

# FLOW MATCHING ACHIEVES ALMOST MINIMAX OPTIMAL CONVERGENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Flow matching (FM) has gained significant attention as a simulation-free generative model. Unlike diffusion models, which are based on stochastic differential equations, FM employs a simpler approach by solving an ordinary differential equation with an initial condition from a normal distribution, thus streamlining the sample generation process. This paper discusses the convergence properties of FM for large sample size under the  $p$ -Wasserstein distance, a measure of distributional discrepancy. We establish that FM can achieve an almost minimax optimal convergence rate for  $1 \leq p \leq 2$ , presenting the first theoretical evidence that FM can reach convergence rates comparable to those of diffusion models. Our analysis extends existing frameworks by examining a broader class of mean and variance functions for the vector fields and identifies specific conditions necessary to attain almost optimal rates.

## 1 INTRODUCTION

Flow matching (FM) (Lipman et al., 2023; Albergo and Vanden-Eijnden, 2023; Liu et al., 2023b) is a recent simulation-free generative model that produces samples of the target distribution by solving an ordinary differential equation (ODE) initialized with a source normal distribution. The vector field to define the ODE is trained by neural networks with the teaching data of random conditional vectors. This approach bypasses the computationally intensive Monte Carlo sampling required in the diffusion model, which is currently the standard in generative modeling. Various variations have been proposed to refine the learning of vector fields, such as OT-CFM (Tong et al., 2024), rectified flow (Liu et al., 2023b), consistent velocity field (Yang et al., 2024), equivariant flow (Klein et al., 2023), etc. A series of studies also emerge from the viewpoint of interpolating distributions (Albergo et al., 2023c;a).

FM has already been applied to various domains with promising performance. Among many others, the rectified flow method has been extended to high-resolution text image generation (Esser et al., 2024), and there are also many works on the application of FM to molecule generation (Hoogeboom et al., 2022; Guan et al., 2023; Bose et al., 2023; Dunn and Koes, 2024), text generation (Hu et al., 2024), speech generation (Le et al., 2023), motion synthesis (Hu et al., 2023), etc.

Although the methods have been developed on the solid theoretical basis of the flows and continuity equation, their statistical behaviors remain less understood. Recent works have established the convergence of the FM estimator to the true distribution under some distributional metrics (Albergo and Vanden-Eijnden, 2023; Benton et al., 2023b). Beyond the convergence, more detailed understandings, such as convergence rates, are still an open question. In contrast, diffusion models have gained various theoretical understandings, including the convergence rate in terms of the number of steps (Chen et al., 2023; Benton et al., 2023a) and the sample size (Oko et al., 2023; Zhang et al., 2024). Among others, Oko et al. (2023) has shown that diffusion models achieve the minimax optimal convergence rate for a large sample size under the total variation metric and the almost minimax optimal rate under the 1-Wasserstein distance, where the max is taken over the true densities of the Besov space. This result theoretically supports the high generation ability of diffusion models.

This paper aims to bridge this gap by demonstrating that FM can achieve an almost minimax optimal convergence rate for a large sample size under the  $p$ -Wasserstein distance  $W_p$  for  $1 \leq p \leq 2$ , suggesting that FM has a theoretical ability comparable to diffusion models. This problem is significant for comparing the ability of FM methods and diffusion models, and revealing the difference between SDE and ODE in the generative models. Drawing on the methodologies of Oko et al. (2023),

our analysis not only extends to a broader class of mean and variance parameters of Gaussian smoothing for conditional vector fields, but also specifies the conditions on these parameters under which the almost minimax optimal convergence rate can be achieved.

The contributions of this paper are as follows.

- We establish that a widely used class of conditional FM methods achieves an almost minimax optimal convergence rate under the  $p$ -Wasserstein distance ( $1 \leq p \leq 2$ ), marking the first theoretical demonstration of such optimal performance of FM.
- We provide an analytical derivation of the convergence rate under various settings of the parameters, mean and variance, to make a path that connects a source and target point.
- We reveal that the variance parameter, which specifies the contribution of the source normal distribution, must be decreased around the target at a specific rate to attain an almost minimax optimal convergence rate.

## 2 FLOW MATCHING

Throughout the paper, data are in the  $d$ -dimensional space  $\mathbb{R}^d$ . The  $d$ -dimensional normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $V$  is denoted by  $N_d(\boldsymbol{\mu}, V)$ . For a probability  $P_a$  with index  $a$ , the lowercase  $p_a$  denotes its probability density function (p.d.f.).

### 2.1 REVIEW OF FLOW MATCHING

This subsection provides a general review of FM, following Lipman et al. (2023) and Tong et al. (2024). The aim of FM is to generate samples from the true probability  $P_{true}$ . FM methods realize it by a flow  $\varphi_{[\tau]}(\boldsymbol{x})$  ( $\tau \in [0, 1]$ )<sup>1</sup> that maps a sample from the standard normal distribution  $N_d(0, I_d)$  to that of  $P_{true}$ . The flow  $\varphi_{[\tau]}(\boldsymbol{x})$  is defined by a solution to the ODE

$$\frac{d}{d\tau} \boldsymbol{x}_{[\tau]} = \boldsymbol{v}_{[\tau]}(\boldsymbol{x}_{[\tau]}) \quad (\tau \in [0, 1])$$

given by a desired vector field  $\boldsymbol{v}_{[\tau]}$ . FM generates a sample by solving the ODE with an initial point  $\boldsymbol{x}_{[0]}$  from  $P_{[0]} = N_d(0, I_d)$ ; in other words, the distribution at time  $\tau$  is the pushforward  $P_{[\tau]} = \varphi_{[\tau]\#} P_{[0]}$ . The pushforward  $P_{[1]}$  is expected to approximate  $P_{true}$ . In practice, we need to construct the vector field given training data  $\{\boldsymbol{x}^i\}_{i=1}^n$  of size  $n$ , which is i.i.d. samples from  $P_{true}$ .

The relation between the vector field  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x})$  and the p.d.f.  $p_{[\tau]}(\boldsymbol{x})$  is given by the *continuity equation*:

$$\frac{\partial}{\partial \tau} p_{[\tau]}(\boldsymbol{x}) + \operatorname{div}(p_{[\tau]}(\boldsymbol{x}) \boldsymbol{v}_{[\tau]}(\boldsymbol{x})) = 0.$$

Typically, a neural network (NN) is used to construct  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x})$ . However, it is not obvious how to prepare the desired  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x})$  to teach NN. In FM methods, conditional random vectors  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x}_{[\tau]}|\boldsymbol{z})$  given  $\boldsymbol{z}$ , which are to be easily prepared, are used to teach NN; a location  $\boldsymbol{x}_{[\tau]}$  is sampled by a conditional probability  $P_{[\tau]}(\boldsymbol{x}_{[\tau]}|\boldsymbol{z})$  and the vector  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x}_{[\tau]}|\boldsymbol{z})$  is assigned at  $\boldsymbol{x}_{[\tau]}$  as teaching data. Typically, the condition is given by  $\boldsymbol{z} = (\boldsymbol{x}_{[0]}, \boldsymbol{x}_{[1]})$  with  $\boldsymbol{x}_{[0]} \sim P_{[0]}$  and  $\boldsymbol{x}_{[1]} \sim P_{true}$ , and we use this throughout. The vector  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x}_{[\tau]}|\boldsymbol{z})$  is made so that it satisfies the conditional continuity equation:

$$\frac{\partial}{\partial \tau} p_{[\tau]}(\boldsymbol{x}|\boldsymbol{z}) + \operatorname{div}(p_{[\tau]}(\boldsymbol{x}|\boldsymbol{z}) \boldsymbol{v}_{[\tau]}(\boldsymbol{x}|\boldsymbol{z})) = 0. \quad (1)$$

A typical construction of  $\boldsymbol{x}_{[\tau]}$  is to use a path  $\boldsymbol{x}_{[\tau]}$  ( $\tau \in [0, 1]$ ) from  $\boldsymbol{x}_{[0]}$  to  $\boldsymbol{x}_{[1]}$  and define the conditional vector by its time derivative  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x}|\boldsymbol{z}) := \frac{d}{d\tau} \boldsymbol{x}_{[\tau]}$  (see Sec. 2.2). For a deterministic path,  $P_{[\tau]}(\boldsymbol{x}_{[\tau]}|\boldsymbol{z})$  is the delta function at a point in the path  $\boldsymbol{x}_{[\tau]}$ .

Note that, given  $(\boldsymbol{x}, \tau)$ , the vector  $\boldsymbol{v}_{[\tau]}(\boldsymbol{x}|\boldsymbol{z})$  is random by the choice of  $\boldsymbol{z} = (\boldsymbol{x}_{[0]}, \boldsymbol{x}_{[1]})$ ; different vectors may be assigned to the same location  $(\boldsymbol{x}, \tau)$ . Most importantly, by averaging over  $\boldsymbol{z} \sim Q$ , where  $Q$  is the joint distribution with marginals  $P_{[0]}$  and  $P_{[1]}$ , we can see that the p.d.f. of  $\boldsymbol{x}$  at time  $\tau$ ,

$$p_{[\tau]}(\boldsymbol{x}) = \int p_{[\tau]}(\boldsymbol{x}|\boldsymbol{z}) dQ(\boldsymbol{z}), \quad (2)$$

<sup>1</sup>We use  $[\tau]$  to denote the time  $\tau \in [0, 1]$  in this section and preserve  $\boldsymbol{x}_t$  for the reverse time indexing, which is adopted from Section 4 to align with the notation of diffusion models.

and the averaged vector field  $\mathbf{v}_{[\tau]}(\mathbf{x})$  given by

$$\mathbf{v}_{[\tau]}(\mathbf{x}) := \int \mathbf{v}_{[\tau]}(\mathbf{x}|\mathbf{z})p_{[\tau]}(\mathbf{z}|\mathbf{x})d\mathbf{z}, \quad p_{[\tau]}(\mathbf{z}|\mathbf{x}) := \frac{p_{[\tau]}(\mathbf{x}|\mathbf{z})q(\mathbf{z})}{p_t(\mathbf{x})} \quad (3)$$

satisfy the continuity equation

$$\frac{\partial}{\partial \tau} p_{[\tau]}(\mathbf{x}) + \operatorname{div} (p_{[\tau]}(\mathbf{x})\mathbf{v}_{[\tau]}(\mathbf{x})) = 0. \quad (4)$$

This provides the theoretical basis for FM methods; the averaged vector field  $\mathbf{v}_{[\tau]}$  transports  $N_d(0, I_d)$  to  $P_{true}$ . The average  $\mathbf{v}_{[\tau]}$  is learned by a NN  $\phi(\mathbf{x}, \tau)$  with noisy training data  $\{(\mathbf{x}_{[\tau]}, \tau), \mathbf{v}_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{z})\}$ . Empirically, the conditional  $\mathbf{v}_{[\tau]}(\mathbf{x}|\mathbf{z})$  is given by the random sample  $\mathbf{z} = (\mathbf{x}_{[0]}, \mathbf{x}_{[1]})$  and the location  $\mathbf{x}_{[\tau]} \sim P_{[\tau]}(\mathbf{x}|\mathbf{z})$  (or path) with uniform  $\tau$ . The NN is trained with the mean square error (MSE):

$$\min_{\phi} \mathbb{E} \|\phi(\mathbf{x}_{[\tau]}, \tau) - \mathbf{v}_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{z})\|^2. \quad (5)$$

Note that  $\phi(\mathbf{x}_{[\tau]}, \tau)$  does not depend on  $\mathbf{z}$ . Since the MSE minimizer is the conditional expectation of the teaching data, the empirical minimizer  $\hat{\phi}$  is an estimator of  $\mathbf{v}_{[\tau]}(\mathbf{x})$ . Using the estimator  $\hat{\phi}$  and the corresponding flow  $\hat{\phi}_{[\tau]}$  given by ODE, we obtain the estimator  $\hat{P}_{[1]}$  for  $P_{true}$  by sampling. In practice, to reduce the variance of  $(\mathbf{x}_{[0]}, \mathbf{x}_{[1]})$  and simplify the ODE solution, the optimal transport for pairing  $\mathbf{x}_{[0]}$  and  $\mathbf{x}_{[1]}$  is applied effectively (Tong et al., 2024; Pooladian et al., 2023).

## 2.2 PATH CONSTRUCTION

This paper focuses on the following class of paths to construct the conditional vector field. Let  $\mathbf{x}_{[0]} \sim P_{[0]} = N_d(0, I_d)$ , and  $\mathbf{x}_{[1]}$  be a sample of  $P_{[1]}$  (or the empirical distribution  $\hat{P}_{train} = (1/n) \sum_{i=1}^n \delta_{\mathbf{x}^i}$  in practice). A conditional path is defined by

$$\mathbf{x}_{[\tau]} := \sigma_{[\tau]}\mathbf{x}_{[0]} + m_{[\tau]}\mathbf{x}_{[1]} \quad (0 \leq \tau \leq 1), \quad (6)$$

where  $\sigma_t$  and  $m_t$  are non-negative coefficients. We assume that  $\sigma_t$  and  $m_t$  are monotonic,  $\sigma_{[\tau]} \rightarrow 1, m_{[\tau]} \rightarrow 0$  as  $\tau \rightarrow 0^+$ , and  $\sigma_{[\tau]} \rightarrow 0, m_{[\tau]} \rightarrow 1$  as  $\tau \rightarrow 1^-$ . Let  $\sigma'_{[\tau]}$  ( $m'_{[\tau]}$ , resp.) be the time derivative of  $\sigma_{[\tau]}$  ( $m_{[\tau]}$ , resp.). With sampling  $\tau \sim \text{Unif}[0, 1]$ , a random conditional vector is assigned at  $(\mathbf{x}_{[\tau]}, \tau) \in \mathbb{R}^d \times [0, 1]$  by

$$\mathbf{v}_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{x}_{[0]}, \mathbf{x}_{[1]}) := \sigma'_{[\tau]}\mathbf{x}_{[0]} + m'_{[\tau]}\mathbf{x}_{[1]}. \quad (7)$$

Note that, due to  $\mathbf{x}_{[0]} \sim N_d(0, I_d)$ , the distribution of  $\mathbf{x}_{[\tau]}$  given  $\mathbf{x}_{[1]}$  equals  $P_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{x}_{[1]}) = N_d(m_{[\tau]}\mathbf{x}_{[1]}, \sigma_{[\tau]}^2 I_d)$ , and thus we call  $m_{[\tau]}$  and  $\sigma_{[\tau]}^2$  the mean and variance parameters, respectively.

Since (6) leads  $\mathbf{x}_{[0]} = \frac{\mathbf{x}_{[\tau]} - m_{[\tau]}\mathbf{x}_{[1]}}{\sigma_{[\tau]}}$ , the conditional vector (7) is written as

$$\mathbf{v}_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{x}_{[1]}) = \sigma'_{[\tau]} \frac{\mathbf{x}_{[\tau]} - m_{[\tau]}\mathbf{x}_{[1]}}{\sigma_{[\tau]}} + m'_{[\tau]}\mathbf{x}_{[1]}. \quad (8)$$

This class covers some popular constructions of conditional vector fields in the literature.

