring regularity on the entire space. We now aim to clip the input of the resulting approximation to $[-R, R]^d$. In this step, we utilize standard results on concatenation and parallelization of ReLU neural networks (Nakada and Imaizumi, 2020). For every $x \in \mathbb{R}^d$ define the clipping function:

$$f_{\text{clip}}(x) = \min(\max(x, -R), R) = x + \text{ReLU}(-R - x) - \text{ReLU}(x - R).$$

Observe that f_{clip} is a two-layer ReLU neural network. By applying parallel stacking, we extend f_{clip} to clip a d -dimensional vector component-wise onto $[-R, R]^d$. Hence, the final approximation is of the form $\tilde{f}(x) = \tilde{f}_1(f_{\text{clip}}(x)) + \tilde{f}_2(f_{\text{clip}}(x)) - \frac{d}{2} \log(2\pi\gamma\sigma_T^2)$ satisfies the bound

$$\|\tilde{f} - \log \nu_{\parallel}^*\|_{L^\infty([-R, R]^d)} \leq \|\tilde{f}_1 - f_1\|_{L^\infty([-R, R]^d)} + \|\tilde{f}_2 - f_2\|_{L^\infty([-R, R]^d)} \lesssim \varepsilon.$$

By rescaling ε to an absolute constant we can prove the desired approximation error bound (19). Furthermore, from (27) and (28) it follows that $\tilde{f} \in \text{NN}(L, W, S, B)$ with

$$\begin{aligned} L &\lesssim \log^2(1/\varepsilon) \log(\sigma_{\text{data}}^{-2} R), & B &\lesssim R, \\ \|W\|_\infty \vee S &\lesssim (\log(1/\varepsilon))^{5d/2+5} (R\sigma_{\text{data}}^{-1})^d \log(\sigma_{\text{data}}^{-2} R). \end{aligned}$$

Furthermore, for all $x \in \mathbb{R}^d$, we have that

$$\min_{u \in [-R, R]^d} \log \nu_{\parallel}^*(u) - \varepsilon \leq \tilde{f}(x) \leq \max_{u \in [-R, R]^d} \log \nu_{\parallel}^*(u) + \varepsilon \leq M_{\parallel} + \varepsilon.$$

The proof is now complete. □

B. Proof of Lemma 4.1

The proof proceeds in three steps.

Step 1: factorization of the Gaussian density and the OU kernel. Let $S := (S_1^\top, \dots, S_J^\top)^\top \in \mathbb{R}^{D \times D}$. Since $S_j S_k^\top = 0$ for $j \neq k$ and $S_j S_j^\top = I_{d_j}$ for every j , we have $SS^\top = I_D$. As S is square, it is orthogonal. Therefore, for every $v \in \mathbb{R}^D$,

$$\|v\|^2 = \|Sv\|^2 = \sum_{j=1}^J \|S_j v\|^2.$$

Applying this identity with $v = x$ and $v = y - e^{-bT}x$, and using $d_1 + \dots + d_J = D$, we obtain

$$\begin{aligned} \rho_0(x) &= (2\pi)^{-D/2} \exp\left\{-\frac{\|x\|^2}{2}\right\} = \prod_{j=1}^J (2\pi)^{-d_j/2} \exp\left\{-\frac{\|S_j x\|^2}{2}\right\} = \prod_{j=1}^J \rho^{(j)}(S_j x), \\ \mathbf{q}_T(y|x) &= \prod_{j=1}^J \mathbf{q}_T^{(j)}(S_j y | S_j x), & \mathbf{p}_T(y) &= \prod_{j=1}^J \mathbf{p}_T^{(j)}(S_j y). \end{aligned}$$

Step 2: a product ansatz solves the full Schrödinger system. For each $j \in \{1, \dots, J\}$, the pair $(\nu_0^{(j)}, \nu_T^{(j)})$ satisfies the Schrödinger system for the d_j -dimensional sub-problem, namely

$$\begin{aligned} \rho^{(j)}(u) &= \nu_0^{(j)}(u) \int_{\mathbb{R}^{d_j}} \mathbf{q}_T^{(j)}(v|u) \nu_T^{(j)}(v) dv, & u &\in \mathbb{R}^{d_j}, \\ \mathbf{p}_T^{(j)}(v) &= \nu_T^{(j)}(v) \int_{\mathbb{R}^{d_j}} \mathbf{q}_T^{(j)}(v|u) \nu_0^{(j)}(u) du, & v &\in \mathbb{R}^{d_j}. \end{aligned}$$

Define

$$\widehat{\nu}_0(x) := \prod_{j=1}^J \nu_0^{(j)}(S_j x), \quad \widehat{\nu}_T(y) := \prod_{j=1}^J \nu_T^{(j)}(S_j y).$$

Consider the induced joint density

$$\widehat{\pi}(x, y) := \widehat{\nu}_0(x) \mathbf{q}_T(y | x) \widehat{\nu}_T(y), \quad x, y \in \mathbb{R}^D.$$

Since all factors are non-negative, Tonelli's theorem applies. Using the orthogonal change of variables $(u_1, \dots, u_J) = Sx$, whose Jacobian equals one, we get

$$\begin{aligned} \int_{\mathbb{R}^D} \widehat{\pi}(x, y) dx &= \prod_{j=1}^J \left[\nu_T^{(j)}(S_j y) \int_{\mathbb{R}^{d_j}} \mathbf{q}_T^{(j)}(S_j y | u_j) \nu_0^{(j)}(u_j) du_j \right] \\ &= \prod_{j=1}^J \mathbf{p}_T^{(j)}(S_j y) = \mathbf{p}_T(y). \end{aligned}$$

Similarly, using the orthogonal change of variables $(v_1, \dots, v_J) = Sy$, we obtain

$$\begin{aligned} \int_{\mathbb{R}^D} \widehat{\pi}(x, y) dy &= \prod_{j=1}^J \left[\nu_0^{(j)}(S_j x) \int_{\mathbb{R}^{d_j}} \mathbf{q}_T^{(j)}(v_j | S_j x) \nu_T^{(j)}(v_j) dv_j \right] \\ &= \prod_{j=1}^J \rho^{(j)}(S_j x) = \mathbf{p}_0(x). \end{aligned}$$

Hence, $(\widehat{\nu}_0, \widehat{\nu}_T)$ is a pair of Schrödinger potentials for the original problem on \mathbb{R}^D .

Step 3: identification with (ν_0^*, ν_T^*) . By Step 2, $(\widehat{\nu}_0, \widehat{\nu}_T)$ is a non-negative pair of Schrödinger potentials for the original problem on \mathbb{R}^D . By the uniqueness statement recalled above, see Theorem 2.12 in (Léonard, 2013), any such pair coincides with (ν_0^*, ν_T^*) up to reciprocal multiplicative constants. Hence there exists $c > 0$ such that

$$\widehat{\nu}_0 = c \nu_0^*, \quad \widehat{\nu}_T = \nu_T^* / c.$$

Replacing $(\nu_0^{(1)}, \nu_T^{(1)})$ with $(c\nu_0^{(1)}, \nu_T^{(1)}/c)$ leaves the first sub-problem unchanged, so we may absorb the constant c into the gauge of the first factor. Taking logarithms then yields

$$\log \nu_T^*(y) = \sum_{j=1}^J \log \nu_T^{(j)}(S_j y), \quad \log \nu_0^*(x) = \sum_{j=1}^J \log \nu_0^{(j)}(S_j x).$$

This completes the proof. □

C. Proofs of the results from Section 5

C.1. Proof of Lemma 5.1

The proof is quite technical, so we divide it into two steps.