- Affine path: one of the most popular constructions is the following,

$$\mathbf{x}_{[\tau]} := (1 - \tau)\mathbf{x}_{[0]} + \tau\mathbf{x}_{[1]}, \quad \mathbf{v}_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{x}_{[1]}) = \mathbf{x}_{[1]} - \mathbf{x}_{[0]}.$$

This corresponds to  $m_{[\tau]} = \tau$  and  $\sigma_{[\tau]} = 1 - \tau$ . In Lipman et al. (2023)  $\mathbf{x}_{[0]}$  and  $\mathbf{x}_{[1]}$  are generated independently, while in Tong et al. (2024) they are taken by the optimal transport in a minibatch. **This case is covered by our result, which does not depend on the construction of joint distribution.**

- Diffusion: Lipman et al. (2023) presents the diffusion path, which corresponds to the deterministic probability flow (Song et al., 2020). The conditional density is given by  $p_{[\tau]}(\mathbf{x}_{[\tau]}|\mathbf{x}_{[1]} = \mathbf{y}) = N_d(m_{[\tau]}\mathbf{y}, \sigma_{[\tau]}^2 I_d)$ . The setting  $\sigma_{[\tau]}^2 = 1 - m_{[\tau]}^2$  and  $\sigma_{[\tau]} \sim \sqrt{1 - \tau}$  is typically used.

### 3 CONVERGENCE RATE OF FLOW MATCHING

We assume that the true density  $p_{[1]}$  is included in the Besov space  $B_{p',q'}^s$  ( $s > 0$ ,  $0 < p', q' \leq \infty$ ) on the cube  $[-1, 1]^d$  (we assume the unit cube, but the results can be extended to any size straightforwardly). The parameter  $s$  specifies the degree of smoothness and is most relevant in this paper. The definition of the Besov space is deferred to Appendix A.1. We use the  $r$ -Wasserstein distance  $W_r$  to measure the accuracy of the estimator. The distance  $W_r$  of the probabilities  $P_1$  and  $P_2$  on  $\mathbb{R}^d$  is defined by

$$W_r(P_1, P_2) := \left( \inf_{Q \in \Gamma(P_1, P_2)} \int \|\mathbf{x}_1 - \mathbf{x}_2\|^r dQ(\mathbf{x}_1, \mathbf{x}_2) \right)^{1/r}, \quad (9)$$

where  $\Gamma(P_1, P_2)$  denotes the joint distribution of  $(\mathbf{x}_1, \mathbf{x}_2)$  with marginals  $P_1$  and  $P_2$ . It is well known that  $W_r(P_1, P_2) \leq W_{r'}(P_1, P_2)$  holds for  $r' \geq r \geq 1$ .

As discussed in Sec. 3.1, to obtain an accurate estimator, we need to adopt early stopping of ODE and use  $\hat{P}_{[1-T_0]}$  with small  $T_0$ . Our aim is to derive a bound of  $W_p(\hat{P}_{[1-T_0]}, P_{true})$  for a large sample  $n \rightarrow \infty$ . The informal version of our main result is summarized in the following theorem.

**Theorem 1 (Informal).** *Suppose that the target probability  $P_{[1]}$  has p.d.f.  $p_{[1]}$  in the Besov space  $B_{p',q'}^s([-1, 1]^d)$  of smoothness degree  $s$ , and that  $n$  training data  $\{x^{(i)}\}_{i=1}^n$  is i.i.d. samples from  $P_{[1]}$ . Assume that  $\sigma_{[\tau]} \sim (1 - \tau)^\kappa$  ( $\tau \rightarrow 1^-$ ) with  $\kappa \geq 1/2$ , the conditional vector field is given by (6) and (7), and that time-divided neural networks are used (see Sec. 4.4). Then, under several assumptions, the FM estimator  $\hat{P}_{[1-T_0]}$  with  $T_0 = n^{-R_0}$  with appropriate  $R_0$  satisfies, for any  $\delta > 0$ ,*

$$\mathbb{E}[W_2(\hat{P}_{[1-T_0]}, P_{true})] = O\left(n^{-\frac{s+(2\kappa)^{-1}-\delta}{2s+d}}\right) \quad (n \rightarrow \infty), \quad (10)$$

where  $\mathbb{E}$  denotes the expectation over the training data.

It is known that a lower bound of the minimax convergence rate exists for the Wasserstein distance for probability estimation. We use the notation  $\gtrsim$  to mean the lower bound up to a constant factor.

**Proposition 2 (Niles-Weed and Berthet (2022)).** *Let  $p', q' \geq 1$ ,  $s > 0$ ,  $r \geq 1$ , and  $d \geq 2$ . Then,*

$$\inf_{\hat{P}} \sup_{p \in B_{p',q'}^s([-1, 1]^d)} \mathbb{E}[W_r(\hat{P}, P)] \gtrsim n^{-\frac{s+1}{2s+d}} \quad (n \rightarrow \infty),$$

where  $\hat{P}$  runs over all estimators based on  $n$  i.i.d. samples from  $P$ .

For  $\kappa = 1/2$ , by Theorem 1 and Proposition 2, the upper bound  $n^{-\frac{s+1-\delta}{2s+d}}$  is almost the optimal convergence rate up to an arbitrarily small  $\delta > 0$ . In addition, this convergence rate coincides with that of the diffusion model given in Oko et al. (2023) for  $W_1$ . The above result indicates that the flow matching is as good as the diffusion model regarding the minimax convergence rate under  $W_1$ , where the max in minimax means sup over the Besov space.

#### 3.1 KERNEL DENSITY ESTIMATION AND EARLY STOPPING OF ODE

In practice, with conditional density  $P_{[\tau]}(\mathbf{x}|\mathbf{x}_{[1]}) = N_d(m_{[\tau]}\mathbf{x}_{[1]}, \sigma_{[\tau]}^2 I_d)$ , the parameter  $\sigma_{[1]}$  is often set as a small positive value  $\sigma_{[1]} = \sigma_{min} > 0$  so that (7) is well defined up to  $\tau = 1$  (e.g. Lipman et al., 2023). If  $\mathbf{x}_{[1]}$  is sampled from  $\hat{P}_{train} = \frac{1}{n} \sum_{j=1}^n \delta_{\mathbf{x}^j}$ , the obtained distribution equals to

$$\hat{p}_{[1]}(\mathbf{x}) = \int p_{[1]}(\mathbf{x}|\mathbf{x}_{[1]}) d\hat{P}_{train}(\mathbf{x}_{[1]}) = \frac{1}{n} \sum_{j=1}^n \frac{1}{(2\pi\sigma_{min}^2)^{d/2}} \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}^j\|^2}{2\sigma_{min}^2}\right),$$

which is exactly the kernel density estimator (KDE) with the Gaussian kernel of bandwidth  $\sigma_{min}$ . If the ODE is solved up to  $\tau = 1$  rigorously, the pushforward realizes this KDE. As is well known (Scott, 1992), the convergence rate of this KDE under MSE is  $O(n^{-4/(4+d)})$  at best by choosing the optimal  $\sigma_{min}$  depending on  $n$ , which is much slower than the optimal rate  $n^{-2s/(2s+d)}$  under MSE for the true density in  $B_{p',q'}^s(I^d)$  (Liu et al., 2023a). Based on this consideration, we discuss the early stopping of the ODE, where we stop at  $\tau = 1 - T_0$  with small  $T_0 > 0$  and consider the convergence rate of the estimator  $\hat{p}_{[1-T_0]}$ . Notice that  $\hat{p}_{[1-T_0]}$  differs from KDE, since it is given by the pushforward of the trained vector field. For diffusion models, Oko et al. (2023) and Zhang et al. (2024) also discuss the estimator obtained by stopping the reverse SDE at  $T_0 > 0$  to derive the convergence rate.

### 216 3.2 RELATED WORKS

217  
218 Among many literatures on the statistical convergence of diffusion models, the most relevant to this  
219 work is Oko et al. (2023). Although our analysis is based on Oko et al. (2023) and derives comparable  
220 results, there are significant differences. First, we analyze the more general settings for  $m_{[\tau]}$  and  $\sigma_{[\tau]}$   
221 in the conditional distribution  $P_{[\tau]}(\mathbf{x}|\mathbf{y}) = N_d(m_{[\tau]}\mathbf{y}, \sigma_{[\tau]}^2 I)$ ; Oko et al. (2023) considers only the  
222 case of  $\sigma_t \sim \sqrt{t}$  and  $m_t \sim 1 - t$  (in reverse time  $t$ ), which is a typical choice for diffusion models.  
223 Consequently, we have shown that for  $\sigma_{[\tau]} \sim (1 - \tau)^\kappa$  with  $\kappa \geq 1/2$ , only  $\kappa = 1/2$  achieves the  
224 almost minimax optimal convergence rate. Second, due to the difference between the ODE and  
225 diffusion processes, the proof technique for relating the Wasserstein metric and the  $L_2$ -risk is very  
226 different. Our technique is based on Alekseev-Gröbner lemma to derive the bound for  $r$ -Wasserstein  
227 with  $1 \leq r \leq 2$ , while Oko et al. (2023) obtained the bound only for 1-Wasserstein. Third, this is  
228 the first theoretical result for FM showing a convergence rate that is almost optimal. Although FM  
229 has been recently used in many applications with competitive results to diffusion models, theoretical  
230 comparisons in terms of convergence rates have been lacking. The results of this paper show that **both**  
231 **FM and DM can attain the same almost minimax optimal convergence rate for generalization error.**

232 For FM, there are some recent works on convergence. Albergo and Vanden-Eijnden (2023) and  
233 Benton et al. (2023b) relate the Wasserstein distance to the  $L_2$ -risk of the vector fields and show  
234 convergence for a large sample size, but did not derive a convergence rate. Jiao et al. (2024) discusses  
235 convergence rates of FM applied in the latent space of the autoencoder and considers the discretization  
236 effect of the numerical ODE solution in their analysis. However, they did not include the degree of  
237 smoothness in developing the convergence rate. Albergo et al. (2023b) present a unifying view of the  
238 theory of diffusion models and FM with the upper bounds of discrepancy measures.

## 239 4 THEORETICAL DETAILS

240  
241 This section rigorously presents the main result with the assumptions and shows the proof outline. In  
242 the sequel, we use *reverse time index*  $t = 1 - \tau$  ( $\tau \in [0, 1]$ );  $t = 0$  for  $P_{true}$  and  $t = 1$  for  $N_d(0, I_d)$ ,  
243 which align with the notations of the diffusion models. We use  $\text{poly}(\log n)$  to indicate the term of  
244  $O(\log^r n)$ -order with some natural number  $r$ , and  $\tilde{O}(n^\alpha)$  to mean the order up to  $\text{poly}(\log n)$  factor.

### 245 4.1 PROBLEM SETTING

246  
247 With reverse time  $t$ , the definitions (7), (2), and (3) are modified by replacing  $[1 - t]$  with  $t$ ;

$$248 P_t = P_{[1-t]}, \quad P_0 = P_{true}, \quad P_1 = N_d(0, I_d).$$

249 The flows  $\varphi_t$  and  $\hat{\varphi}_t$  are defined by solving the ODE from  $t = 1$  in the reverse time direction:

$$250 \frac{d}{dt}\varphi_t(\mathbf{x}) = \mathbf{v}_t(\varphi_t(\mathbf{x})), \quad \frac{d}{dt}\hat{\varphi}_t(\mathbf{x}) = \hat{\mathbf{v}}_t(\varphi_t(\mathbf{x})), \quad (11)$$

251 where  $\mathbf{v}_t(\mathbf{x})$  and  $\hat{\mathbf{v}}_t(\mathbf{x})$  are the vector field (3) and its neural estimate, respectively. The distributions  
252 at  $t \in [0, 1]$  are given by

$$253 P_t = (\varphi_t)_\# P_1, \quad \hat{P}_t = (\hat{\varphi}_t)_\# P_1, \quad (12)$$

254 where  $(\varphi_t)_\#$  and  $(\hat{\varphi}_t)_\#$  denote the pushforward by the respective flows  $\varphi_t$  and  $\hat{\varphi}_t$ .

### 255 4.2 ASSUMPTIONS

256  
257 In the remainder of this paper,  $\delta > 0$  is an arbitrarily small positive value. As in Oko et al. (2023),  
258 we introduce  $N$  to specify the number of basis functions of the  $B$ -spline for approximating  $p_t(\mathbf{x})$   
259 and  $\mathbf{v}_t(\mathbf{x})$ . This number  $N$  depends on the sample size  $n$  ( $N = n^{\frac{d}{2s+d}}$  is used), balancing the  
260 approximation error and complexity of the  $B$ -spline and NN. We set the stopping time  $T_0 = N^{-R_0}$   
261 as discussed in Sec. 3.1 ( $R_0$  is specified later), and solve the ODE from 1 to  $T_0$ . For simplicity, the  $d$   
262 dimensional cube  $[-1, 1]^d$  and the reduced cube  $[-1 + N^{-(1-\kappa\delta)}, 1 - N^{-(1-\kappa\delta)}]^d$  are denoted by  $I^d$   
263 and  $I_N^d$ , respectively, where  $\kappa > 0$  is specified below in (A3). We make the following assumptions.

264 (A1) The target probability  $P_0$  has support  $I^d$  and its p.d.f.  $p_0$  satisfies  $p_0 \in B_{p',q'}^s(I^d)$  and  $p_0 \in$   
265  $B_{p',q'}^{\tilde{s}}(I^d \setminus I_N^d)$  with  $\tilde{s} > \max\{6s - 1, 1\}$ .

(A2) There exists  $C_0 > 0$  such that  $C_0^{-1} \leq p_0(\mathbf{x}) \leq C_0$  for all  $\mathbf{x} \in I^d$ .

(A3) There is  $\kappa \geq 1/2$ ,  $b_0 > 0$ ,  $\tilde{\kappa} > 0$ , and  $\tilde{b}_0 > 0$  such that

$$\sigma_t = b_0 t^\kappa, \quad 1 - m_t = \tilde{b}_0 t^{\tilde{\kappa}}$$

for sufficiently small  $t \geq T_0$ . Also, there are  $D_0 > 0$  and  $K_0 > 0$  such that

$$D_0^{-1} \leq \sigma_t^2 + m_t^2 \leq D_0, \quad |\sigma_t'| + |m_t'| \leq N^{K_0} \quad (\forall t \in [T_0, 1]).$$

(A4) If  $\kappa = 1/2$ , there is  $b_1 > 0$  and  $D_1 > 0$  such that for any  $0 \leq \gamma < R_0$

$$\int_{T_0}^{N^{-\gamma}} \{(\sigma_t')^2 + (m_t')^2\} dt \leq D_1 (\log N)^{b_1}.$$

(A5) There is a constant  $C_L > 0$  such that  $\|\frac{\partial}{\partial \mathbf{x}} \int \mathbf{y} p_t(\mathbf{y}|\mathbf{x}) d\mathbf{y}\|_{op} \leq C_L$  for any  $t \in [T_0, 1]$ .