Step 1: exact expression for the k -th derivative. Recalling the definition of $\mathbf{p}_{||}$ (see (9)), we have that

$$\nabla^k \left(\log \mathbf{p}_{||}(u) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2} \right) = \nabla^k \log \int_{\mathbb{R}^d} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z).$$

Proceeding similarly to the proof of Lemma 5.2 (Step 1), we deduce that, for every $k \in \mathbb{N}$,

$$\begin{aligned} & \nabla^k \log \int_{\mathbb{R}^d} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z) \\ &= (\sigma_{\text{data}}^2 Z(u))^{-k} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \exp \left\{ \frac{u^\top (z_1 + \dots + z_k)}{\sigma_{\text{data}}^2} - \sum_{j=1}^k \frac{\|z_j\|^2}{2\sigma_{\text{data}}^2} \right\} \prod_{j=1}^k d\mu(z_k), \end{aligned}$$

where

$$Z(u) = \int_{\mathbb{R}^d} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z)$$

and $P_k(z_1, \dots, z_k)$ satisfies the recurrence

$$P_1 = z_1, \quad P_2 = \frac{1}{2}(z_1 - z_2) \otimes (z_1 - z_2),$$

and

$$P_{k+1} = P_k \otimes (z_1 + \dots + z_k - kz_{k+1}), \quad k \in \mathbb{N}. \quad (29)$$

Unrolling (29) yields the closed form

$$P_k = \frac{1}{2}(z_1 - z_2) \otimes (z_1 - z_2) \otimes \dots \otimes \left(\sum_{j=1}^{k-1} z_j - (k-1)z_k \right).$$

Equivalently, defining

$$d\pi_u(z) := \frac{1}{Z(u)} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z), \quad (30)$$

we observe that

$$\begin{aligned} & \nabla^k \log \int_{\mathbb{R}^d} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z) \\ &= \frac{1}{2} \mathbb{E} \left[(X_1 - X_2) \otimes (X_1 - X_2) \otimes \dots \otimes \left(\sum_{j=1}^{k-1} X_j - (k-1)X_k \right) \right], \end{aligned} \quad (31)$$

where X_1, \dots, X_k are independent samples from π_u .

Step 2: reduction to sub-Gaussian conditional. We now reformulate the task using the notation above. Completing the exponent of (30) gives

$$d\pi_u(z) = \frac{1}{Z(u)} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z) \propto \exp \left\{ -\frac{\|u - z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z).$$

Let us introduce $S := \eta + Z$, where $\eta \sim \mathcal{N}(0, \sigma_{\text{data}}^2 I_d)$ is independent of $Z \sim \mu$. Then π_u can be considered as the conditional Law of Z given $S = u$. Applying Lemma 5.3 to the pair (η, Z) with $R = \|u\| + 2\tau\sqrt{d}$, one obtains

$$\|\eta - \mathbb{E}[\eta \mid S = u]\|_{\psi_2(\cdot \mid S=u)} \leq 6R + 6\sigma_{\text{data}}, \quad \|\mathbb{E}[\eta \mid S = u]\| \leq 3R + 3\sigma_{\text{data}}.$$

Here and further in the proof, $\|\cdot\|_{\psi_2(\cdot \mid S=u)}$ stands for the Orlicz norm with respect to the conditional distribution given $S = u$. Let us denote

$$\mu_T := \mathbb{E}[Z \mid S = u] \quad \text{and} \quad \tilde{\sigma} := \|Z - \mu_T\|_{\psi_2 \mid S=u}.$$

Since $Z = u - \eta$ whenever $S = u$, we observe that

$$\tilde{\sigma} \leq 6R + 6\sigma_{\text{data}}, \quad \text{and} \quad \|\mu_T\| \leq \|u\| + 3R + 3\sigma_{\text{data}} \leq 4R + 4\sigma_{\text{data}}. \quad (32)$$

Introducing $Z_j = (X_j - \mu_T)/\tilde{\sigma}$, $j \in \{1, \dots, k\}$, we find that

$$\begin{aligned} & \mathbb{E}(X_1 - X_2) \otimes (X_1 - X_2) \otimes \dots \otimes \left(\sum_{j=1}^{k-1} X_j - (k-1)X_k \right) \\ &= \tilde{\sigma}^k \mathbb{E}(Z_1 - Z_2) \otimes (Z_1 - Z_2) \otimes \dots \otimes \left(\sum_{j=1}^{k-1} Z_j - (k-1)Z_k \right). \end{aligned}$$

Since Z_1, \dots, Z_k are independent, centered-sub-Gaussian with the conditional Orlicz norm at most 1, we can apply Lemma D.2 and obtain that

$$\left\| \mathbb{E}(X_1 - X_2) \otimes (X_1 - X_2) \otimes \dots \otimes \left(\sum_{j=1}^{k-1} X_j - (k-1)X_k \right) \right\|_{\infty} \leq \left(\frac{6\tilde{\sigma}\sqrt{C}}{\sigma_{\text{data}}^2} \right)^k k!,$$

where $C \geq 1$ is an absolute constant. Substituting this bound into (31) and using (32), we deduce that

$$\begin{aligned} \left\| \nabla^k \log \int_{\mathbb{R}^d} \exp \left\{ \frac{u^\top z}{\sigma_{\text{data}}^2} - \frac{\|z\|^2}{2\sigma_{\text{data}}^2} \right\} d\mu(z) \right\|_{\infty} &\leq \left(\frac{36(R + \sigma_{\text{data}})\sqrt{C}}{\sigma_{\text{data}}^2} \right)^k k! \\ &= \left(\frac{36(\|u\| + 2\tau\sqrt{d} + \sigma_{\text{data}})\sqrt{C}}{\sigma_{\text{data}}^2} \right)^k k!. \end{aligned}$$

□

C.2. Proof of Lemma 5.2

Let us fix arbitrary $i_1, \dots, i_k \in \{1, \dots, d\}$. By the definition of the ℓ_{∞} -norm, it is enough to show that

$$\begin{aligned} \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right)_{i_1, \dots, i_k} \right| &\leq \exp \left\{ k(2M_{\parallel} + 3) + \frac{ke^{-2bT}\|y\|^2}{2(\gamma\sigma_T^2)^2} \right\} \\ &\cdot \left(\frac{2^{3+d/2}e^{-bT}\sqrt{Ck}}{\gamma\sigma_T^2} \right)^k k!. \end{aligned}$$

The proof of this inequality is quite intricate, so we divide it into several steps to enhance readability.

Step 1: exact expression for the k -th derivative. Let us introduce

$$\Phi(x, y) = \exp \left\{ \frac{\|y\|^2}{2\gamma\sigma_T^2} \right\} \cdot \frac{\mathbf{q}_{\parallel}(y|x)\rho_{\parallel}(x)}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)}, \quad Z(y) = \int_{\mathbb{R}^d} \Phi(x, y) dx, \quad x, y \in \mathbb{R}^d,$$

and show that, for any $y \in \mathbb{R}^d$,

$$\begin{aligned} & \nabla^k \left(\frac{\|y\|^2}{2\gamma\sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right) \\ &= \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k, \end{aligned}$$

where

$$P_1 = \frac{e^{-bT}x_1}{\gamma\sigma_T^2}, \quad P_2 = P_2(x_1, x_2) = \frac{1}{2} \left(\frac{e^{-bT}}{\gamma\sigma_T^2} \right)^2 (x_1 - x_2) \otimes (x_1 - x_2),$$

and

$$P_{k+1} = P_{k+1}(x_1, \dots, x_{k+1}) = \left(\frac{e^{-bT}}{\gamma\sigma_T^2} \right) P_k(x_1, \dots, x_k) \otimes (x_1 + \dots + x_k - kx_{k+1}). \quad (33)$$

Indeed, for $k = 1$, it holds that

$$\begin{aligned} \nabla \left(\frac{\|y\|^2}{2\gamma\sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) &= \left\langle \nabla \log \int_{\mathbb{R}^d} \Phi(x, y) dx, v \right\rangle \\ &= \left(\int_{\mathbb{R}^d} \frac{\partial \Phi(x, y)}{\partial y} dx \right) / \left(\int_{\mathbb{R}^d} \Phi(x, y) dx \right) \\ &= \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right) \int_{\mathbb{R}^d} x \Phi(x, y) dx. \end{aligned}$$