The higher degree of smoothness is assumed around the boundary of  $I^d$  in (A1) for a technical reason to compensate for the nondifferentiability of  $p_0(x)$  at the boundary by (A2). In (A3), it may be more natural to require  $\sigma_t^2 + m_t^2 = 1$  so that signal power can be maintained. However, in this paper, to pursue the flexibility of choosing  $\sigma_t$  and  $m_t$ , we allow bounded changes of  $\sigma_t^2 + m_t^2$  over  $t$ . (A4) is required to limit the complexity of the neural network model (see Lemma 5). In (A3),  $\kappa$  is assumed to be not less than  $1/2$ , because for  $\kappa < 1/2$ , the integral  $\int_{T_0} (\sigma_t')^2 dt$  with  $T_0 = N^{-R_0}$  diverges to infinity as  $N \rightarrow +\infty$ , which causes the divergence of the complexity bound in Lemma 5. Note that the boundary case  $\kappa = 1/2$  is, in fact, popularly used for the diffusion model. In this case,  $(\sigma_t')^2$  is the order  $1/t$  for  $t \rightarrow 0^+$  and the integral from  $T_0$  is of the order  $\log N$ , which still diverges to infinity as  $n \rightarrow \infty$ . As discussed in Section 4.4, we consider this integral only for a short time interval, and we will see that the  $W_2$  distance converges to zero as  $n \rightarrow \infty$ . (A5) is made to bound the Lipschitz factor in Theorem 3 under (A3) (see Lemma 10).

### 4.3 GENERALIZATION BOUND

It is known (Albergo and Vanden-Eijnden, 2023; Benton et al., 2023b) that, given two vector fields, the  $W_2$ -distance of the pushforwards of the same distribution by the corresponding flows admits an upper bound by the  $L_2$ -risk of the vector fields;

**Theorem 3.** *Let  $\mathbf{v}_t(\mathbf{x})$  and  $\hat{\mathbf{v}}_t(\mathbf{x})$  be vector fields such that  $\mathbf{x} \mapsto \hat{\mathbf{v}}_t(\mathbf{x})$  is  $L_t$ -Lipschitz for each  $t$ , and  $P_t$  and  $\hat{P}_t$  be the pushforwards of distribution  $P_0$  by the corresponding flows at time  $t$  from  $t = 0$ . Then, for any  $t \in [0, 1]$ , we have*

$$W_2(\hat{P}_t, P_t) \leq \sqrt{t} \left( \int_0^t \int e^{2 \int_s^t e^{L_u} du} \|\hat{\mathbf{v}}_s(\mathbf{x}) - \mathbf{v}_s(\mathbf{x})\|^2 dP_s(\mathbf{x}) d\mathbf{x} ds \right)^{1/2}. \quad (13)$$

See Appendix B for the proof. From Theorem 3, we can consider the  $L_2$ -error  $\mathbb{E}[\int \int \|\hat{\mathbf{v}}(\mathbf{x}, s) - \mathbf{v}(\mathbf{x}, s)\|^2 dP_s(\mathbf{x}) d\mathbf{x} ds]$  of the vector field to obtain the bound of the  $W_2$  distance of the distributions. From the fact  $W_r \leq W_{r'}$  ( $1 \leq r \leq r'$ ), the same upper bound holds for  $W_r$  for  $1 \leq r \leq 2$ .

We first review a general method for bounding the generalization. We consider training within the general time interval  $[T_\ell, T_u]$  where  $T_0 \leq T_\ell < T_u \leq 1$ . For an estimator  $\phi(\mathbf{x}, t)$  of the true vector field  $\mathbf{v}_t(\mathbf{x})$ , we define the loss function  $\ell_\phi^{T_\ell, T_u}(\mathbf{x})$  for  $\mathbf{x} \in I^d$  by

$$\ell_\phi^{T_\ell, T_u}(\mathbf{x}) := \int_{T_\ell}^{T_u} \int \|\phi(\mathbf{x}_t, t) - \mathbf{v}_t(\mathbf{x}_t|\mathbf{x})\|^2 p_t(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}_t dt, \quad (14)$$

where  $\mathbf{x}$  is the condition of  $\mathbf{v}_t(\mathbf{x}_t|\mathbf{x})$ . Although the definition depends on  $T_\ell$  and  $T_u$ , we omit them when there is no confusion. Given the training data  $\{\mathbf{x}^i\}_{i=1}^n$ , the vector field is trained with the teaching data  $\mathbf{v}_t(\mathbf{x}_t|\mathbf{x}^i)$  at the location  $(\mathbf{x}_t, t)$  ( $t \in [T_\ell, T_u]$ ), which is sampled from  $p_t(\mathbf{x}_t|\mathbf{x}^i)$  and the uniform distribution  $U([T_\ell, T_u])$ . Note that given  $\mathbf{x}^i$ , we can generate any number of  $(\mathbf{x}_t, t)$ . Thus, the sampling error in (14) is negligible and the training by a NN can be regarded as minimization of

$$\frac{1}{n} \sum_{i=1}^n \ell_\phi(\mathbf{x}^i). \quad (15)$$

See Oko et al. (2023, Section 4) for the discussion of the effect of sampling.

Let  $\widehat{\phi}$  be the minimizer of (15) among the function class  $\mathcal{S}$ . The generalization error is then given by

$$\mathcal{E}_{gen} := \mathbb{E} \left[ \int \int \int_{T_\ell}^{T_u} \|\widehat{\phi}(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}|\mathbf{y}) dt d\mathbf{x} p_0(\mathbf{y}) d\mathbf{y} \right] = \mathbb{E} \left[ \int \ell_{\widehat{\phi}}(\mathbf{y}) p_0(\mathbf{y}) d\mathbf{y} \right]. \quad (16)$$

Let  $\mathcal{L} := \{\ell_\phi \mid \phi \in \mathcal{S}\}$  and  $\mathcal{N}(\mathcal{L}, \|\cdot\|_{L^\infty(I^d)}, \varepsilon)$  be the covering number of the function class  $\mathcal{L}$  with the  $\|\cdot\|_{L^\infty(I^d)}$ -norm. Then, a standard argument on the generalization error analysis derives the following upper bound (see Oko et al. (Theorem C.4, 2023) and also Hayakawa and Suzuki (2020)).

**Theorem 4.** *The generalization error of the minimizer of (15) among  $\phi \in \mathcal{S}$  is upper bounded by*

$$\begin{aligned} \mathcal{E}_{gen} \leq & 2 \inf_{\phi \in \mathcal{S}} \int \int_{T_\ell}^{T_u} \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) dt d\mathbf{x} \\ & + \frac{\sup_{\phi \in \mathcal{S}} \|\ell_\phi\|_{L^\infty(I^d)}}{n} \left( \frac{37}{9} \log \mathcal{N}(\mathcal{L}, \|\cdot\|_{L^\infty(I^d)}, \varepsilon) + 32 \right) + 3\varepsilon. \end{aligned} \quad (17)$$

From Theorems 3 and 4, it suffices to consider the approximation error (1st term) and complexity (2nd term) in (17) for deriving the  $W_2$  distributional bound.

#### 4.3.1 COMPLEXITY TERM IN GENERALIZATION BOUND

We first consider the complexity term, where the class  $\mathcal{S}$  is given by NN. A class of NN  $\mathcal{M}(L, W, S, B)$  with height  $L$ , width  $W$ , sparsity constraint  $S$ , and norm constraint  $B$  is defined as

$$\begin{aligned} \mathcal{M}(L, W, S, B) := & \{\psi_{A^{(L)}, b^{(L)}} \circ \dots \circ \psi_{A^{(2)}, b^{(2)}}(A^{(1)}\mathbf{x} + b^{(1)}) \mid A^{(i)} \in \mathbb{R}^{W_{i+1} \times W_i}, b^{(i)} \in \mathbb{R}^{W_{i+1}}, \\ & \sum_{i=1}^L (\|A^{(i)}\|_0 + \|b^{(i)}\|_0) \leq S, \max_i \|A^{(i)}\|_\infty \vee \|b^{(i)}\|_\infty \leq B\}, \end{aligned}$$

where  $\psi_{A,b}(z) = A \text{ReLU}(z) + b$ . As shown in Theorems 7 and 8 later, it suffices to consider the NNs that satisfy

$$\|\phi(\mathbf{x}, t)\|_\infty \leq D(|\sigma'_t| \sqrt{\log n} + |m'_t|)$$

for some constant  $D$ . Also, we can see in Lemma A.2) that  $\mathbf{x} \mapsto \mathbf{v}_t(\mathbf{x})$  is Lipschitz continuous with Lipschitz constant proportional to  $1/t$  under (A3) and (A5). Reflecting these facts, we define the following NN class for training the vector field:

$$\begin{aligned} \mathcal{H}_n := & \{\phi \in \mathcal{M}(L, W, S, B) \mid \|\phi(\cdot, t)\|_\infty \leq D(|\sigma'_t| \sqrt{\log n} + |m'_t|) \text{ for } \forall t \in [T_0, 1], \\ & \mathbf{x} \mapsto \phi(\mathbf{x}, t) \text{ is } L_t\text{-Lipschitz for each } t \in [T_0, 1] \text{ where } L_t = \tilde{C}_L/t\}, \end{aligned} \quad (18)$$

where  $D$  and  $\tilde{C}_L$  are some positive constants.

The supremum norm and the covering number in Theorem 4 are given in the following lemmas.

**Lemma 5.** *Let  $T_0 \leq T_\ell < T_u \leq 1$ . Under Assumption (A4), there is  $C_s > 0$  such that*

$$\sup_{\phi \in \mathcal{H}_n} \|\ell_\phi\|_{L^\infty(I^d)} \leq C_s (\log n)^{b+1}, \quad (19)$$

where  $b = b_1$  in (A4) for  $\kappa = 1/2$ , and  $b = 0$  for  $\kappa > 1/2$ .

See Appendix C.1 for the proof. To obtain this bound, we need to impose the upper bound of  $\phi$  as in (18). In practice, the vectors in the teaching data satisfy this upper bound, and thus  $\phi$  will naturally satisfy the same bound by the least square error solution. The following bound of the covering number for neural networks is given by Suzuki (Lemma 3, 2019).

**Lemma 6.** *For the function class  $\mathcal{H}_n$ , the covering number satisfies*

$$\log \mathcal{N}(\mathcal{L}, \|\cdot\|_{L^\infty(I^d)}, \varepsilon) \leq SL \log(\varepsilon^{-1} \|W\|_\infty B n).$$

#### 4.3.2 APPROXIMATION ERROR FOR SMALL $t$

Recall that  $N$  specifies the number of basis functions of the  $B$ -spline for the approximation. We derive upper bounds of the approximation error of the NN model  $\mathcal{M}(L, W, S, B)$ , where  $L, W, S,$

and  $B$  are specified in terms of  $N$ . We will separate  $[T_0, 1]$  into two intervals,  $[T_0, 3T_*]$  and  $[T_*, 1]$ , where  $T_* := N^{-(\kappa^{-1}-\delta)/d}$ , and provide different upper bounds. The reason for this choice of division point  $T_*$  is sketched as follows and is detailed in C.2. In the approximation of the vector field, we use the  $B$ -spline approximation of densities as in Oko et al. (2023). To show a fast convergence rate, the first interval is more subtle because  $p_t(x)$  is rougher. In approximating the density on the smoother boundary region, we divide the region into small cubes, each of which uses  $N^\delta$  bases for  $B$ -spline approximation. To make the total number of  $B$ -spline bases comparable with  $N$ , the width  $a_0$  of the region should be  $a_0 = N^{(1-\kappa\delta)/d}$ . On the other hand, in Theorem 7, we need a concentration of an integral around the boundary region for a better approximation by the higher smoothness, and this limits the variance of the Gaussian kernel so that  $\sigma_t = t^\kappa \leq a_0$ . This derives  $t \leq N^{-(\kappa^{-1}-\delta)/d}$ . As a result, the division point is small enough as  $T_* := N^{-(\kappa^{-1}-\delta)/d}$ .

The approximation bound for  $t \in [T_0, 3T_*]$  with  $T_* := N^{-\frac{\kappa^{-1}-\delta}{d}}$  is given in the following Theorem.

**Theorem 7.** *Under assumptions (A1)-(A5), there is a neural network  $\phi_1 \in \mathcal{M}(L, W, S, B)$  and a constant  $C_6$ , which is independent of  $t$ , such that, for sufficiently large  $N$ ,*

$$\int \|\phi_1(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \leq C_6 \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\frac{2s}{d}}, \quad (20)$$

for any  $t \in [T_0, 3T_*]$ , where

$$L = O(\log^4 N), \|W\|_\infty = O(N \log^6 N), S = O(N \log^8 N), B = \exp(O(\log N \log \log N)).$$

Additionally, we can take  $\phi_1$  to satisfy

$$\|\phi_1(\mathbf{x}, t)\| \leq \tilde{C}_6 \{(\sigma'_t) \sqrt{\log n} + |m'_t|\},$$

where  $\tilde{C}_6$  is a constant independent of  $t$ .

See Appendix C.4 for the proof.

#### 4.3.3 APPROXIMATION ERROR FOR LARGE $t$

We derive a bound of the approximation error on any interval  $[2t_*, 1]$ , where  $t_* \geq T_* = N^{-\frac{\kappa^{-1}-\delta}{d}}$ . This is used to discuss the optimal convergence rate in Section 4.4.

**Theorem 8.** *Fix  $t_* \in [T_*, 1]$  and take an arbitrary  $\eta > 0$ . Under the assumptions (A1)-(A5), there is a neural network  $\phi_2 \in \mathcal{M}(L, W, S, B)$  and  $C_7 > 0$ , which does not depend on  $t$ , such that the bound*

$$\int \|\phi_2(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \leq C_7 \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\eta} \quad (21)$$

holds for all  $t \in [2t_*, 1]$ , and the NN model satisfies  $L = O(\log^4 N)$ ,  $\|W\|_\infty = O(N)$ ,  $S = O(t_*^{-d\kappa} N^{\delta\kappa})$ , and  $B = \exp(O(\log N \log \log N))$ . Moreover,  $\phi_2$  can be taken so that

$$\|\phi_2(\cdot, t)\|_\infty \leq \tilde{C}_7 \{(\sigma'_t) \log N + |m'_t|\}$$

with constant  $\tilde{C}_7 > 0$  independent of  $t$ .

See Appendix C.5 for the proof. The approximation error  $N^{-\eta}$  is arbitrarily small and is not dominant, while  $S$  may be dominant in the complexity term.

## 4.4 CONVERGENCE RATE UNDER WASSERSTEIN DISTANCE

We can consider a generalization bound based on Theorems 3, 4, 7, and 8 deriving the bounds for  $[T_0, 2T_*]$  and  $[2T_*, 1]$ . However, if we apply Theorem 8 for  $[2t_*, 1]$  with  $t_* = T_* = N^{-(\kappa^{-1}-\delta)/d}$ , the dominant factor of the log covering number in (17) is the sparsity  $S = O(t_*^{-d\kappa} N^{\delta\kappa}) = O(N)$ . From Theorems 3 and 4, the complexity part gives  $O((N/n)^{1/2})$  term in the  $W_2$  generalization error. If we plug  $N = n^{(2s+d)/d}$ , which is optimal for the MSE generalization, we have  $O(n^{-s/(2s+d)})$  for the upper bound of  $W_2$  generalization, which is slower than the lower bound in Proposition 2. To achieve a better convergence rate, we will make use of the factor  $\sqrt{t}$  in front of (13) by dividing the



interval  $[T_0, 1]$  into small pieces and using a NN for each small interval, as in Oko et al. (2023) for diffusion models.