Differentiating the expression in the right-hand side once again, we obtain that

$$\begin{aligned} \nabla^2 \left(\frac{\|y\|^2}{2\gamma\sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \\ = \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^2 \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} x_1 \otimes (x_1 - x_2) \prod_{j=1}^2 \Phi(x_j, y) dx_1 dx_2. \end{aligned}$$

Due to exchangeability of x_1 and x_2 in the expression in the right-hand side, it holds that

$$\begin{aligned} \nabla^2 \left(\frac{\|y\|^2}{2\gamma\sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \\ = \frac{1}{2} \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^2 \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} x_2 \otimes (x_2 - x_1) \prod_{j=1}^2 \Phi(x_j, y) dx_1 dx_2 \\ + \frac{1}{2} \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^2 \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} x_1 \otimes (x_1 - x_2) \prod_{j=1}^2 \Phi(x_j, y) dx_1 dx_2 \\ = \frac{1}{2} \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^2 \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} (x_1 - x_2) \otimes (x_1 - x_2) \prod_{j=1}^2 \Phi(x_j, y) dx_1 dx_2. \end{aligned}$$

Proceeding by the induction, we observe that, for any $k \geq 2$,

$$\begin{aligned} \nabla^{k+1} \left(\frac{\|y\|^2}{2\gamma\sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \\ = \left[\left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \right] \\ = \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^k \left[\int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes \frac{\partial}{\partial y} \left(\prod_{j=1}^k \Phi(x_j, y) \right) dx_1 \dots dx_k \right. \\ \left. - \frac{k}{Z(y)} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \otimes \nabla Z(y) \right]. \end{aligned}$$

Since

$$\begin{aligned}
 & \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes \frac{\partial}{\partial y} \left(\prod_{j=1}^k \Phi(x_j, y) \right) dx_1 \dots dx_k \\
 &= \left(\frac{e^{-bT}}{\gamma \sigma_T^2} \right) \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes (x_1 + \dots + x_k) \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \\
 &= \frac{e^{-bT}/(\gamma \sigma_T^2)}{Z(y)} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes (x_1 + \dots + x_k) \prod_{j=1}^{k+1} \Phi(x_j, y) dx_1 \dots dx_{k+1}
 \end{aligned}$$

and

$$\begin{aligned}
 & \frac{1}{Z(y)} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \otimes \nabla Z(y) \\
 &= \frac{1}{Z(y)} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \otimes \int_{\mathbb{R}^d} \frac{\partial \Phi(x_{k+1}, y)}{\partial y} dx_{k+1} \\
 &= \frac{e^{-bT}/(\gamma \sigma_T^2)}{Z(y)} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes x_{k+1} \prod_{j=1}^{k+1} \Phi(x_j, y) dx_1 \dots dx_{k+1},
 \end{aligned}$$

we conclude that

$$\begin{aligned}
 & \nabla^{k+1} \left(\frac{\|y\|^2}{2\gamma \sigma_T^2} + \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \\
 &= \left(\frac{e^{-bT}/(\gamma \sigma_T^2)}{Z(y)} \right)^{k+1} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \otimes (x_1 + \dots + x_k - kx_{k+1}) \prod_{j=1}^{k+1} \Phi(x_j, y) dx_1 \dots dx_{k+1} \\
 &= \left(\frac{e^{-bT}/(\gamma \sigma_T^2)}{Z(y)} \right)^{k+1} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_{k+1} \prod_{j=1}^{k+1} \Phi(x_j, y) dx_1 \dots dx_{k+1}.
 \end{aligned}$$

It remains to note that $\nabla^k \|y\|^2$ equals to zero for all $k \geq 3$. For this reason,

$$\begin{aligned}
 & \nabla^k \left(\log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \\
 &= \left(\frac{e^{-bT}/(\gamma \sigma_T^2)}{Z(y)} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} P_k \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k
 \end{aligned} \tag{34}$$

for all $k \geq 3$ and $y \in \mathbb{R}^d$.

Step 2: lower bound on $Z(y)$. Let us recall that, according to the definition,

$$\begin{aligned}
 Z(y) &= \int_{\mathbb{R}^d} \Phi(x, y) dx = \int_{\mathbb{R}^d} \exp \left\{ \frac{\|y\|^2}{2\gamma \sigma_T^2} \right\} \cdot \frac{\mathbf{q}_{||}(y|x) \rho_{||}(x)}{\mathcal{T}_{||}[\nu_{||}](x)} dx \\
 &= (2\pi)^{-d} (\gamma \sigma_T^2)^{-d/2} \int_{\mathbb{R}^d} \exp \left\{ \frac{e^{-bT} y^\top x}{\gamma \sigma_T^2} - \frac{\|x\|^2}{2} \left(1 + \frac{e^{-2bT}}{\gamma \sigma_T^2} \right) \right\} \frac{dx}{\mathcal{T}_{||}[\nu_{||}](x)}.
 \end{aligned}$$

Taking into account that $\nu_{||} \leq e^{M_{||}}$ due to the conditions of the lemma, we obtain that $\mathcal{T}_{||}[\nu_{||}](x) \leq e^{M_{||}}$ for all $x \in \mathbb{R}^d$, and then

$$Z(y) \geq (2\pi)^{-d} (\gamma\sigma_T^2)^{-d/2} e^{-M_{||}} \int_{\mathbb{R}^d} \exp \left\{ \frac{e^{-bT} y^\top x}{\gamma\sigma_T^2} - \frac{\|x\|^2}{2} \left(1 + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right) \right\} dx.$$

The integral in the right-hand side admits a closed form. Namely, it is straightforward to check that

$$\begin{aligned} & \int_{\mathbb{R}^d} \exp \left\{ \frac{e^{-bT} y^\top x}{\gamma\sigma_T^2} - \frac{\|x\|^2}{2} \left(1 + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right) \right\} dx \\ &= (2\pi)^{d/2} \left(1 + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right)^{-d/2} \exp \left\{ \frac{e^{-2bT} \|y\|^2}{2\gamma\sigma_T^2(\gamma\sigma_T^2 + e^{-2bT})} \right\}. \end{aligned}$$

Hence, it holds that

$$Z(y) \geq (2\pi)^{-d/2} (\gamma\sigma_T^2 + e^{-2bT})^{-d/2} e^{-M_{||}} \exp \left\{ \frac{e^{-2bT} \|y\|^2}{2\gamma\sigma_T^2(\gamma\sigma_T^2 + e^{-2bT})} \right\}. \quad (35)$$

Step 3: upper bound on $\Phi(x, y)$. An upper bound on $\Phi(x, y)$ is a bit more complicated. Let us introduce

$$\mathcal{K}(T) = (1 - e^{-2bT})^{-5e^2\sqrt{d}} \exp \left\{ 2e^2\sqrt{d} \arcsin(e^{-bT}) \right\} \quad (36)$$

and

$$\alpha(T) = \frac{2be^2\mathcal{K}(T)}{\gamma\sqrt{d}} \arcsin(e^{-bT}). \quad (37)$$

Applying Lemma B.3 from (Puchkin et al., 2025) and taking into account that

$$\int_{\mathbb{R}^d} \log(\nu_{||}(y)) e^{-\|y\|^2/(2\gamma\sigma_T^2)} dy = 0,$$

we obtain that

$$\left(\frac{e^{M_{||}}}{\mathcal{T}_{||}[\nu_{||}](x)} \right) \leq \exp \left\{ \frac{\alpha(T)\|x\|^2}{2} + (M_{||} + 2 \log \mathcal{K}(T))\mathcal{K}(T) \right\}.$$