Notice that, when time  $t$  is far from 0, the convolution of  $p_t(\mathbf{x}|\mathbf{y})$  with larger  $\sigma_t$  results in a smoother target vector field  $\mathbf{v}_t(\mathbf{x})$ , which is easier to approximate with a low-complexity model. On the other hand, when  $t$  approaches 0, with the fixed number of  $B$ -spline bases  $N$ , the approximation error bound  $\{(\sigma'_t)^2 \sqrt{\log N} + (m'_t)^2\} N^{-2s/d}$  can increase for large  $\sigma'_t$  or  $m'_t$  (e.g.  $\sigma_t \sim t^\kappa$  with  $\kappa < 1$ ). We therefore need a more complex model (that is, larger  $N$ ) than the one needed for larger  $t$ . Thus, it will be more efficient if the number of  $B$ -spline bases  $N$ , which controls the approximation error and complexity, is chosen adaptively to the time region  $t$ .

Specifically, we make a partition  $T_0 = t_0 < t_1 < t_2 < \dots < t_K = 1$  such that  $t_j = 2t_{j-1}$  for  $1 \leq j \leq K-1$  with  $2t_{K-1} \geq 1$ , and build a neural network for each  $[t_{j-1}, t_j]$  ( $j = 1, \dots, K$ ). Note that we train each network for interval  $[t_{j-1}, t_j]$  with  $n$  training data  $(\mathbf{x}_i)_{i=1}^n$ . We assume that  $t_{j_*}$  with  $T_* \leq t_{j_*} \leq 3T_*$  serves as the boundary to apply two different error bounds. The total number of intervals  $K$  is  $O(\log N) = O(\log n)$ , since  $2^K T_0 \geq 1$  with  $T_0 = N^{-R_0}$  can be achieved by  $K \geq R_0 \log_2 N$ . The constant  $R_0$  is fixed as  $R_0 \geq \frac{s+1}{\min(\kappa, \bar{\kappa})}$  so that  $W_2(P_{T_0}, P_0)$  is negligible (see the proof sketch of Theorem 9). In this setting, we have the following main result.

**Theorem 9** (Main result). *Assume (A1)-(A5) and  $d \geq 2$ . If the above time-partition is applied and a neural network is trained for each time division, for arbitrarily small  $\delta > 0$  and  $1 \leq r \leq 2$ , we have*

$$\mathbb{E}[W_r(\hat{P}_{T_0}, P_{true})] = O\left(n^{-\frac{s+(2\kappa)^{-1}-\delta}{2s+d}}\right) \quad (n \rightarrow \infty). \quad (22)$$

*Proof Sketch.* Let  $J_j := [t_{j-1}, t_j]$  ( $j = 1, \dots, K$ ). We use a smaller neural network model for larger  $j$ . Specifically, the number of  $B$ -spline bases for  $J_j$  is  $N'_j := t_{j-1}^{-d\kappa} N^{\delta\kappa}$  for  $j > j_*$ , while  $N'_j = N$  for  $j \leq j_*$ , where  $j_*$  is defined as above. Note that  $N'_j \rightarrow \infty$  as  $N \rightarrow \infty$  due to  $\delta > 0$ , and that  $N'_j \leq N^{1-\delta\kappa} N^{\delta\kappa} = N$  for  $j \geq j_*$  due to  $t_j \geq N^{-\frac{\kappa-1}{d}-\delta}$ , which means a lower complexity. See also Figure E.1 in Appendix.

Next, we consider the bound of  $W_2$ -distance based on the partition. For each of  $j = 1, \dots, K$ , we introduce a vector field  $\tilde{\mathbf{v}}_t^{(j)}$  such that it coincides with the target  $\mathbf{v}_t$  for  $t \in [t_j, 1]$  and with the learned  $\hat{\mathbf{v}}_t$  for  $t \in [T_0, t_j]$ . Let  $Q^{(j)}$  be the pushforward from  $P_1 = N_d(0, I_d)$  to  $t = T_0$  by the flow of the vector field  $\tilde{\mathbf{v}}_t^{(j)}$ . Then,  $Q^{(0)} = P_{T_0}$ , the pushforward by the flow of the target  $\mathbf{v}_t$  from  $t = 1$  to  $T_0$ , and also  $Q^{(K)} = \hat{P}_{T_0}$ . Note also that  $\tilde{\mathbf{v}}^{(j)}$  and  $\tilde{\mathbf{v}}^{(j-1)}$  differ only in  $J_j$  by  $\mathbf{v}_t(\mathbf{x}) - \hat{\mathbf{v}}_t(\mathbf{x})$ . Therefore, from Theorem 3 and the Lipschitz assumption on  $\mathcal{H}_n$ , we have

$$\begin{aligned} W_2(P_0, \hat{P}_{T_0}) &\leq W_2(P_0, P_{T_0}) + \sum_{j=1}^K W_2(Q^{(j-1)}, Q^{(j)}) \\ &\leq \theta_n + C \sum_{j=1}^K \sqrt{t_j} \left\{ \int_{t_{j-1}}^{t_j} e^{2 \int_t^{t_j} (\tilde{C}/u) du} \int \|\hat{\phi}(\mathbf{x}, t) - \mathbf{v}(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} dt \right\}^{1/2}, \end{aligned}$$

where  $\theta_n^2 = db_0^2 n^{-\frac{2R_0\kappa}{2s+d}} + \int \|y\|^2 dP_0(\mathbf{y}) \tilde{b}_0^2 n^{-\frac{2R_0\bar{\kappa}}{2s+d}}$ , which is derived from Lemma 11 and (A3). We take a constant  $R_0 \geq (s+1)/\min(\kappa, \bar{\kappa})$  so that  $\theta_n$  is of  $O(n^{-\frac{s+1}{2s+d}})$  and thus  $\theta_n$  is negligible. It is easy to see that the factor  $e^{\int_t^{t_j} (\tilde{C}/u) du}$  is bounded by a constant because of  $t_j = 2t_{j-1}$  by definition.

For simplicity, let  $t_* := t_{j_*}$ . From Theorems 7, 8, and 4, the generalization bound of  $\hat{P}_{T_0}$  is given by

$$\begin{aligned} &\mathbb{E} \left[ W_2(P_0, \hat{P}_{T_0}) \right] \\ &\leq \theta_n + \sum_{j=1}^{j_*} \sqrt{t_*} \left\{ C_6 \int_{t_{j-1}}^{t_j} \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-2s/d} dt + \frac{N}{n} O(\text{poly}(\log n)) \right\}^{1/2} \\ &\quad + \sum_{j=k_*}^K \sqrt{t_j} \left\{ C_7 \int_{t_{j-1}}^{t_j} \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\eta} dt + \frac{t_j^{-d\kappa} N^{\delta\kappa}}{n} O(\text{poly}(\log n)) \right\}^{1/2} \\ &\leq \theta_n + C''' \sqrt{t_*} t_*^{\kappa-1/2} N^{-s/d} O(\text{poly}(\log n)) + C''' \sqrt{\frac{N}{n}} O(\text{poly}(\log n)) \\ &\quad + C''' \sum_{j=k_*}^K \left\{ \sqrt{t_j} N^{-\eta/2} O(\text{poly}(\log n)) + \sqrt{t_j} \frac{t_j^{-d\kappa/2} N^{\delta\kappa/2}}{\sqrt{n}} O(\text{poly}(\log n)) \right\} \end{aligned}$$

$$\begin{aligned} &\leq \theta_n + C''' t_*^\kappa n^{-s/(2s+d)} O(\text{poly}(\log n)) + C''' \sqrt{t_*} n^{-s/(2s+d)} O(\text{poly}(\log n)) \\ &\quad + C''' \sum_{j=k_*}^K \left\{ \sqrt{t_j} N^{-\eta/2} O(\text{poly}(\log n)) + t_j^{\frac{1-d\kappa}{2}} n^{\frac{d\delta\kappa}{2(2s+d)} - \frac{1}{2}} O(\text{poly}(\log n)) \right\}. \end{aligned}$$

In the third inequality, we use  $\int_{t_{j-1}}^{t_j} \{(\sigma'_t)^2 + (m'_t)^2\} dt = O(\text{poly}(\log n))$  for  $\kappa = 1/2$ , and the fact that it is bounded by a constant for  $\kappa > 1/2$ . Since  $\eta$  is arbitrarily large and  $\kappa \geq 1/2$ , neglecting the factors of  $\text{poly}(\log n)$ , the candidates of the dominant terms in the above expression are  $t_*^{1/2} n^{-\frac{s}{2s+d}}$  in the third term and  $t_*^{\frac{1-d\kappa}{2}} n^{\frac{d\delta\kappa}{2(2s+d)} - \frac{1}{2}}$  in the last summation. By balancing these two terms, the upper bound can be minimized by setting

$$t_* = C_* n^{-\frac{\kappa-1-\delta}{2s+d}}, \quad (23)$$

for some contact  $C_*$ , and the dominant term of the upper-bound is given by

$$\tilde{O}\left(n^{-\frac{s+(2\kappa)^{-1}-\delta/2}{2s+d}}\right). \quad (24)$$

This proves the claim by replacing  $\delta/2$  with  $\delta$  and absorbing the  $\text{poly}(\log)$  factor in  $\delta$ .  $\square$

From Proposition 2 and Theorem 9, if  $\kappa = 1/2$ , the FM method achieves an almost optimal rate up to the  $\text{poly}(\log n)$  factor and arbitrary small  $\delta > 0$ . On the other hand, for  $\kappa > 1/2$ , the obtained upper bound is not optimal. This suggests that the choice of  $\sigma_t \sim \sqrt{t}$  around  $t \rightarrow 0^+$  is theoretically reasonable. This is also a popular choice for the diffusion model.

#### 4.5 DISCUSSION

In the derivation of the almost minimax optimal convergence rate, the use of neural networks for divided time intervals is a limitation of the current theoretical analysis. As discussed before Theorem 9, without this partition, the current analysis gives only  $\tilde{O}(n^{-\frac{s}{2s+d}})$ , which is not optimal for  $W_2$ . It is obviously an important question of how to avoid such a time division. In Oko et al. (2023), the optimal convergence rate for the diffusion model has been proved without time division for the total variation, which is  $\tilde{O}(n^{-\frac{s}{2s+d}})$ . The bound is based on Girsanov's theorem, which gives an upper bound of the KL divergence of SDE by the  $L_2$  losses of the drift estimation. To the best of our knowledge, no bounds are known for ODE with respect to the KL or total variation for the difference of vector fields. This is an important future direction for understanding the ability of FM theoretically.

## 5 CONCLUSION

This paper has rigorously analyzed the convergence rate of flow matching, demonstrating for the first time that FM can achieve the almost minimax optimal convergence rate under the 2-Wasserstein distance. This result positions FM as a competitive alternative to diffusion models in terms of asymptotic convergence rates, which concurs with empirical results in various applications. Our findings further reveal that the convergence rate is significantly influenced by the variance decay rate in the Gaussian conditional kernel, where  $\sigma_t \sim \sqrt{t}$  is shown to yield the optimal rate. Although there are several popular proposals for the mean and variance functions, theoretical justification or comparison has not been explored intensively. The current result on the upper bound (Theorem 9) provides theoretical insight on the influence of the choice of these functions.

Although this study offers substantial theoretical contributions, these insights are still grounded in specific modeling assumptions that limit broader applicability. In addition to the time-partition discussed in Sec. 4.5, this paper focuses primarily on assumptions utilizing Gaussian conditional kernels. However, other FM implementations might employ different path constructions, as suggested by recent proposals Kerrigan et al. (2023); Isobe et al. (2024). The theoretical implications of these alternative approaches remain an essential area for future research.

## REFERENCES

Michael S Albergo, Nicholas M Boffi, Michael Lindsey, and Eric Vanden-Eijnden. Multimarginal generative modeling with stochastic interpolants. *arXiv*, October 2023a.

- 540 Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying  
541 framework for flows and diffusions, 2023b. URL <https://arxiv.org/abs/2303.08797>.  
542
- 543 Michael S Albergo, Mark Goldstein, Nicholas M Boffi, Rajesh Ranganath, and Eric Vanden-Eijnden.  
544 Stochastic interpolants with data-dependent couplings. *arXiv [cs.LG]*, October 2023c.
- 545 Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic  
546 interpolants. In *The Eleventh International Conference on Learning Representations, 2023*.  
547
- 548 Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly  $d$ -Linear  
549 convergence bounds for diffusion models via stochastic localization. *arXiv [stat.ML]*, 2023a.
- 550 Joe Benton, George Deligiannidis, and Arnaud Doucet. Error bounds for flow matching methods.  
551 *arXiv*, May 2023b.  
552
- 553 Joey Bose, Tara Akhound-Sadegh, Guillaume Hugué, Kilian Fatras, Jarrid Rector-Brooks, Cheng-  
554 Hao Liu, Andrei Cristian Nica, Maksym Korablyov, Michael M Bronstein, and Alexander Tong.  
555 SE(3)-Stochastic flow matching for protein backbone generation. *arXiv*, October 2023.
- 556 Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling:  
557 User-Friendly bounds under minimal smoothness assumptions. *arXiv*, June 2023.  
558
- 559 Ian Dunn and David Ryan Koes. Mixed continuous and categorical flow matching for 3D de novo  
560 molecule generation. *arXiv*, April 2024.
- 561 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
562 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion En-  
563 glish, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow  
564 transformers for High-Resolution image synthesis. *arXiv*, March 2024.  
565
- 566 Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3D equiv-  
567 ariant diffusion for Target-Aware molecule generation and affinity prediction. In *The Eleventh  
568 International Conference on Learning Representations, 2023*.
- 569 Satoshi Hayakawa and Taiji Suzuki. On the minimax optimality and superiority of deep neural  
570 network learning over sparse parameter spaces. *Neural Netw.*, 123:343–361, March 2020.  
571
- 572 Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion  
573 for molecule generation in 3D. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba  
574 Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference  
575 on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 8867–  
576 8887. PMLR, November 2022.
- 577 Vincent Hu, Di Wu, Yuki Asano, Pascal Mettes, Basura Fernando, Björn Ommer, and Cees Snoek.  
578 Flow matching for conditional text generation in a few sampling steps. In Yvette Graham and  
579 Matthew Purver, editors, *Proceedings of the 18th Conference of the European Chapter of the  
580 Association for Computational Linguistics (Volume 2: Short Papers)*, pages 380–392, St. Julian’s,  
581 Malta, March 2024. Association for Computational Linguistics.
- 582 Vincent Tao Hu, Wenzhe Yin, Pingchuan Ma, Yunlu Chen, Basura Fernando, Yuki M Asano, Efstratios  
583 Gavves, Pascal Mettes, Bjorn Ommer, and Cees G M Snoek. Motion flow matching for human  
584 motion synthesis and editing. *arXiv*, December 2023.  
585
- 586 Noboru Isobe, Masanori Koyama, Kohei Hayashi, and Kenji Fukumizu. Extended flow matching: a  
587 method of conditional generation with generalized continuity equation. *arXiv*, February 2024.
- 588 Yuling Jiao, Yanming Lai, Yang Wang, and Bokai Yan. Convergence analysis of flow matching in  
589 latent space with transformers. *J. Mach. Learn. Res.*, 23:1–65, April 2024.  
590
- 591 Gavin Kerrigan, Giosue Migliorini, and Padhraic Smyth. Functional flow matching. *arXiv*, 2023.  
592
- 593 Leon Klein, Andreas Krämer, and Frank No’e. Equivariant flow matching. *Neural Inf Process Syst.*,  
abs/2306.15030, June 2023.