This yields that

$$\begin{aligned} (2\pi)^d (\gamma\sigma_T^2)^{d/2} e^{M_{||}} \Phi(x, y) &= \exp \left\{ \frac{e^{-bT} y^\top x}{\gamma\sigma_T^2} - \frac{\|x\|^2}{2} \left(1 + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right) \right\} \cdot \frac{e^{M_{||}}}{\mathcal{T}_{||}[\nu_{||}](x)} \\ &\leq \exp \left\{ \frac{e^{-bT} y^\top x}{\gamma\sigma_T^2} - \frac{1}{2} \left(1 - \alpha(T) + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right) \|x\|^2 + (M_{||} + 2 \log \mathcal{K}(T))\mathcal{K}(T) \right\}. \end{aligned} \quad (38)$$

Step 4: summing up the two upper bounds. The expression (34) yields that

$$\begin{aligned} & \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right)_{i_1, \dots, i_k} \right| \\ &= \left| \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} (P_k)_{i_1, \dots, i_k} \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k \right| \\ &\leq \left(\frac{e^{-bT}/(\gamma\sigma_T^2)}{Z(y)} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} |(P_k)_{i_1, \dots, i_k}| \prod_{j=1}^k \Phi(x_j, y) dx_1 \dots dx_k. \end{aligned}$$

In view of (35) and (38), it holds that

$$\begin{aligned} & \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right)_{i_1, \dots, i_k} \right| \\ & \leq (2\pi)^{-kd/2} \left(\frac{\gamma\sigma_T^2 + e^{-2bT}}{\gamma\sigma_T^2} \right)^{kd/2} \exp \left\{ k(M_{\parallel} + 2 \log \mathcal{K}(T))\mathcal{K}(T) - \frac{ke^{-2bT}\|y\|^2}{2\gamma\sigma_T^2(\gamma\sigma_T^2 + e^{-2bT})} \right\} \\ & \quad \cdot \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} |(P_k)_{i_1, \dots, i_k}| \prod_{j=1}^k \exp \left\{ \frac{e^{-bT} y^\top x_j}{\gamma\sigma_T^2} - \frac{1}{2} \left(1 - \alpha(T) + \frac{e^{-2bT}}{\gamma\sigma_T^2} \right) \|x_j\|^2 \right\} dx_1 \dots dx_k. \end{aligned}$$

Let us note that, due to the definition of P_k (see (33)), it is invariant under the substitutions

$$x_j := x_j - \frac{e^{-bT} y}{\check{\sigma}^2}, \quad 1 \leq j \leq k, \quad \text{where} \quad \check{\sigma}^2 = \gamma\sigma_T^2 + e^{-2bT} - \gamma\sigma_T^2 \alpha(T). \quad (39)$$

This implies that

$$\begin{aligned} & \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right)_{i_1, \dots, i_k} \right| \\ & \leq (2\pi)^{-kd/2} \left(\frac{\gamma\sigma_T^2 + e^{-2bT}}{\gamma\sigma_T^2} \right)^{kd/2} \exp \left\{ k(M_{\parallel} + 2 \log \mathcal{K}(T))\mathcal{K}(T) - \frac{ke^{-2bT}\|y\|^2}{2\gamma\sigma_T^2(\gamma\sigma_T^2 + e^{-2bT})} \right\} \\ & \quad \cdot \left(\frac{e^{-bT}}{\gamma\sigma_T^2} \right)^k \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} |(P_k)_{i_1, \dots, i_k}| \prod_{j=1}^k \exp \left\{ \frac{e^{-2bT}\|y\|^2}{2\gamma\sigma_T^2\check{\sigma}^2} - \frac{\check{\sigma}^2\|x_j\|^2}{2\gamma\sigma_T^2} \right\} dx_1 \dots dx_k. \end{aligned}$$

Using the identity

$$\frac{e^{-2bT}\|y\|^2}{2\gamma\sigma_T^2\check{\sigma}^2} - \frac{e^{-2bT}\|y\|^2}{2\gamma\sigma_T^2(\gamma\sigma_T^2 + e^{-2bT})} = \frac{e^{-2bT}\alpha(T)\|y\|^2}{2(\gamma\sigma_T^2 + e^{-2bT})\check{\sigma}^2}$$

and rearranging the terms, we obtain that

$$\begin{aligned} & \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right)_{i_1, \dots, i_k} \right| \\ & \leq \exp \left\{ k(M_{\parallel} + 2 \log \mathcal{K}(T))\mathcal{K}(T) + \frac{ke^{-2bT}\alpha(T)\|y\|^2}{2(\gamma\sigma_T^2 + e^{-2bT})\check{\sigma}^2} \right\} \\ & \quad \cdot \left(\frac{e^{-bT}}{\gamma\sigma_T^2} \right)^k \left(\frac{\gamma\sigma_T^2 + e^{-2bT}}{\check{\sigma}^2} \right)^{kd/2} \mathcal{J}_k, \end{aligned} \quad (40)$$

where

$$\mathcal{J}_k = (2\pi)^{-kd/2} \left(\frac{\gamma\sigma_T^2}{\check{\sigma}^2} \right)^{-kd/2} \int_{\mathbb{R}^d} \dots \int_{\mathbb{R}^d} |(P_k)_{i_1, \dots, i_k}| \prod_{j=1}^k \exp \left\{ -\frac{\check{\sigma}^2\|x_j\|^2}{2\gamma\sigma_T^2} \right\} dx_1 \dots dx_k. \quad (41)$$

Step 4: evaluation of \mathcal{J}_k . The integral \mathcal{J}_k defined in (41) can be considered as the expectation of a function of k i.i.d. standard Gaussian random vectors. Indeed, let $\eta_1, \dots, \eta_k \sim \mathcal{N}(0, \gamma\sigma_T^2\check{\sigma}^{-2}I_d)$, where σ_T and $\check{\sigma}$ are given by (3) and (39), respectively. Unrolling the recursion (33) for P_k , we arrive at

$$\mathcal{J}_k = \mathbb{E} \left| (\eta_1 - \eta_2)_{i_1} (\eta_1 - \eta_2)_{i_2} (\eta_1 + \eta_2 - 2\eta_3)_{i_3} \dots (\eta_1 + \dots + \eta_{k-1} - (k-1)\eta_k)_{i_k} \right|.$$

Let $\xi_1, \dots, \xi_k \sim \mathcal{N}(0, I_d)$, and let us introduce

$$S_j = \sum_{i=1}^j \xi_i, \quad j \in \{1, \dots, k\}.$$

Then

$$\mathcal{J}_k = \left(\frac{\gamma\sigma_T^2}{\check{\sigma}^2} \right)^k \mathbb{E} \left| (S_1 - \xi_2)_{i_1} (S_1 - \xi_2)_{i_2} (S_2 - 2\xi_3)_{i_3} \dots (S_{k-1} - (k-1)\xi_k)_{i_k} \right|.$$

Applying Lemma D.2, we deduce that

$$\mathbb{E} \left| (S_1 - \xi_2)_{i_1} (S_1 - \xi_2)_{i_2} (S_2 - 2\xi_3)_{i_3} \dots (S_{k-1} - (k-1)\xi_k)_{i_k} \right| \leq (6\sqrt{C})^k k!,$$

where $C \geq 1$ is a universal constant. Due to (40), this implies that

$$\begin{aligned} & \left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{\parallel}(y|x) \frac{\rho_{\parallel}(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}](x)} \right)_{i_1, \dots, i_k} \right| \\ & \leq \exp \left\{ k(M_{\parallel} + 2 \log \mathcal{K}(T)) \mathcal{K}(T) + \frac{ke^{-2bT} \alpha(T) \|y\|^2}{2(\gamma\sigma_T^2 + e^{-2bT})\check{\sigma}^2} \right\} \\ & \cdot \left(\frac{e^{-bT}}{\check{\sigma}^2} \right)^k \left(\frac{\gamma\sigma_T^2 + e^{-2bT}}{\check{\sigma}^2} \right)^{kd/2} (6\sqrt{C})^k k!. \end{aligned} \quad (42)$$

Step 5: simplifying the bound (42). On this step, we aim to make the bound (42) more user-friendly using the inequalities (16). First, note that

$$\arcsin x \leq 2x \quad \text{for all } x \in [0, 1] \quad \text{and} \quad -\log(1-x) \leq 2x \quad \text{for all } x \in [0, 1/2].$$

Then, in view of the definition of $\mathcal{K}(T)$ (see (36)), we have that

$$\begin{aligned} 1 & \leq \mathcal{K}(T) \leq \exp \left\{ -5e^2\sqrt{d} \log(1 - e^{-2bT}) + 2e^2\sqrt{d} \arcsin(e^{-bT}) \right\} \\ & \leq \exp \{ 14e^2\sqrt{d}e^{-bT} \} \leq 2, \end{aligned} \quad (43)$$

where the last inequality uses the conditions on bT from (16). This yields that

$$(M_{\parallel} + 2 \log \mathcal{K}(T)) \mathcal{K}(T) \leq 2M_{\parallel} + 4 \log 2 \leq 2M_{\parallel} + 3.$$

On the other hand, $\alpha(T)$, given by (37), satisfies

$$\alpha(T) = \frac{e^2 \mathcal{K}(T) (1 - e^{-2bT}) \arcsin(e^{-bT})}{\gamma\sigma_T^2 \sqrt{d}} \leq \frac{4e^{2-bT}}{\gamma\sigma_T^2 \sqrt{d}} \leq \frac{1}{2}.$$

The last inequality is also due to (16). This implies that

$$\frac{\alpha(T)}{(\gamma\sigma_T^2 + e^{-2bT})\check{\sigma}^2} = \frac{\alpha(T)}{(\gamma\sigma_T^2 + e^{-2bT})(\gamma\sigma_T^2 + e^{-2bT} - \alpha(T)\gamma\sigma_T^2)} \quad (44)$$

$$\leq \frac{\alpha(T)}{1 - \alpha(T)} \cdot \frac{1}{(\gamma\sigma_T^2)^2} \leq \frac{1}{(\gamma\sigma_T^2)^2} \quad (45)$$

and that

$$\frac{\gamma\sigma_T^2 + e^{-2bT}}{\check{\sigma}^2} = \frac{\gamma\sigma_T^2 + e^{-2bT}}{\gamma\sigma_T^2 + e^{-2bT} - \gamma\sigma_T^2 \alpha(T)} \leq \frac{1}{1 - \alpha(T)} \leq 2. \quad (46)$$

Substituting (44) and (46) into (42), we conclude that

$$\left| \left(\nabla^k \log \int_{\mathbb{R}^d} \mathbf{q}_{||}(y|x) \frac{\rho_{||}(x) dx}{\mathcal{T}_{||}[\nu_{||}](x)} \right)_{i_1, \dots, i_k} \right| \leq \exp \left\{ k(2M_{||} + 3) + \frac{ke^{-2bT} \|y\|^2}{2(\gamma\sigma_T^2)^2} \right\} \cdot \left(\frac{3 \cdot 2^{2+d/2} e^{-bT} \sqrt{C}}{\gamma\sigma_T^2} \right)^k k!.$$

The proof is now complete. □

C.3. Proof of Lemma 5.3

Step 1: Gaussian conjugacy. Let us introduce

$$f_S(t) = \mathbb{E} \varphi_\sigma(t - Y), \quad \text{where} \quad \varphi_\sigma(u) = (2\pi\sigma^2)^{-d/2} e^{-\|u\|^2/(2\sigma^2)}$$

Note that $f_S(s) > 0$ for all $s \in \mathbb{R}^d$. For any $\lambda \in \mathbb{R}^d$ and a test function $g : \mathbb{R}^d \rightarrow \mathbb{R}$, the change of variables $u = x + Y$ gives

$$\mathbb{E} \left[g(S) e^{\lambda^\top X} \right] = \mathbb{E} \left[\int_{\mathbb{R}^d} e^{\lambda^\top x} g(x + Y) \varphi_\sigma(x) dx \right] = \int_{\mathbb{R}^d} g(u) \mathbb{E} \left[e^{\lambda^\top (u - Y)} \varphi_\sigma(u - Y) \right] du.$$

Hence,

$$f_S(u) \mathbb{E} \left[e^{\lambda^\top X} \mid S = u \right] = \mathbb{E} \left[e^{\lambda^\top (u - Y)} \varphi_\sigma(u - Y) \right].$$

Using the identity

$$e^{\lambda^\top v} \varphi_\sigma(v) = e^{\sigma^2 \|\lambda\|^2/2} \varphi_\sigma(v - \sigma^2 \lambda),$$

we obtain that

$$\mathbb{E} [e^{\lambda^\top X} \mid S = s] = e^{\sigma^2 \|\lambda\|^2/2} \frac{f_S(s - \sigma^2 \lambda)}{f_S(s)}, \quad \lambda \in \mathbb{R}^d. \quad (47)$$

Step 2: two-sided control of f_S . Trivially,

$$f_S(t) \leq \|\varphi_\sigma\|_\infty = (2\pi\sigma^2)^{-d/2}.$$

On the other hand, by Markov's inequality,

$$\mathbb{P}(\|Y\| \leq 2\tau\sqrt{d}) = \mathbb{P}(\|Y\|^2 \leq 4d\tau^2) \geq 3/4.$$

On the event $\{\|Y\| \leq 2\tau\sqrt{d}\}$, we have $\|s - Y\| \leq \|s\| + 2\tau\sqrt{d} = R$. Because φ_σ is radially decreasing,

$$\varphi_\sigma(s - Y) \geq (2\pi\sigma^2)^{-d/2} e^{-R^2/(2\sigma^2)}.$$

Thus,

$$f_S(s) \geq \mathbb{E} \left[\varphi_\sigma(s - Y) \mathbb{1}\{\|Y\| \leq 2\sqrt{d}\tau\} \right] \geq \frac{3}{4} (2\pi\sigma^2)^{-d/2} e^{-R^2/(2\sigma^2)}.$$

Plugging both bounds into (47) yields

$$\mathbb{E} [e^{\lambda^\top X} \mid S = s] \leq \frac{4}{3} e^{\sigma^2 \|\lambda\|^2/2 + R^2/(2\sigma^2)}$$

for all $\lambda \in \mathbb{R}^d$.

1595 **Step 3: uncentered ψ_2 -norm via Gaussian integration.** Fix an arbitrary unit vector $v \in S^{d-1}$ and let $X_v := v^\top X$.
 1596 Taking $\lambda = \alpha v$ gives the one-dimensional marginal bound

$$1597 \mathbb{E}[e^{\alpha X_v} \mid S = s] \leq \frac{4}{3} e^{R^2/(2\sigma^2)} e^{\sigma^2 \alpha^2/2}$$

1600 for all $\alpha \in \mathbb{R}$. To convert this inequality into an upper bound on the ψ_2 -norm, we introduce a standard Gaussian random
 1601 variable $\xi \sim \mathcal{N}(0, 1)$, which is independent of X and S . Using the Gaussian integral identity

$$1602 e^{X_v^2/(8\sigma^2)} = \mathbb{E}_\xi \left[e^{\xi X_v/(2\sigma)} \right],$$

1603 we can rewrite the conditional expectation as

$$1604 \mathbb{E} \left[e^{X_v^2/(8\sigma^2)} \mid S = s \right] = \mathbb{E}_{X|S} \left[\mathbb{E}_\xi e^{\xi X_v/(2\sigma)} \mid S = s \right] = \mathbb{E}_\xi \mathbb{E}_{X|S} \left[e^{\xi X_v/(2\sigma)} \mid S = s \right].$$