- 594 Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vi-  
595 mal Manohar, Yossi Adi, Jay Mahadeokar, and Wei-Ning Hsu. Voicebox: Text-Guided multilingual  
596 universal speech generation at scale. *arXiv*, November 2023.
- 597 Yaron Lipman, Ricky T Q Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow  
598 matching for generative modeling. In *The Eleventh International Conference on Learning Repre-*  
599 *sentations*, 2023.
- 600 Linxi Liu, Dangna Li, and Wing Hung Wong. Convergence rates of a class of multivariate density  
601 estimation methods based on adaptive partitioning. *J. Mach. Learn. Res.*, 24(50):1–64, 2023a.
- 602 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate  
603 and transfer data with rectified flow. In *The Eleventh International Conference on Learning*  
604 *Representations*, 2023b.
- 605 Ryumei Nakada and Masaaki Imaizumi. Adaptive approximation and generalization of deep neural  
606 network with intrinsic dimensionality. *J. Mach. Learn. Res.*, 21(174):1–38, 2020.
- 607 Jonathan Niles-Weed and Quentin Berthet. Minimax estimation of smooth densities in Wasserstein  
608 distance. *Ann. Stat.*, 50(3):1519–1540, 2022.
- 609 Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution  
610 estimators. volume 202, pages 26517–26582. PMLR, 4 2023. URL [https://proceedings.](https://proceedings.mlr.press/v202/oko23a.html)  
611 [mlr.press/v202/oko23a.html](https://proceedings.mlr.press/v202/oko23a.html).
- 612 Philipp Petersen and Felix Voigtlaender. Optimal approximation of piecewise smooth functions using  
613 deep ReLU neural networks. *Neural Netw.*, 108:296–330, December 2018.
- 614 Aram-Alexandre Pooladian, Carles Domingo-Enrich, Ricky T. Q. Chen, and Brandon Amos. Neural  
615 optimal transport with Lagrangian costs. In *ICML Workshop on New Frontiers in Learning,*  
616 *Control, and Dynamical Systems*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=myb0FKB8C9)  
617 [myb0FKB8C9](https://openreview.net/forum?id=myb0FKB8C9).
- 618 Johannes Schmidt-Hieber. Deep ReLU network approximation of functions on a manifold. *arXiv*,  
619 August 2019.
- 620 David W Scott. *Multivariate density estimation: Theory, practice, and visualization*. Wiley Series in  
621 Probability and Statistics. John Wiley & Sons, Nashville, TN, August 1992.
- 622 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
623 Poole. Score-based generative modeling through stochastic differential equations. In *International*  
624 *Conference on Learning Representations*, October 2020.
- 625 Taiji Suzuki. Adaptivity of deep ReLU network for learning in besov and mixed smooth besov spaces:  
626 optimal rate and curse of dimensionality. In *International Conference on Learning Representations*,  
627 2019.
- 628 Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrid Rector-  
629 Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models  
630 with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024.
- 631 Ling Yang, Zixiang Zhang, Zhilong Zhang, Xingchao Liu, Minkai Xu, Wentao Zhang, Chenlin Meng,  
632 Stefano Ermon, and Bin Cui. Consistency flow matching: Defining straight flows with velocity  
633 consistency. *arXiv [cs.CV]*, July 2024.
- 634 Kaihong Zhang, Heqi Yin, Feng Liang, and Jingbo Liu. Minimax optimality of score-based diffusion  
635 models: Beyond the density lower bound assumptions. In *International Conference on Machine*  
636 *Learning*, pages 60134–60178. PMLR, July 2024.

644  
645  
646  
647

## Appendix

This section summarizes some basic mathematical facts, which can be easily derived and used in the proof of our main results, and known facts developed in [Oko et al. \(2023\)](#).

### A BASIC MATHEMATICAL RESULTS

#### A.1 DEFINITION OF BESOV SPACE

Besov space is an extension of the Sobolev space, allowing for non-integer orders of smoothness, and is effective in measuring both the local regularity and the global behavior of functions. It is formally defined as follows.

Let  $\Omega$  be a domain in  $\mathbb{R}^d$ . For a function  $f \in L_{p'}(\Omega)$  for some  $p' \in (0, \infty]$ , the  $r$ -th modulus of smoothness of  $f$  is defined by

$$w_{r,p'}(f, t) = \sup_{\|\mathbf{h}\|_2 \leq t} \|\Delta_{\mathbf{h}}^r(f)\|_{p'},$$

where

$$\Delta_{\mathbf{h}}^r(f)(\mathbf{x}) = \begin{cases} \sum_{j=0}^r \binom{r}{j} (-1)^{r-j} f(\mathbf{x} + j\mathbf{h}) & \text{if } \mathbf{x} + j\mathbf{h} \in \Omega \text{ for all } j, \\ 0 & \text{otherwise.} \end{cases}$$

For  $0 < p', q' \leq \infty$ ,  $s > 0$ ,  $r := |s| + 1$ , let the seminorm  $|\cdot|_{B_{p',q'}^s}$  be defined by

$$|f|_{B_{p',q'}^s} := \begin{cases} \left( \int_0^\infty (t^{-s} w_{r,p'}(f, t))^{q'} \frac{dt}{t} \right)^{\frac{1}{q'}} & (q' < \infty), \\ \sup_{t>0} t^{-s} w_{r,p'}(f, t) & (q' = \infty). \end{cases}$$

The norm of the Besov space  $B_{p',q'}^s(\Omega)$  is defined by

$$\|f\|_{B_{p',q'}^s} := \|f\|_{p'} + |f|_{B_{p',q'}^s},$$

and

$$B_{p',q'}^s(\Omega) := \{f \in L_{p'}(\Omega) \mid \|f\|_{B_{p',q'}^s} < \infty\}.$$

The parameter  $s$  serves as the order of smoothness. If  $p' = q'$  and  $s$  is an integer,  $B_{p',q'}^s$  coincides with the Sobolev space. For details of Besov spaces, see [Triebel \(1992\)](#), for example.

#### A.2 LIPSCHITZ CONDITION

This section shows a lemma that provides a Lipschitz constant of  $\mathbf{x} \mapsto \mathbf{v}_t(\mathbf{x})$  under the assumptions (A3) and (A5).

**Lemma 10.** *Let  $\mathbf{v}_t(\mathbf{x})$  be a vector field defined by (3) and (8), i.e.,*

$$\mathbf{v}_t(\mathbf{x}) = \int \left\{ \sigma'_t \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} + m'_t \mathbf{y} \right\} p_t(\mathbf{y}|\mathbf{x}) d\mathbf{y}$$

where

$$p_t(\mathbf{y}|\mathbf{x}) = \frac{p_t(\mathbf{x}|\mathbf{y})p_0(\mathbf{y})}{\int p_t(\mathbf{x}|\tilde{\mathbf{y}})p_0(\tilde{\mathbf{y}})d\tilde{\mathbf{y}}}, \quad p_t(\mathbf{x}|\mathbf{y}) = \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}}.$$

Then, under (A3) and (A5),  $\mathbf{v}_t(\mathbf{x})$  is Lipschitz continuous with Lipschitz constant  $\tilde{C}_L/t$  for any sufficiently small  $t$ , where  $\tilde{C}_L$  is independent of  $t$ .

*Proof.* By the definition of  $\mathbf{v}_t$  and the form of  $\sigma_t$  and  $m_t$  in (A3), we can compute explicitly

$$\frac{\partial \mathbf{v}_t(\mathbf{x})}{\partial \mathbf{x}} = \frac{\kappa}{t} I_d + \int \left\{ \frac{\kappa}{t} (\mathbf{x} - m_t \mathbf{y}) + \frac{\bar{\kappa}}{t} \mathbf{y} \right\} \frac{\partial p_t(\mathbf{y}|\mathbf{x})}{\partial \mathbf{x}} d\mathbf{y}.$$

As  $\frac{\partial}{\partial \mathbf{x}} \int p_t(\mathbf{y}|\mathbf{x}) d\mathbf{y} = \frac{\partial}{\partial \mathbf{x}} 1 = 0$ , we further obtain

$$\frac{\partial \mathbf{v}_t(\mathbf{x})}{\partial \mathbf{x}} = \frac{\kappa}{t} I_d + \frac{\bar{\kappa} - m_t \kappa}{t} \frac{\partial}{\partial \mathbf{x}} \int \mathbf{y} p_t(\mathbf{y}|\mathbf{x}) d\mathbf{y}.$$

The claim is obvious under (A5).  $\square$

### 702 A.3 WASSERSTEIN DISTANCE FOR CONVOLUTION

703 The following lemma is used in the proof sketch of Theorem 9, where  $W_2(P_t, P_{T_0})$  is bounded.

704 **Lemma 11.** *Let  $P$  be a probability distribution on  $\mathbb{R}^d$  with  $V := \int \|\mathbf{y}\|^2 dP(\mathbf{y}) < \infty$  and  $P_{m,\sigma}$  be*  
 705 *given by the density  $\int \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp(-\frac{\|\mathbf{x}-m\mathbf{y}\|^2}{2\sigma^2}) dP(\mathbf{y})$ . Then,*

$$706 W_2(P, P_{m,\sigma}) \leq \sqrt{(1-m)^2 V + d\sigma^2}.$$

707 *Proof.* The proof is elementary, but we include it for completeness. Let  $Y$  and  $Z$  be independent  
 708 random variables with probability  $P$  and  $N_d(0, I_d)$ , respectively. Let  $X := mY + \sigma Z$ , then the  
 709 distribution of  $X$  is  $P_{m,\sigma}$ . Considering a coupling  $(X, Y)$ ,

$$710 W_2(P, P_{m,\sigma})^2 \leq E\|X - Y\|^2$$

$$711 = E\|(m-1)Y + \sigma Z\|^2$$

$$712 = (1-m)^2 V + d\sigma^2,$$

713 which completes the proof.  $\square$

### 714 A.4 APPROXIMATION OF A FUNCTION IN BESOV SPACE

715 In this subsection, we present several approximation results developed in Suzuki (2019); Oko et al.  
 716 (2023). Although these results are already known, we include them here for ease of reference.

717 Let  $\mathcal{N}(x)$  be the function defined by  $\mathcal{N}(x) = 1$  for  $x \in [0, 1]$  and 0 otherwise. The *cardinal B-spline*  
 718 of order  $\ell \in \mathbb{N}$  is defined by

$$719 \mathcal{N}_\ell(x) := \mathcal{N} * \mathcal{N} * \dots * \mathcal{N}(x),$$

720 which is the convolution  $(\ell + 1)$  times of  $\mathcal{N}$ . Here, the convolution  $f * g$  is defined by

$$721 (f * g)(x) = \int f(x-y)g(y)dy.$$

722 For a multi-index  $k \in \mathbb{N}^d$  and  $j \in \mathbb{Z}^d$ , the *tensor product B-spline basis* in  $\mathbb{R}^d$  of order  $\ell$  is defined by

$$723 M_{k,j}^d(\mathbf{x}) := \prod_{i=1}^d \mathcal{N}_\ell(2^{k_i} x_i - j_i).$$

724 The following theorem says that a function  $f$  in the Besov space is approximated by a superposition  
 725 of  $M_{k,j}^d(\mathbf{x})$  of the form

$$726 f_N(\mathbf{x}) = \sum_{(k,j)} \alpha_{k,j} M_{k,j}^d(\mathbf{x}).$$

727 In the sequel, we fix the order  $\ell$  of the B-spline.  $(a)_+$  denotes  $\max\{0, a\}$ .

728 **Theorem 12** (Oko et al. (2023); Suzuki (2019)). *Let  $C > 0$  and  $0 < p', q', r \leq \infty$ . Under*  
 729  *$s > d(1/p' - 1/r)_+$  and  $0 < s < \min\{\ell, \ell - 1 + 1/p'\}$ , where  $\ell \in \mathbb{N}$  is the order of the cardinal*  
 730 *B-spline bases, for any  $f \in B_{p',q'}^s([-C, C]^d)$ , there exists  $f_N$  that satisfies*

$$731 \|f - f_N\|_{L_r([-C, C]^d)} \lesssim C^s N^{-s/d} \|f\|_{B_{p',q'}^s([-C, C]^d)}$$

732 for  $N \gg 1$  and has the following form:

$$733 f_N(\mathbf{x}) = \sum_{k=0}^K \sum_{j \in J(k)} \alpha_{k,j} M_{k,j}^d(\mathbf{x}) + \sum_{k=K+1}^{K^*} \sum_{i=1}^{n_k} \alpha_{k,j_i} M_{k,j_i}^d(\mathbf{x})$$

734 with

$$735 \sum_{k=0}^K |J(k)| + \sum_{k=K+1}^{K^*} n_k = N,$$

736 where  $J(k) = \{-C2^k - \ell, -C2^k - \ell + 1, \dots, C2^k - 1, C2^k\}$ ,  $(j_i)_{i=1}^{n_k} \subset J(k)$ ,  $K =$   
 737  $O(d^{-1} \log(N/C^d))$ ,  $K^* = (O(1) + \log(N/C^d))\nu^{-1} + K$ ,  $n_k = O((N/C^d)2^{-\nu(k-K)})$  ( $k =$   
 738  $K + 1, \dots, K^*$ ) for  $\nu = (s - \omega)/(2\omega)$  with  $\omega = d(1/p' - 1/r)_+$ . Moreover, we can take  $\alpha_{k,j}$  so  
 739 that  $|\alpha_{k,j}| \leq N^{(\nu^{-1} + d^{-1})(d/p' - s)_+}$ .

Based on the above theorem, the following result shows the accuracy of approximating the true density  $p_0$  by the  $B$ -spline functions with support restriction.

**Theorem 13** (Oko et al. (2023)). *Under Assumptions (A1)-(A5), there exists  $f_N$  of the form*

$$f_N(\mathbf{x}) = \sum_{i=1}^N \alpha_i \mathbf{1}[\|\mathbf{x}\|_\infty \leq 1] M_{k_i, j_i}^d(\mathbf{x}) + \sum_{i=N+1}^{3N} \alpha_i \mathbf{1}\left[\|\mathbf{x}\|_\infty \leq 1 - N^{-\frac{\kappa-1-\delta}{d}}\right] M_{k_i, j_i}^d(\mathbf{x}), \quad (25)$$

that satisfies

$$\|p_0 - f_N\|_{L^2(I^d)} \leq C_a N^{-s/d}, \quad (26)$$

$$\|p_0 - f_N\|_{L^2(I^d \setminus I_N^d)} \leq C_a N^{-\tilde{s}/d}, \quad (27)$$

for some  $C_a > 0$  and  $f_N(\mathbf{x}) = 0$  for any  $x$  with  $\|\mathbf{x}\|_\infty \geq 1$ . Here,  $-2^{(k_i)_m} - \ell \leq (j_i)_m \leq 2^{(k_i)_m}$  ( $i = 1, 2, \dots, N$ ;  $m = 1, 2, \dots, d$ ),  $|k_i| \leq K^* = (O(1) + \log N)\nu^{-1} + O(d^{-1} \log N)$  for  $\nu = (2s - \omega)/(2\omega)$  with  $\omega = d(1/p - 1/2)_+$ . The notations  $(k_i)_m$  and  $(j_i)_m$  are the  $m$ -th component of the multi-indices  $k_i$  and  $j_i$ , respectively. Moreover, we can take  $|\alpha_i| \leq N(\nu^{-1} + d^{-1})(d/p - s)_+$ .

*Proof.* See Oko et al. (2023, Lemma B.4).  $\square$

The following result shows the accuracy of approximating the ‘‘smoothed’’  $B$ -spline basis function by a neural network. This is essential to consider the approximation of  $p_t(\mathbf{x})$  and  $\mathbf{v}_t(\mathbf{x})$  by a neural network taken for  $p_0(\mathbf{x})$ .

**Theorem 14.** *Let  $C > 0$ ,  $k \in \mathbb{Z}_+$ ,  $j \in \mathbb{Z}^d$ ,  $\ell \in \mathbb{Z}_+$  with  $-C2^k - \ell \leq j_i \leq C2^k$  ( $i = 1, 2, \dots, d$ ), and  $C_{b,1} = 1$  or  $1 - N^{-(1-\delta)}$ . For any  $\varepsilon$  ( $0 < \varepsilon < 1/2$ ), there exists a neural network  $\phi_3^{k,j}(\mathbf{x}, t)$  and  $\phi_4^{k,j}(\mathbf{x}, t)$  in  $\mathcal{M}(L, W, S, B)$  with*

$$\begin{aligned} L &= O(\log^4 \varepsilon^{-1} + \log^2 C + k^2), \\ \|W\|_\infty &= O(\log^6 \varepsilon^{-1} + \log^3 C + k^3), \\ S &= O(\log^8 \varepsilon^{-1} + \log^4 C + k^4), \\ B &= \exp(O(\log \varepsilon^{-1} \log \log \varepsilon^{-1} + \log C + k)) \end{aligned} \quad (28)$$

such that

$$\left| \phi_3^{k,j}(\mathbf{x}, t) - \int_{\mathbb{R}^d} \mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d(\mathbf{y}) \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{y} \right| \leq \varepsilon \quad (29)$$

and

$$\left| \phi_4^{k,j}(\mathbf{x}, t) - \int_{\mathbb{R}^d} \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d(\mathbf{y}) \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{y} \right| \leq \varepsilon \quad (30)$$

hold for all  $\mathbf{x} \in [-C, C]^d$ . Furthermore, we can choose the networks so that  $\|\phi_3^{k,j}\|_\infty$  and  $\|\phi_4^{k,j}\|_\infty$  are of class  $O(1)$ .

*Proof.* See Oko et al. (2023, Lemma B.3).  $\square$

#### A.5 APPROXIMATION OF GAUSSIAN INTEGRALS

The following lemma is an elementary fact about Gaussian integrals used in the proof of Theorem 7 in Section C.4. We include it here for completeness.

**Lemma 15.** *Let  $\mathbf{x} \in \mathbb{R}^d$ ,  $0 < \varepsilon < 1/2$ , and  $\alpha \in \{0, 1\}$ . For any function  $F(\mathbf{y})$  supported on  $I^d$ , there is  $C_b > 0$  that depends only on  $d$  such that*

$$\left| \int_{I^d} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} F(\mathbf{y}) d\mathbf{y} - \int_{A_{\mathbf{x}}} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} F(\mathbf{y}) d\mathbf{y} \right| \leq \varepsilon, \quad (31)$$

where

$$A_{\mathbf{x}} := \left\{ \mathbf{y} \in I^d \mid \left\| \mathbf{y} - \frac{\mathbf{x}}{m_t} \right\|_\infty \leq C_b \frac{\sigma_t \sqrt{\log N}}{m_t} \right\}.$$

*Proof.* See Oko et al. (2023, Lemma F.9).  $\square$

## B PROOF OF THEOREM 3

Although the proof of Theorem 3 is basically the same as the proof of Benton et al. (2023b, Theorem 1), a slight difference appears since the current bound shows for arbitrary time  $t$ . We include the proof here for completeness.

Let  $\mathbf{v}(\mathbf{x}, t)$  and  $\widehat{\mathbf{v}}(\mathbf{x}; t)$  be smooth vector fields and  $\varphi_{s,t}$  and  $\widehat{\varphi}_{s,t}$  be respective flows;

$$\begin{aligned} \frac{d}{dt}\varphi_{s,t}(\mathbf{x}) &= \mathbf{v}(\varphi_{s,t}(\mathbf{x}), t), & \varphi_{s,s}(\mathbf{x}) &= \mathbf{x}, \\ \frac{d}{dt}\widehat{\varphi}_{s,t}(\mathbf{x}) &= \widehat{\mathbf{v}}(\widehat{\varphi}_{s,t}(\mathbf{x}), t), & \widehat{\varphi}_{s,s}(\mathbf{x}) &= \mathbf{x} \end{aligned}$$

**Lemma 16** (Alekseev-Gröbner). *Under the above notations, for any  $T \geq 0$ ,*

$$\widehat{\varphi}_{0,T}(\mathbf{x}_0) - \varphi_{0,T}(\mathbf{x}_0) = \int_0^T (\nabla_x \widehat{\varphi}_{s,T}(\mathbf{x})|_{x=\varphi_{0,s}(\mathbf{x}_0)}) (\widehat{\mathbf{v}}(\varphi_{0,s}(\mathbf{x}_0), s) - \mathbf{v}(\varphi_{0,s}(\mathbf{x}_0), s)) ds. \quad (32)$$

Note that on the right-hand side, the vector fields  $\widehat{\mathbf{v}}$  and  $\mathbf{v}$  are evaluated at the same point  $\varphi_{0,s}(\mathbf{x}_0)$ .

*Proof of Theorem 3.* Let  $\widehat{P}_t := (\widehat{\varphi}_{0,t})_{\#} P_0$  and  $P_t := (\varphi_{0,t})_{\#} P_0$  be pushforwards. By the definition of the 2-Wasserstein metric,

$$W_2(\widehat{P}_t, P_t) \leq \left( \int \|\widehat{\varphi}_{0,t}(\mathbf{x}_0) - \varphi_{0,t}(\mathbf{x}_0)\|^2 dP_0(\mathbf{x}_0) \right)^{1/2}. \quad (33)$$

From Lemma 16,

$$\|\widehat{\varphi}_{0,t}(\mathbf{x}_0) - \varphi_{0,t}(\mathbf{x}_0)\| \leq \int_0^t \left\| \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|_{op} \|\widehat{\mathbf{v}}(\varphi_{0,s}(\mathbf{x}_0), s) - \mathbf{v}(\varphi_{0,s}(\mathbf{x}_0), s)\| ds.$$

As a general relation of the largest singular value, we have

$$\frac{\partial}{\partial t} \left\| \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|_{op} \leq \left\| \frac{\partial}{\partial t} \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|_{op}.$$

Then, it follows from  $\frac{\partial}{\partial t} \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) = \nabla_y \widehat{\mathbf{v}}(y, t)|_{y=\varphi_{0,s}(\mathbf{x})} \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x})$  that the inequality

$$\frac{\partial}{\partial t} \left\| \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|_{op} \leq L_t \left\| \nabla_x \varphi_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|$$

holds by the  $L_t$ -Lipschitzness of  $\widehat{\mathbf{v}}(\cdot, t)$ . Accordingly, noting that  $\|\nabla_x \widehat{\varphi}_{s,s}(\mathbf{x})\|_{op} = \|\nabla_x \mathbf{x}\|_{op} = 1$ , the standard ODE argument leads to

$$\left\| \nabla_x \widehat{\varphi}_{s,t}(\mathbf{x}) \Big|_{x=\varphi_{0,s}(\mathbf{x}_0)} \right\|_{op} \leq e^{\int_s^t L_u du},$$

and therefore

$$\|\widehat{\varphi}_{0,t}(\mathbf{x}_0) - \varphi_{0,t}(\mathbf{x}_0)\| \leq \int_0^t e^{\int_s^t L_u du} \|\widehat{\mathbf{v}}(\varphi_{0,s}(\mathbf{x}_0), s) - \mathbf{v}(\varphi_{0,s}(\mathbf{x}_0), s)\| ds.$$

From Cauchy-Schwarz inequality,

$$\begin{aligned} & \|\widehat{\varphi}_{0,t}(\mathbf{x}_0) - \varphi_{0,t}(\mathbf{x}_0)\|^2 \\ & \leq \int_0^t 1^2 ds \int_0^t e^{2 \int_s^t L_u du} \|\widehat{\mathbf{v}}(\varphi_{0,s}(\mathbf{x}_0), s) - \mathbf{v}(\varphi_{0,s}(\mathbf{x}_0), s)\|^2 ds \\ & = t \int_0^t e^{2 \int_s^t L_u du} \|\widehat{\mathbf{v}}(\varphi_{0,s}(\mathbf{x}_0), s) - \mathbf{v}(\varphi_{0,s}(\mathbf{x}_0), s)\|^2 ds. \end{aligned}$$

Since the distribution of  $\varphi_{0,s}(\mathbf{x}_0)$  with  $x_0 \sim P_0$  is given by  $P_s$ , combining the above bound with (33) completes the proof.  $\square$



## C PROOF OF MAIN THEOREMS

### C.1 PROOF OF LEMMA 5: SUP OF LOSS FUNCTIONS

*Proof.*

$$\begin{aligned} \ell_\phi(\mathbf{x}) &= \int_{T_\ell}^{T_u} \int \|\phi(\mathbf{x}_t, t) - \mathbf{v}_t(\mathbf{x}_t|\mathbf{x})\|^2 p_t(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}_t dt \\ &\leq 2 \int_{T_\ell}^{T_u} \int \|\phi(\mathbf{x}_t, t)\|^2 p_t(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}_t dt + 2 \int_{T_\ell}^{T_u} \int \|\mathbf{v}_t(\mathbf{x}_t|\mathbf{x})\|^2 p_t(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}_t dt \end{aligned}$$

From the definition of  $\mathcal{H}_n$  and Assumptions (A3-A4), the first term is bounded by

$$4C^2 \int_{T_0}^1 ((\sigma'_t)^2 \log n + (m'_t)^2) dt \leq 8\tilde{C} ((\log n)^{b+1} + (\log n)^b),$$

where  $\tilde{C} > 0$  is a constant, and  $b = 0$  for  $\kappa > 1/2$  and  $b = b_0$  for  $\kappa = 1/2$ . From the definition

$$\mathbf{v}_t(\mathbf{x}_t|\mathbf{x}) = \sigma'_t \frac{\mathbf{x}_t - m_t \mathbf{x}}{\sigma_t} + m'_t \mathbf{x}$$

and the fact  $\|\mathbf{x}\| \leq 1$ , the second term is upper bounded by

$$\begin{aligned} &4 \int_{T_0}^1 \int \left\{ (\sigma'_t)^2 \frac{\|\mathbf{x}_t - m_t \mathbf{x}\|^2}{\sigma_t^2} + (m'_t)^2 \|\mathbf{x}\|^2 \right\} p_t(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}_t dt \\ &\leq 4 \int_{T_0}^1 \{d(\sigma'_t)^2 + (m'_t)^2\} dt \leq 4d\tilde{D}(\log n)^b \end{aligned}$$

for some constants  $\tilde{D} > 0$ . Note that the first inequality is given by  $p_t(\mathbf{x}_t|\mathbf{x}) = N_d(m_t \mathbf{x}, \sigma_t^2 I_d)$ . This completes the proof.  $\square$

### C.2 DIVISION POINT

The appropriate division point of the time interval  $[T_0, 1]$  arises from the width  $a_0$  of the smoother region around the boundary of  $I^d$ , which is assumed in (A1). Suppose that the smoothness of  $p_0(x)$  is  $\check{s}$ , higher than  $s$ , in the region  $I^d \setminus [-1 + a_0, 1 - a_0]^d$ . We should make the smoother region as small as possible to ensure that  $p_0$  is essentially in  $B_{p,1}^s(I^d)$ . We set  $a_0 = N^{-\gamma}$  and consider the partition of  $I^d \setminus [-1 + a_0, 1 - a_0]^d$  by  $N^{d\gamma} - (N^\gamma - 2)^d$  cubes of size  $a_0 = N^{-\gamma}$  (see Figure C.1). Suppose that we use  $N^{\delta'}$  bases ( $N \gg 1$ ) of  $B$ -spline for each small cube with arbitrarily small  $\delta' > 0$ . The condition  $\delta' > 0$  guarantees that the approximation can be arbitrarily accurate in each small cube. Since  $p_0$  restricted to each cube has a smooth degree  $\check{s}$ , Theorem 12 tells that it can be approximated by a  $B$ -spline function with the accuracy

$$a_0^{-\check{s}} N^{-\check{s}\delta'/d}.$$

The total number of  $B$ -spline bases is then

$$(N^{d\gamma} - (N^\gamma - 2)^d) N^{\delta'} \sim N^{(d-1)\gamma + \delta'},$$

To make the number of bases equal to or less than  $N$ , which is the number used for the  $B_{p',q'}^s$ -region of  $p_0$ , we set  $\gamma = (1 - \delta\kappa)/d$ , i.e.,  $a_0 = N^{-(1-\delta\kappa)/d}$  (we set  $\delta' = \delta\kappa$  for notational simplicity).

As seen in the proof of Theorem 7, to obtain the desired bound, the deviation  $\sigma_t$  must satisfy  $\sigma_t \leq a_0$  to bound the integral around the boundary. When  $t$  is small so that  $\sigma_t \sim t^\kappa$ ,  $\sigma_t \leq a_0$  means  $t \lesssim T_* := N^{-(\kappa^{-1}-\delta)/d}$ , which gives a constant on  $T_*$ . Consequently, we divide the time interval into  $[T_0, T_*]$  and  $[T_*, 1]$ , and show the different bounds for the approximation error.

### C.3 BASIC BOUNDS OF $p_t(x)$ AND $v_t(x)$

Recall that

$$p_t(\mathbf{x}) = \int_{[-1,1]^d} \frac{1}{(2\pi\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} p_0(\mathbf{y}) d\mathbf{y}.$$

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

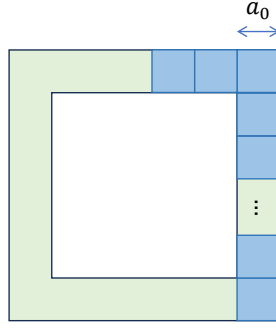


Figure C.1: Division of the cube into smoother small regions and the general region.