1605 Using our moment generating function bound conditionally on ξ with $\alpha = \xi/(2\sigma)$, so that $\sigma^2 \alpha^2/2 = \xi^2/8$, we obtain

$$1606 \mathbb{E}_{X|S} \left[e^{\xi X_v/(2\sigma)} \mid S = s \right] \leq \frac{4}{3} e^{R^2/(2\sigma^2)} e^{\xi^2/8}.$$

1607 Taking the outer expectation over ξ , we get

$$1608 \mathbb{E}[e^{X_v^2/(8\sigma^2)} \mid S = s] \leq \frac{4}{3} e^{R^2/(2\sigma^2)} \mathbb{E}_\xi \left[e^{\xi^2/8} \right] = \frac{4e^{R^2/(2\sigma^2)}}{3\sqrt{1-1/4}} = \left(\frac{4}{3}\right)^{3/2} e^{R^2/(2\sigma^2)}.$$

1609 Let $u \geq 1$ be a constant. Since the mapping $x \mapsto x^{1/u^2}$ is concave for $x \geq 0$, Jensen's inequality yields that

$$1610 \mathbb{E}[e^{X_v^2/(8\sigma^2 u^2)} \mid S = s] \leq \left(\mathbb{E}[e^{X_v^2/(8\sigma^2)} \mid S = s] \right)^{1/u^2} \leq \exp \left\{ \frac{R^2/\sigma^2 + 3 \ln(4/3)}{2u^2} \right\}.$$

1611 We demand this bound to be at most $e^{1/2} < 2$, which requires $u^2 \geq R^2/\sigma^2 + 3 \ln(4/3)$. Since $3 \ln(4/3) \leq 1$, we choose
 1612 $u^2 = R^2/\sigma^2 + 1$, giving $8\sigma^2 u^2 = 8(R^2 + \sigma^2)$. Consequently,

$$1613 \mathbb{E}[e^{X_v^2/(8(R^2 + \sigma^2))} \mid S = s] \leq e^{1/2} \leq 2.$$

1614 By the definition of the Orlicz norm, the uncentered variable obeys

$$1615 \|X_v\|_{\psi_2(\cdot|_{S=s})} \leq \sqrt{8(R^2 + \sigma^2)} \leq \sqrt{8}(R + \sigma) \leq 3R + 3\sigma.$$

1616 **Step 4: centering via the triangle inequality.** Let $\mu_v := \mathbb{E}[X_v \mid S = s]$ denote the projection of the conditional mean.
 1617 Let us take $K := 3R + 3\sigma$. By applying Jensen's inequality to the strictly convex mapping $x \mapsto \exp(x^2/K^2)$, we have

$$1618 \exp(\mu_v^2/K^2) \leq \mathbb{E}[\exp(X_v^2/K^2) \mid S = s] \leq 2,$$

1619 which implies $|\mu_v| \leq K\sqrt{\ln 2}$. Since the ψ_2 -norm of any constant c is given by $|c|/\sqrt{\ln 2}$, we find

$$1620 \|\mu_v\|_{\psi_2(\cdot|_{S=s})} \leq K = 3R + 3\sigma.$$

1621 The triangle inequality yields that

$$1622 \|X_v - \mu_v\|_{\psi_2(\cdot|_{S=s})} \leq \|X_v\|_{\psi_2(\cdot|_{S=s})} + \|\mu_v\|_{\psi_2(\cdot|_{S=s})} \leq 2K = 6R + 6\sigma.$$

1623 Because this bound is uniform over all chosen test directions $v \in S^{d-1}$, taking the supremum over v establishes the final
 1624 sub-Gaussian vector estimate:

$$1625 \|X - \mu(s)\|_{\psi_2(\cdot|_{S=s})} \leq 6R + 6\sigma.$$

1626 This completes the proof.

1627 □

D. Auxiliary results

Lemma D.1. *Grant Assumption 3.1. Assume that there exists $M > 0$ such that $\log \nu_T^*(y) \leq M$ for all $y \in \mathbb{R}^D$ and*

$$\mathcal{T}_\infty[\log \nu_{\parallel}^*] = \mathcal{T}_\infty[\log \nu_{\perp}^*] = 0,$$

where $\log \nu_{\parallel}^*$ and $\log \nu_{\perp}^*$ are defined in Lemma 3.3. Assume further that bT satisfies the conditions (16). Then, there exist constants $A > 0$ and $B > 0$ such that

$$-A\|z\|^2 - B \leq \log \nu_{\parallel}^*(z) \leq M_{\parallel}, \quad \text{for all } z \in \mathbb{R}^d,$$

where $A \lesssim \sigma_{\text{data}}^{-2}$ and $B \lesssim d\tau^2/\sigma_{\text{data}}^2 + d|\log(\gamma\sigma_T^2)| + M\sqrt{d}e^{-bT}$, and $M_{\parallel} \lesssim M + D\sigma_{\text{data}}^{-2}$ with absolute hidden constants.

Lemma D.2. *Let $k \in \mathbb{N}$ and let Z_1, \dots, Z_k be independent random vectors in \mathbb{R}^d satisfying $\|Z_j\|_{\psi_2} \leq 1$ for all $j \in \{1, \dots, k\}$. Then there exists an absolute constant $C \geq 1$ such that for any $i_1, \dots, i_k \in \{1, \dots, d\}$ the following inequality holds:*

$$\mathbb{E} \left| (Z_{1,i_1} - Z_{2,i_1})(Z_{1,i_2} - Z_{2,i_2}) \dots \left(\sum_{j=1}^{k-1} Z_{j,i_k} - (k-1)Z_{k,i_k} \right) \right| \leq (6\sqrt{C})^k k!.$$

D.1. Proof of Lemma D.1

Write $G_{\parallel} := G$ for the matrix from Assumption 3.1, and let $G_{\perp} \in \mathbb{R}^{D \times (D-d)}$ denote any matrix with orthonormal columns satisfying $G_{\perp}^{\top} G_{\parallel} = 0$, so that $[G_{\parallel} \ G_{\perp}]$ is an orthonormal basis of \mathbb{R}^D and $\|G_{\parallel}^{\top} y\|^2 + \|G_{\perp}^{\top} y\|^2 = \|y\|^2$ for every $y \in \mathbb{R}^D$.

Step 1: latent decomposition and gauge. By Lemma 3.3, $\log \nu_T^*(y) = \log \nu_{\parallel}^*(G_{\parallel}^{\top} y) + \log \nu_{\perp}^*(G_{\perp}^{\top} y)$. Moreover, since log potentials are separable

$$\mathcal{T}_\infty[\log \nu_T^*] = \mathcal{T}_\infty[\log \nu_{\parallel}^*] + \mathcal{T}_\infty[\log \nu_{\perp}^*],$$

we may use this freedom to enforce both being equal to 0, thus, preserving $\mathcal{T}_\infty[\log \nu_T^*] = 0$. In this gauge, the closed-form expression of Lemma 3.3 reads

$$\log \nu_{\perp}^*(v) = -\frac{f_{\perp}}{2} \left(\|v\|^2 - \frac{(D-d)\gamma}{2b} \right), \quad v \in \mathbb{R}^{D-d}, \quad (48)$$

where $f_{\perp} = 1/(\sigma_{\text{data}}^2 - \sigma^2) - 1/(\gamma\sigma_T^2)$ and σ solves (25), $\log \nu_{\perp}^*(0) = f_{\perp}(D-d)\gamma/(4b)$. Conditions (16) together with (25) imply $\sigma^2 \leq \sigma_{\text{data}}^2/2$, so that $0 < f_{\perp} \leq 2/\sigma_{\text{data}}^2$. Evaluating Lemma 3.3 at $y = Gu$ for $u \in \mathbb{R}^d$,

$$\log \nu_{\parallel}^*(u) = \log \nu_T^*(Gu) - \log \nu_{\perp}^*(0) \leq M_{\parallel} := M - \frac{f_{\perp}(D-d)\gamma}{4b}, \quad (49)$$

and consequently $M_{\parallel} \leq M + f_{\perp}(D-d)\gamma/(4b) \lesssim M + D\sigma_{\text{data}}^{-2}$.