**Lemma 17.** *There exists  $C_1 = C_1(d, C_0) > 0$  such that*

$$C_1^{-1} \exp\left(-\frac{(\|\mathbf{x}\|_\infty - m_t)_+^2}{\sigma_t^2}\right) \leq p_t(\mathbf{x}) \leq C_1 \exp\left(-\frac{(\|\mathbf{x}\|_\infty - m_t)_+^2}{2\sigma_t^2}\right) \quad (34)$$

for any  $\mathbf{x} \in \mathbb{R}^d$  and  $t \in [T_0, 1]$ .

*Proof.* The proof is elementary and the same as Oko et al. (2023, Lemma A.2). We omit it.  $\square$

**Lemma 18.** (i) *Let  $k \in \mathbb{N}$  be arbitrary. There is  $C_2 = C_2(d, k, C_0) > 0$  such that*

$$\left\| \partial_{x_{i_1}} \cdots \partial_{x_{i_k}} p_t(\mathbf{x}) \right\| \leq \frac{C_2}{\sigma_t^k} \quad (\forall t \in [T_0, 1]).$$

(ii) *There is  $C_3 = C_3(d, C_0) > 0$  such that*

$$\|v_t(\mathbf{x})\| \leq C_3 \left\{ \sigma_t' \left( \frac{(\|\mathbf{x}\|_\infty - m_t)_+}{\sigma_t} \vee 1 \right) + |m_t'| \right\}$$

for any  $\mathbf{x} \in \mathbb{R}^d$  and  $t \in [T_0, 1]$ .

*Proof.* (i) is standard and the same as (Lemma A.3 Oko et al., 2023). We omit it.

For (ii), let  $g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) = p_t(\mathbf{x} | \mathbf{y}) = \frac{1}{(2\pi\sigma_t)^{d/2}} \exp\left\{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}\right\}$ . Note that

$$v_t(\mathbf{x}) = \frac{\int v_t(\mathbf{x} | \mathbf{y}) g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y}}{p_t(\mathbf{x})}.$$

The norm of the numerator is upper bounded by

$$\begin{aligned} & \left\| \int \left\{ \sigma_t' \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} + m_t' \mathbf{y} \right\} g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} \right\| \\ & \leq \sigma_t' \left\| \int \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} \right\| + |m_t'| \int \|\mathbf{y}\| g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y}. \end{aligned}$$

Since  $\text{Supp}(p_0) = [-1, 1]^d$  by assumption (A1), the second term in the last line is upper bounded by

$$|m_t'| \int g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} = |m_t'| p_t(\mathbf{x}). \quad (35)$$

To bound the first term, we use the restriction of the integral region as in Lemma F.9, Oko et al. (2023). Namely, letting  $\varepsilon := C_1^{-1} \exp\left(-\frac{(\|\mathbf{x}\|_\infty - m_t)_+^2}{2\sigma_t^2}\right)$ , the lower bound of  $p_t(\mathbf{x})$ , the integral is approximated as for any  $j = 1, \dots, d$ ,

$$\left| \int_{\mathbb{R}^d} \left( \frac{x_j - m_t y_j}{\sigma_t} \right)^\alpha g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} - \int_{A_x} \left( \frac{x_j - m_t y_j}{\sigma_t} \right)^\alpha g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} \right| \leq \varepsilon,$$

where  $A_x := \prod_{i=1}^d \left[ \frac{x_i}{m_t} - \frac{C\sigma_c}{m_t} \sqrt{\log(1/\varepsilon)}, \frac{x_i}{m_t} + \frac{C\sigma_t}{m_t} \sqrt{\log(1/\varepsilon)} \right]$ ,  $C$  is a positive constant depending only on  $d$ , and  $\alpha \in \{0, 1\}$ . Then, the first term is upper-bounded by

$$\sigma'_t \left\| \int_{A_x} \frac{x - m_t \mathbf{y}}{\sigma_t} g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} + \varepsilon \mathbf{1} \right\|.$$

Noting that  $\mathbf{y} \in A_x$  is equivalent to  $|x_i - m_t y_i|/\sigma_t \leq C\sqrt{\log(1/\varepsilon)}$ , the above quantity is further upper-bounded by

$$\sqrt{d} C \sigma'_t \sqrt{\log(1/\varepsilon)} \int_{A_x} g(\mathbf{x}; m_t \mathbf{y}, \sigma_t^2) p_0(\mathbf{y}) d\mathbf{y} + \sqrt{d} \varepsilon \leq C' \left( \sigma'_t \sqrt{\log(1/\varepsilon)} p_t(\mathbf{x}) + \varepsilon \right).$$

Noting that  $\varepsilon = C_1^{-1} \exp(-\frac{(\|\mathbf{x}\|_\infty - m_t)_+^2}{\sigma_t^2})$  and  $p_t(\mathbf{x}) \geq \varepsilon$ , we obtain

$$\begin{aligned} \|v_t(\mathbf{x})\| &\leq \frac{C'(\sigma'_t \sqrt{\log(1/\varepsilon)} p_t(\mathbf{x}) + \varepsilon) + |m'_t| p_t(\mathbf{x})}{p_t(\mathbf{x})} \\ &\leq C'' \left\{ \sigma'_t \left( \frac{(\|\mathbf{x}\|_\infty - m_t)_+}{\sigma_t} \right) + \sigma'_t + |m'_t| \right\} \\ &\leq C''' \left\{ \sigma'_t \left( \frac{(\|\mathbf{x}\|_\infty - m_t)_+}{\sigma_t} \vee 1 \right) + |m'_t| \right\}. \end{aligned}$$

This completes the proof.  $\square$

The following lemma shows an upper bound of  $v_t(\mathbf{x})$  when  $\mathbf{x}$  is in a bounded region of  $\sigma_t \sqrt{1/\varepsilon}$ , and presents bounds of relevant integrals.

**Lemma 19.** *Let  $\varepsilon > 0$  be sufficiently small.*

(i) *For any  $C_4 > 0$ , we have*

$$\|v_t(\mathbf{x})\| \leq C_4 (\sigma'_t \sqrt{\log(1/\varepsilon)} + |m'_t|) \quad (36)$$

*for any  $\mathbf{x}$  with  $\|\mathbf{x}\|_\infty \leq m_t + C_4 \sigma_t \sqrt{\log(1/\varepsilon)}$  and  $t \in [T_0, 1]$ .*

(ii) *For any  $C_5 > 0$ , there is  $\tilde{C} > 0$  such that*

$$\int_{\|\mathbf{x}\|_\infty \geq m_t + C_5 \sigma_t \sqrt{\log(1/\varepsilon)}} p_t(\mathbf{x}) \|v_t(\mathbf{x})\|^2 d\mathbf{x} \leq \tilde{C} \left\{ (\sigma'_t)^2 \log^{\frac{d}{2}}(\varepsilon^{-1}) + (m'_t)^2 \log^{\frac{d-2}{2}}(\varepsilon^{-1}) \right\} \varepsilon^{\frac{C_5^2}{2}} \quad (37)$$

*and*

$$\left| \int_{\|\mathbf{x}\|_\infty \geq m_t + C_5 \sigma_t \sqrt{\log(1/\varepsilon)}} p_t(\mathbf{x}) d\mathbf{x} \right| \leq \tilde{C} \log^{\frac{d-2}{2}}(\varepsilon^{-1}) \varepsilon^{\frac{C_5^2}{2}} \quad (38)$$

*hold for any  $\varepsilon > 0$  and  $t \in [T_0, 1]$ .*

*Proof.* (i) is obvious from Lemma 18, since  $(\|\mathbf{x}\|_\infty - m_t)_+/\sigma_t \leq C_4 \sqrt{\log(1/\varepsilon)}$  under the assumption of  $\mathbf{x}$ .

We show (ii). It follows from Lemmas 17 and 18 that

$$p_t(\mathbf{x}) \|v_t(\mathbf{x})\|^2 \leq 2C_1 C_3 \exp\left(-\frac{(\|\mathbf{x}\|_\infty - m_t \mathbf{y})_+^2}{2\sigma_t^2}\right) \left\{ (\sigma'_t)^2 \frac{(\|\mathbf{x}\|_\infty - m_t)_+^2}{\sigma_t^2} + (m'_t)^2 \right\}.$$

Let  $r := (\|\mathbf{x}\|_\infty - m_t)_+/\sigma_t$  and  $B_i := \{\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d \mid |x_i| = \max_{1 \leq j \leq d} |x_j|\}$ . The space is divided into  $d$  regions  $\cup_{i=1}^d B_i$  with measure-zero intersections. In  $B_1$ , the variables  $x_2, \dots, x_d$  satisfy  $|x_j| \leq m_t + C_4 \sigma_t \sqrt{\log(1/\varepsilon)}$ , and integral (37) is upper bounded by

$$\begin{aligned} &2C_1 C_3 d \int_{C_4 \sqrt{\log(1/\varepsilon)}} e^{-\frac{r^2}{2}} \left\{ (\sigma'_t)^2 r^2 + (m'_t)^2 \right\} (\sigma_t r + m_t)^{d-1} dr \\ &\leq C' \int_{C_4 \sqrt{\log(1/\varepsilon)}} e^{-\frac{r^2}{2}} \left\{ (\sigma'_t)^2 r^{d+1} + (m'_t)^2 r^{d-1} \right\} dr, \end{aligned}$$

where we use  $(\sigma_t r + m_t)^{d-1} \leq (r+1)^{d-1} \leq 2^{d-1} r^{d-1}$ .

For  $\ell \in \mathbb{N} \cup \{0\}$ , define

$$\psi_\ell(z) := \int_z^\infty r^\ell e^{-\frac{r^2}{2}} dr.$$

It is easy to see  $\psi_1(z) = e^{-\frac{z^2}{2}}$  and  $\psi_\ell(z) = x^{\ell-1} e^{-\frac{z^2}{2}} + (\ell-1)\psi_{\ell-1}(z)$  by partial integral. Using these formulas, we can see that for  $z \geq 1$

$$\psi_\ell(z) \leq B_\ell z^{\ell-1} e^{-\frac{z^2}{2}},$$

where  $B_\ell$  is a positive constant that depends only on  $\ell$ .

Thus, we obtain an upper bound

$$\tilde{C} \left\{ (\sigma'_t)^2 \log^{\frac{d}{2}} \left( \frac{1}{\varepsilon} \right) + (m'_t)^2 \log^{\frac{d-2}{2}} \left( \frac{1}{\varepsilon} \right) \right\} \varepsilon^{\frac{C_3^2}{2}},$$

which proves (37). The assertion (38) is similar.  $\square$

#### C.4 BOUNDS OF THE APPROXIMATION ERROR FOR SMALL $t$

This subsection shows the proof of Theorem 7.

##### (I) Restriction of the integral.

We first show that the left-hand side of (20) can be approximated by the integral over the bounded region

$$D_{t,N} := \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_\infty \leq m_t + C_4 \sigma_t \sqrt{\log N}\}.$$

To see this, observe that from Lemma 18 (ii), for  $\mathbf{x} \in D_{t,N}$  we have

$$\|\mathbf{v}_t(\mathbf{x})\|^2 \leq 2C_3^2 \{(\sigma'_t)^2 \log N + |m'_t|^2\}. \quad (39)$$

From the bound of  $\mathbf{v}_t(\mathbf{x})$ , we can restrict the neural network  $\phi(\mathbf{x}, t)$  so that it satisfies the same upper bound. Therefore, (39) is applied to  $\|\phi(\mathbf{x}, t)\|^2$  also. Combining this fact with Lemma 19 (ii), we have

$$\begin{aligned} & \int_{D_{t,N}^c} \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq 4C_3^2 \{(\sigma'_t)^2 \log N + |m'_t|^2\} \int_{D_{t,N}^c} \frac{1}{(\sqrt{2\pi}\sigma)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma^2}} d\mathbf{x} \\ & \leq 4C_3^2 \tilde{C} \{(\sigma'_t)^2 \log N + |m'_t|^2\} N^{-C_4^2/2} \log^{\frac{d-2}{2}} N. \end{aligned}$$

If  $C_4$  is taken large enough to satisfy  $C_4^2/2 > \frac{2s}{d}$ , it follows that

$$\begin{aligned} & \int \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq \int_{D_{t,N}} \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} + C' \{(\sigma'_t)^2 \log N + |m'_t|^2\} N^{-\frac{2s}{d}} \quad (40) \end{aligned}$$

for some constant  $C' > 0$ . Thus, we can consider the first term on the right-hand side.

Let  $\omega > 0$  be an arbitrary positive number. The integral over  $D_{t,N}$  in (40) can be restricted to the region  $\{\mathbf{x} \mid p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}\}$ . This can be easily seen by

$$\begin{aligned} & \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \leq N^{-\frac{2s+\omega}{d}}] \|\mathbf{v}_t(\mathbf{x}) - \phi(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq 4C_3^2 \int_{D_{t,N}} \{(\sigma'_t)^2 \log N + |m'_t|^2\} N^{-\frac{2s+\omega}{d}} d\mathbf{x} \\ & \leq 4C_3^2 N^{-\frac{2s+\omega}{d}} \{(\sigma'_t)^2 \log N + |m'_t|^2\} 2^d (m_t + C_4 \sigma_t \sqrt{\log N})^d \\ & \leq C'' \{(\sigma'_t)^2 \log N + |m'_t|^2\} N^{-\frac{2s+\omega}{d}} \log^{\frac{d}{2}} N, \end{aligned}$$

where  $C''$  depends only on  $d$ ,  $C_3$ , and  $C_4$ . This bound is of smaller order than the second term on the right-hand side of (40), and thus negligible. In summary, we have

$$\begin{aligned} & \int \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} + C' \{(\sigma'_t)^2 \log N + |m'_t|^2\} N^{-\frac{2s}{d}} \end{aligned} \quad (41)$$

for sufficiently large  $N$ .

**(II) Decomposition of integral.** Here, we give a bound of the integral  $\int \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x}$  over the region  $D_{t,N} \cap \{\mathbf{x} \mid p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}\}$ . The norm  $\|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|$  is bounded in detail.