Step 2: lower paraboloid bound on $\log \nu_{\parallel}^*$. By Lemma 3.3,

$$-\log \nu_{\parallel}^*(u) = -\log p_{\parallel}(u) + \log \int_{\mathbb{R}^d} \frac{q_{\parallel}(u|x)p_0(x) dx}{\mathcal{T}_{\parallel}[\nu_{\parallel}^*](x)}, \quad u \in \mathbb{R}^d. \quad (50)$$

Introduce the probability measure $\tilde{\mu}$ on \mathbb{R}^d via

$$d\tilde{\mu}(z) := \frac{(2\pi\sigma_{\text{data}}^2)^{-d/2}}{p_{\parallel}(0)} \exp\left(-\frac{\|z\|^2}{2\sigma_{\text{data}}^2}\right) d\mu(z).$$

A direct computation shows that

$$p_{\parallel}(u) = p_{\parallel}(0) \exp\left(-\frac{\|u\|^2}{2\sigma_{\text{data}}^2}\right) \int_{\mathbb{R}^d} \exp\left(\frac{u^{\top} z}{\sigma_{\text{data}}^2}\right) d\tilde{\mu}(z).$$

Taking logarithms and applying Jensen's inequality to the resulting integral, we obtain

$$-\log p_{||}(u) = -\log p_{||}(0) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2} - \log \int_{\mathbb{R}^d} \exp\left(\frac{u^\top z}{\sigma_{\text{data}}^2}\right) d\tilde{\mu}(z) \quad (51)$$

$$\leq -\log p_{||}(0) + \frac{\|u\|^2}{2\sigma_{\text{data}}^2} - \frac{u^\top \tilde{m}}{\sigma_{\text{data}}^2} \leq -\log p_{||}(0) + \frac{\|u\|^2}{\sigma_{\text{data}}^2} + \frac{\|\tilde{m}\|^2}{2\sigma_{\text{data}}^2}, \quad (52)$$

where $\tilde{m} := \int z d\tilde{\mu}(z)$. The probability measure $\tilde{\mu}$ is the conditional law $\mathcal{L}(Z | Z + \sigma_{\text{data}} \xi_0 = 0)$ for $\xi_0 \sim \mathcal{N}(0, I_d)$ independent of $Z \sim \mu$. Lemma 5.3 applied to the pair $(\sigma_{\text{data}} \xi_0, Z)$ at $s = 0$, together with $Z = -\sigma_{\text{data}} \xi_0$ on the conditioning event, yields $\|\tilde{m}\| \leq 6\sqrt{d}\tau + 3\sigma_{\text{data}}$, so by $(a+b)^2 \leq 2(a^2 + b^2)$ and $\sigma_{\text{data}} \leq 1$, $\|\tilde{m}\|^2/(2\sigma_{\text{data}}^2) \lesssim d\tau^2/\sigma_{\text{data}}^2 + 1$. Markov's inequality gives $\mu(\{\|z\| \leq 2\sqrt{d}\tau\}) \geq 3/4$, on which event $e^{-\|z\|^2/(2\sigma_{\text{data}}^2)} \geq e^{-2d\tau^2/\sigma_{\text{data}}^2}$, so $p_{||}(0) \geq (3/4)(2\pi\sigma_{\text{data}}^2)^{-d/2} \exp(-2d\tau^2/\sigma_{\text{data}}^2)$. Using $\sigma_{\text{data}} \leq 1$ to bound $(d/2) \log(2\pi\sigma_{\text{data}}^2) \leq (d/2) \log(2\pi) \lesssim d$, we obtain $-\log p_{||}(0) \lesssim d + d\tau^2/\sigma_{\text{data}}^2$. Plugging into (51) yields

$$-\log p_{||}(u) \lesssim \frac{\|u\|^2}{\sigma_{\text{data}}^2} + \frac{d\tau^2}{\sigma_{\text{data}}^2} + d. \quad (53)$$

Set

$$\mathcal{K}(T) := (1 - e^{-2bT})^{-5e^2\sqrt{d}} \exp(2e^2\sqrt{d} \arcsin(e^{-bT})), \quad \alpha(T) := \frac{2be^2\mathcal{K}(T)}{\gamma\sqrt{d}} \arcsin(e^{-bT}).$$

For the second summand of (50), the gauge on $\log \nu_{||}^*$ together with the upper bound (49) allow us to apply Step 3 of the proof of Lemma 5.2 on \mathbb{R}^d with $M_{||}$ in place of M , which gives

$$\frac{e^{M_{||}}}{\mathcal{T}_{||}[\nu_{||}^*](x)} \leq \exp\left\{\frac{1}{2}\alpha(T)\|x\|^2 + (M_{||} + 2\log \mathcal{K}(T))\mathcal{K}(T)\right\}, \quad x \in \mathbb{R}^d, \quad (54)$$

and multiplying by $e^{-M_{||}}$ and integrating against $q_{||}(u|x)p_0(x)$,

$$\begin{aligned} \log \int_{\mathbb{R}^d} \frac{q_{||}(u|x)p_0(x) dx}{\mathcal{T}_{||}[\nu_{||}^*](x)} &\leq M_{||}(\mathcal{K}(T) - 1) + 2\log \mathcal{K}(T) \cdot \mathcal{K}(T) \\ &\quad + \log \int_{\mathbb{R}^d} \exp\left(\frac{1}{2}\alpha(T)\|x\|^2\right) q_{||}(u|x)p_0(x) dx. \end{aligned} \quad (55)$$

The remaining integral is Gaussian. Setting $\check{\sigma}^{-2} := 1 - \alpha(T) + e^{-2bT}/(\gamma\sigma_T^2)$ and completing the square in x ,

$$\int_{\mathbb{R}^d} \exp\left(\frac{1}{2}\alpha(T)\|x\|^2\right) q_{||}(u|x)p_0(x) dx = (2\pi\gamma\sigma_T^2)^{-d/2} \check{\sigma}^d \exp\left\{-\frac{\|u\|^2(1 - \check{\sigma}^2 e^{-2bT}/(\gamma\sigma_T^2))}{2\gamma\sigma_T^2}\right\}. \quad (56)$$

Conditions (16) imply

$$\mathcal{K}(T) \leq 2, \quad \mathcal{K}(T) - 1 \lesssim \sqrt{d}e^{-bT}, \quad \check{\sigma}^2 \leq 2, \quad \frac{\check{\sigma}^2 e^{-2bT}}{\gamma\sigma_T^2} \leq \frac{1}{2}. \quad (57)$$

Taking logarithms in (56), the $\|u\|^2$ contribution is non-positive and may be dropped for an upper bound, $-(d/2) \log(2\pi\gamma\sigma_T^2) + d \log \check{\sigma} \lesssim d |\log(\gamma\sigma_T^2)| + d$. Using (57) to further bound $|M_{||}(\mathcal{K}(T) - 1)| \lesssim |M_{||}|\sqrt{d}e^{-bT}$ and $|2\log \mathcal{K}(T) \cdot \mathcal{K}(T)| \lesssim \sqrt{d}e^{-bT} \lesssim 1$, we deduce from (55) that

$$\log \int_{\mathbb{R}^d} \frac{q_{||}(u|x)p_0(x) dx}{\mathcal{T}_{||}[\nu_{||}^*](x)} \lesssim |M_{||}|\sqrt{d}e^{-bT} + d |\log(\gamma\sigma_T^2)| + d. \quad (58)$$

Combining (53) and (58) via (50),

$$-\log \nu_{\parallel}^*(u) \lesssim \frac{\|u\|^2}{\sigma_{\text{data}}^2} + \frac{d\tau^2}{\sigma_{\text{data}}^2} + d |\log(\gamma\sigma_T^2)| + d + |M_{\parallel}| \sqrt{d} e^{-bT}. \quad (59)$$

Step 3: assembly on \mathbb{R}^D . Combining (59) with (48) via Lemma 3.3, using $f_{\perp} \leq 2/\sigma_{\text{data}}^2$, $\|G_{\parallel}^{\top} y\|^2 + \|G_{\perp}^{\top} y\|^2 = \|y\|^2$, and $|M_{\parallel}| \leq M + f_{\perp}(D-d)\gamma/(4b)$,

$$\begin{aligned} \log \nu_T^*(y) &= \log \nu_{\parallel}^*(G_{\parallel}^{\top} y) + \log \nu_{\perp}^*(G_{\perp}^{\top} y) \\ &\gtrsim -\frac{\|y\|^2}{\sigma_{\text{data}}^2} - \frac{d\tau^2}{\sigma_{\text{data}}^2} - d |\log(\gamma\sigma_T^2)| - d - M\sqrt{d} e^{-bT} + \frac{f_{\perp}(D-d)\gamma}{4b} (1 - \sqrt{d} e^{-bT}). \end{aligned}$$

Conditions (16) imply $\sqrt{d} e^{-bT} \leq 1/2$, so the last term is non-negative and may be dropped. This yields the claimed bound and completes the proof. \square

D.2. Proof of Lemma D.2

Let us fix an arbitrary multi-index $(i_1, \dots, i_k) \subseteq \{1, \dots, d\}^k$. Due to the triangle inequality, it holds that

$$\begin{aligned} &\mathbb{E} \left| (Z_{1,i_1} - Z_{2,i_1})(Z_{1,i_2} - Z_{2,i_2}) \dots \left(\sum_{j=1}^{k-1} Z_{j,i_k} - (k-1)Z_{k,i_k} \right) \right| \quad (60) \\ &\leq \mathbb{E} \sum_{(\omega_0, \dots, \omega_k) \in \{0,1\}^k} |Z_{1,i_1}|^{\omega_0} |Z_{2,i_1}|^{1-\omega_0} \cdot \prod_{j=1}^{k-1} \left(\left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{\omega_j} \cdot j^{1-\omega_j} \cdot |Z_{j+1,i_{j+1}}|^{1-\omega_j} \right). \end{aligned}$$

For a fixed $(\omega_0, \dots, \omega_{k-1}) \in \{0, 1\}^k$, by the Hölder inequality, we have that

$$\begin{aligned} &\mathbb{E} |Z_{1,i_1}|^{\omega_0} |Z_{2,i_1}|^{1-\omega_0} \cdot \prod_{j=1}^{k-1} \left(\left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{\omega_j} \cdot j^{1-\omega_j} \cdot |Z_{j+1,i_{j+1}}|^{1-\omega_j} \right) \quad (61) \\ &\leq \prod_{j=1}^{k-1} j^{1-\omega_j} (\mathbb{E} |Z_{1,i_1}|^{3\omega_0} |Z_{2,i_1}|^{3-3\omega_0})^{1/3} \left(\mathbb{E} \prod_{j=1}^{k-1} |Z_{j+1,i_{j+1}}|^{3(1-\omega_j)} \right)^{1/3} \left(\mathbb{E} \prod_{j=1}^{k-1} \left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{3\omega_j} \right)^{1/3} \\ &= \prod_{j=1}^{k-1} j^{1-\omega_j} (\mathbb{E} |Z_{1,i_1}|^{3\omega_0} \mathbb{E} |Z_{2,i_1}|^{3-3\omega_0})^{1/3} \left(\prod_{j=1}^{k-1} \mathbb{E} |Z_{j+1,i_{j+1}}|^{3(1-\omega_j)} \right)^{1/3} \left(\mathbb{E} \prod_{j=1}^{k-1} \left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{3\omega_j} \right)^{1/3}. \end{aligned}$$

Using properties of sub-Gaussian random variables (see, for instance, Proposition 2.5.2 in (Vershynin, 2018)), we obtain that

$$(\mathbb{E} |Z_{1,i_1}|^{3\omega_0})^{1/3} \leq (3C)^{\omega_0/2}, \quad (\mathbb{E} |Z_{2,i_1}|^{2-2\omega_0})^{1/3} \leq (3C)^{(1-\omega_0)/2},$$

and

$$\left(\mathbb{E} |Z_{j+1,i_{j+1}}|^{3(1-\omega_j)} \right)^{1/3} \leq (3C)^{(1-\omega_j)/2}, \quad \text{for all } 1 \leq j \leq k,$$

where $C \geq 1$ is an absolute constant. Furthermore, let us introduce $\Omega = \omega_1 + \dots + \omega_{k-1}$. Then, applying the Hölder inequality and taking into account that $\|Z_{1,i_{j+1}} + \dots + Z_{j,i_{j+1}}\|_{\psi_2} \lesssim \sqrt{j}$, we deduce that

$$\begin{aligned} \left(\mathbb{E} \prod_{j=1}^{k-1} \left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{3\omega_j} \right)^{1/3} &= \left(\mathbb{E} \prod_{j:\omega_j=1} \left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{3\omega_j} \right)^{1/3} \\ &\leq \left(\prod_{j:\omega_j=1} \mathbb{E} \left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{3\Omega\omega_j} \right)^{1/(3\Omega)} \\ &\leq \prod_{j:\omega_j=1} (3C\Omega j)^{\omega_j/2} = (3C\Omega)^{\Omega/2} \prod_{j=1}^k j^{\omega_j/2}. \end{aligned}$$

Substituting these bounds into (61), we conclude that

$$\begin{aligned} \mathbb{E} |Z_{1,i_1}|^{\omega_0} |Z_{2,i_1}|^{1-\omega_0} \cdot \prod_{j=1}^{k-1} \left(\left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{\omega_j} \cdot j^{1-\omega_j} \cdot |Z_{j+1,i_{j+1}}|^{1-\omega_j} \right) \\ \leq \sqrt{3C} \prod_{j=1}^{k-1} \left(j^{1-\omega_j} \cdot (3C)^{(1-\omega_j)/2} \right) \cdot (3C\Omega)^{\Omega/2} \prod_{j=1}^k j^{\omega_j/2} \\ = (3C)^{k/2} \sqrt{(k-1)!} \cdot \Omega^{\Omega/2} \prod_{j=1}^{k-1} j^{(1-\omega_j)/2}. \end{aligned}$$

Let us note that the latter term in the right-hand side does not exceed

$$\prod_{j=1}^{k-1} j^{(1-\omega_j)/2} \leq \sqrt{\frac{(k-1)!}{\Omega!}}.$$

Since $l! \geq (l/e)^l$ for all $l \in \mathbb{N}$, we obtain that

$$\Omega^{\Omega/2} \prod_{j=1}^{k-1} j^{(1-\omega_j)/2} \leq \Omega^{\Omega/2} \sqrt{\frac{(k-1)!}{\Omega!}} \leq e^{k/2} \sqrt{(k-1)!},$$

and then

$$\mathbb{E} |Z_{1,i_1}|^{\omega_0} |Z_{2,i_1}|^{1-\omega_0} \cdot \prod_{j=1}^{k-1} \left(\left| \sum_{s=1}^j Z_{s,i_{j+1}} \right|^{\omega_j} \cdot j^{1-\omega_j} \cdot |Z_{j+1,i_{j+1}}|^{1-\omega_j} \right) \leq (3Ce)^{k/2} (k-1)!$$

Taking into account that the right-hand side of (60) includes 2^k terms, we conclude that

$$\begin{aligned} \mathbb{E} \left| (Z_{1,i_1} - Z_{2,i_1})(Z_{1,i_2} - Z_{2,i_2}) \dots \left(\sum_{j=1}^{k-1} Z_{j,i_k} - (k-1)Z_{k,i_k} \right) \right| &\leq (2\sqrt{3Ce})^k k! \\ &\leq (6\sqrt{C})^k k!. \end{aligned}$$

The proof is finished. □