First, recall that

$$\mathbf{v}_t(\mathbf{x}) = \frac{\int \mathbf{v}_t(\mathbf{x}|\mathbf{y}) p_t(\mathbf{x}|\mathbf{y}) p_0(\mathbf{y}) d\mathbf{y}}{p_t(\mathbf{x})}, \quad p_t(\mathbf{x}) = \int p_t(\mathbf{x}|\mathbf{y}) p_0(\mathbf{y}) d\mathbf{y}. \quad (42)$$

Based on Theorem 13, we can find  $f_N$  in (25) such that

$$\|p_0 - f_N\|_{L^2(I^d)} \leq C_a N^{-s/d}, \quad \|p_0 - f_N\|_{L^2(I^d \setminus I_N^d)} \leq C_a N^{-\tilde{s}/d} \quad (43)$$

As an approximate of  $p_t(\mathbf{x})$ , define a function  $\tilde{f}_1(\mathbf{x}, t)$  by

$$\tilde{f}_1(\mathbf{x}, t) := \int \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} f_N(\mathbf{y}) d\mathbf{y}. \quad (44)$$

Since we consider the region where  $p_t(\mathbf{x}) \geq N^{-(2s+\omega)/d}$ , we further define

$$f_1(\mathbf{x}, t) := \tilde{f}_1(\mathbf{x}, t) \vee N^{-\frac{2s+\omega}{d}}.$$

In a similar manner, the numerator of (42) can be approximated by

$$\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t),$$

where  $\mathbb{R}^d$ -valued functions  $\mathbf{f}_2$  and  $\mathbf{f}_3$  are defined by

$$\begin{aligned} \mathbf{f}_2(\mathbf{x}, t) &:= \int \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} f_N(\mathbf{y}) d\mathbf{y}. \\ \mathbf{f}_3(\mathbf{x}, t) &:= \int \mathbf{y} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} f_N(\mathbf{y}) d\mathbf{y}. \end{aligned} \quad (45)$$

We have an approximate of  $\mathbf{v}_t(\mathbf{x})$  by

$$\mathbf{f}_4(\mathbf{x}, t) := \frac{\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t)}{f_1(\mathbf{x}, t)} \mathbf{1} \left[ \left| \frac{\mathbf{f}_2}{f_1} \right| \leq C_5 \sqrt{\log N} \right] \mathbf{1} \left[ \left| \frac{\mathbf{f}_3}{f_1} \right| \leq C_5 \right].$$

We want to evaluate

$$\begin{aligned} & \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\phi(\mathbf{x}, t) - \mathbf{f}_4(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \quad + \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & =: I_A + I_B. \end{aligned} \quad (46)$$

**(III) Bound of  $I_A$  (neural network approximation of B-spline)**

We will approximate  $f_1$ ,  $f_2$ , and  $f_3$  by neural networks. For  $k \in \mathbb{Z}_+$  and  $j \in \mathbb{Z}^d$ , let  $E_{k,j,u}^{(a)}$  ( $a = 1, 2, 3, u = 0, 1$ ) denote the functions defined by

$$E_{k,j,u}^{(1)} := \int_{\mathbb{R}^d} \mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d(\mathbf{y}) \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{y}, \quad (47)$$

$$E_{k,j,u}^{(2)} := \int_{\mathbb{R}^d} \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d(\mathbf{y}) \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{y}, \quad (48)$$

and

$$E_{k,j,u}^{(3)} := \int_{\mathbb{R}^d} \mathbf{y} \mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d(\mathbf{y}) \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{y}, \quad (49)$$

where  $C_{b,1} = 1$  for  $u = 0$  and  $C_{b,1} = 1 - N^{-\frac{\kappa-1-\delta}{d}}$  for  $u = 1$ . Then, from Theorem 12,  $f_N$  is written as a linear combination of  $\mathbf{1}[\|\mathbf{y}\|_\infty \leq C_{b,1}] M_{k,j}^d$  with coefficients  $\alpha_{k,j}$ . From Theorem 13, for any  $\varepsilon > 0$ , there are neural networks  $\phi_5, \phi_6$  and  $\phi_7$  such that

$$|f_1(\mathbf{x}, t) - \phi_5(\mathbf{x}, t)| \leq D_5 N \max_i |\alpha_i| \varepsilon$$

$$\|f_2(\mathbf{x}, t) - \phi_6(\mathbf{x}, t)\| \leq D_6 N \max_i |\alpha_i| \varepsilon,$$

$$\|f_3(\mathbf{x}, t) - \phi_7(\mathbf{x}, t)\| \leq D_7 N \max_i |\alpha_i| \varepsilon.$$

Since  $\max_i |\alpha_i| \leq N^{-(\nu^{-1}+d^{-1})(d/p-s)_+}$ , by taking  $\varepsilon$  sufficiently small, for any  $\eta > 0$  we have

$$|f_1(\mathbf{x}, t) - \phi_5(\mathbf{x}, t)| \leq D_5 N^{-\eta}$$

$$\|f_2(\mathbf{x}, t) - \phi_6(\mathbf{x}, t)\| \leq D_6 N^{-\eta}$$

$$\|f_3(\mathbf{x}, t) - \phi_7(\mathbf{x}, t)\| \leq D_7 N^{-\eta}.$$

The operations to obtain the approximation of  $\mathbf{v}_t(\mathbf{x})$  based on  $\phi_5, \phi_6$ , and  $\phi_7$  are given by the following procedures:

$$\begin{aligned} \zeta_1 &:= \text{clip}(\phi_5; N^{-(2s+\omega)/d}, N^{K_0+1}), \\ \zeta_2 &:= \text{recip}(\zeta_1), \\ \zeta_3 &:= \text{mult}(\zeta_2, \phi_6), \\ \zeta_4 &:= \text{clip}(\zeta_3; -C_5\sqrt{\log N}, C_5\sqrt{\log N}), \\ \zeta_5 &:= \text{mult}(\zeta_2, \phi_7), \\ \zeta_6 &:= \text{clip}(\zeta_5; -C_5, C_5), \\ \zeta_7 &:= \text{mult}(\zeta_4, \hat{\sigma}'), \\ \zeta_8 &:= \text{mult}(\zeta_6, \hat{m}'), \\ \phi_8 &:= \zeta_7 + \zeta_8. \end{aligned}$$

As shown in Section D, the neural networks `clip`, `recip`, and `mult` achieve the approximation error  $N^{-\eta}$  with arbitrarily large  $\eta$ , while the complexity of the networks increases only at most polynomials of  $N$  for  $B$  and  $\|W\|_\infty$ , and at most  $\text{poly}(\log N)$  factor for  $L$  and  $S$ . In the upper bound of the generalization error, the network parameters  $B$  and  $\|W\|_\infty$  and the inverse error  $\varepsilon^{-1} = N^\eta$  appear only in the  $\log(\cdot)$  part to the log covering number, and  $S$  and  $B$  appear as a linear factor.

We also need to use the approximations  $\hat{\sigma}'_t$  and  $\hat{m}'_t$  of  $\sigma'_t$  and  $m'_t$ , respectively, by neural networks in construction. However, this can be done in a similar manner to (Section B1, Oko et al., 2023) with all network parameters  $O(\log^r \varepsilon^{-1})$  for approximation accuracy  $\varepsilon$ , and thus they have only  $O(\text{poly}(\log N))$  contributions. We omit the details in this paper. Consequently, the increase of the neural networks to obtain  $\phi_8$  from  $\phi_5, \phi_6$ , and  $\phi_7$  contributes the log covering number only by the  $\text{poly}(\log N)$  factor.

As a result, we obtain

$$I_A = O(N^{-\eta} \text{poly}(\log N)) \quad (50)$$

for arbitrary  $\eta > 0$ , while the required neural network increases the complexity term only by  $O(\text{poly}(\log N))$  factor.

**(IV) Bound of  $I_B$  (B-spline approximation of the true vector field)**

We evaluate here

$$I_B = \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x}. \quad (51)$$

Let  $\mathbf{h}_2(\mathbf{x}, t)$  and  $\mathbf{h}_3(\mathbf{x}, t)$  be functions defined by

$$\begin{aligned} \mathbf{h}_2(\mathbf{x}, t) &:= \int_{\mathbb{R}^d} \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} p_0(\mathbf{y}) d\mathbf{y}, \\ \mathbf{h}_3(\mathbf{x}, t) &:= \int_{\mathbb{R}^d} \mathbf{y} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} p_0(\mathbf{y}) d\mathbf{y}. \end{aligned} \quad (52)$$

Then,

$$\begin{aligned} &\|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\| \\ &= \mathbf{1} \left[ \left\| \frac{\mathbf{f}_2}{f_1} \right\| \leq C_5 \sqrt{\log N} \right] \mathbf{1} \left[ \left\| \frac{\mathbf{f}_3}{f_1} \right\| \leq C_5 \right] \left\| \frac{\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t)}{f_1(\mathbf{x})} - \frac{\sigma'_t \mathbf{h}_2(\mathbf{x}, t) + m'_t \mathbf{h}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\|. \end{aligned} \quad (53)$$

We evaluate the integral (51) by dividing  $D_{t,N}$  into the two domains  $\{\mathbf{x} \mid \|\mathbf{x}\|_\infty \leq m_t\}$  and  $\{\mathbf{x} \mid m_t \leq \|\mathbf{x}\|_\infty \leq m_t + C_4 \sigma_t \sqrt{\log N}\}$ .

Due to condition  $p_t(\mathbf{x}) \geq N^{-(2s+\omega)/d}$ , it suffices to take the network  $f_1(\mathbf{x}, t)$  so that it satisfies  $f_1(\mathbf{x}, t) \geq N^{-(2s+\omega)/d}$  by clipping the function if necessary. We therefore assume in the sequel that  $f_1(\mathbf{x}, t) \geq N^{-(2s+\omega)/d}$  holds.

**(IV-a) case:**  $\|\mathbf{x}\|_\infty \leq m_t$ .

In this case, Lemma 17 shows that  $C_1^{-1} \leq p_t(\mathbf{x}) \leq C_1$  for some  $C_1 > 0$ , which depends only on  $d$  and  $p_0$ . Using this fact and the conditions  $\left\| \frac{\mathbf{f}_2}{f_1} \right\| \leq C_5 \sqrt{\log N}$  and  $\left\| \frac{\mathbf{f}_3}{f_1} \right\| \leq C_5$ , we have

$$\begin{aligned} &\left\| \frac{\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t)}{f_1(\mathbf{x})} - \frac{\sigma'_t \mathbf{h}_2(\mathbf{x}, t) + m'_t \mathbf{h}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| \\ &\leq |\sigma'_t| \left\| \frac{\mathbf{f}_2(\mathbf{x}, t)}{f_1(\mathbf{x})} - \frac{\mathbf{h}_2(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| + |m'_t| \left\| \frac{\mathbf{f}_3(\mathbf{x}, t)}{f_1(\mathbf{x})} - \frac{\mathbf{h}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| \\ &\leq |\sigma'_t| \left\{ \left\| \frac{\mathbf{f}_2(\mathbf{x})}{f_1(\mathbf{x})} - \frac{\mathbf{f}_2(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| + \left\| \frac{\mathbf{f}_2(\mathbf{x}, t)}{p_t(\mathbf{x})} - \frac{\mathbf{h}_2(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| \right\} \\ &\quad + |m'_t| \left\{ \left\| \frac{\mathbf{f}_3(\mathbf{x})}{f_1(\mathbf{x})} - \frac{\mathbf{f}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| + \left\| \frac{\mathbf{f}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} - \frac{\mathbf{h}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| \right\} \\ &\leq C_1 |\sigma'_t| \left\{ C_5 \sqrt{\log N} |p_t(\mathbf{x}) - f_1(\mathbf{x}, t)| + \|\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t)\| \right\} \\ &\quad + C_1 |m'_t| \left\{ C_5 |p_t(\mathbf{x}) - f_1(\mathbf{x}, t)| + \|\mathbf{f}_3(\mathbf{x}, t) - \mathbf{h}_3(\mathbf{x}, t)\| \right\} \\ &\leq \tilde{C} \left\{ (|\sigma'_t| \sqrt{\log N} + |m'_t|) |f_1(\mathbf{x}, t) - p_t(\mathbf{x})| + |\sigma'_t| \|\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t)\| + |m'_t| \|\mathbf{f}_3(\mathbf{x}, t) - \mathbf{h}_3(\mathbf{x}, t)\| \right\}. \end{aligned} \quad (54)$$

We evaluate the integral on  $\{\mathbf{x} \mid \|\mathbf{x}\|_\infty \leq m_t\}$ . From the bound (54), we have

$$\begin{aligned} I_{B,1} &:= \int_{\|\mathbf{x}\|_\infty \leq m_t} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ &\leq C' \left[ \{(\sigma'_t)^2 \log N + (m'_t)^2\} \int_{\|\mathbf{x}\|_\infty \leq m_t} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] |f_1(\mathbf{x}, t) - p_t(\mathbf{x})|^2 p_t(\mathbf{x}) d\mathbf{x} \right. \\ &\quad + (\sigma'_t)^2 \int_{\|\mathbf{x}\|_\infty \leq m_t} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ &\quad \left. + (m'_t)^2 \int_{\|\mathbf{x}\|_\infty \leq m_t} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_3(\mathbf{x}, t) - \mathbf{h}_3(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} \right]. \end{aligned} \quad (55)$$

We write  $J_B^{(1)}$ ,  $J_B^{(2)}$ , and  $J_B^{(3)}$  for the three integrals that appear in the right-hand side of (55).

We will show only the derivation of an upper bound for  $J_B^{(2)}$ , since the other two cases are similar. Recall that by the definition of  $\mathbf{f}_2$  and  $\mathbf{h}_2$ ,

$$\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t) = \int_{I^d} \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y})) d\mathbf{y}.$$

Then, using  $p_t(\mathbf{x}) \leq C_1$ , we have

$$\begin{aligned} J_B^{(2)} &\leq C_1 \int_{\|\mathbf{x}\|_\infty \leq m_t} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \left\| \int_{I^d} \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y})) d\mathbf{y} \right\|^2 d\mathbf{x} \\ &\leq C_1 \int_{\|\mathbf{x}\|_\infty \leq m_t} \left\| \frac{1}{m_t^d} \int_{\mathbb{R}^d} \mathbf{1}[\|\mathbf{y}\|_\infty \leq 1] \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \left( \frac{m_t}{\sqrt{2\pi}\sigma_t} \right)^d e^{-\frac{\|\mathbf{y} - \mathbf{x}/m_t\|^2}{2(\sigma_t/m_t)^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y})) d\mathbf{y} \right\|^2 d\mathbf{x} \\ &\leq \frac{C_1}{m_t^{2d}} \int_{\|\mathbf{x}\|_\infty \leq m_t} \int_{\mathbb{R}^d} \mathbf{1}[\|\mathbf{y}\|_\infty \leq 1] \left\| \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \right\|^2 \left( \frac{m_t}{\sqrt{2\pi}\sigma_t} \right)^d e^{-\frac{\|\mathbf{y} - \mathbf{x}/m_t\|^2}{2(\sigma_t/m_t)^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\ &= \frac{C_1}{m_t^d} \int_{\|\mathbf{x}\|_\infty \leq m_t} \int_{\mathbb{R}^d} \mathbf{1}[\|\mathbf{y}\|_\infty \leq 1] \left\| \frac{\mathbf{x} - m_t \mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x}, \end{aligned} \quad (56)$$

where the third line uses Jensen's inequality for  $\|\cdot\|^2$ . For  $t \in [3N^{-\frac{\kappa-1-\delta}{d}}, 3N^{-\frac{\kappa-1-\delta}{d}}]$  with sufficiently large  $N$ , we can find  $c_0 > 0$  such that  $m_t \geq c_0$  on the time interval  $[T_0, 3N^{-\frac{\kappa-1-\delta}{d}}]$ . We can thus further obtain for some  $C' > 0$

$$\begin{aligned} J_B^{(2)} &\leq C' \int_{I^d} \int_{\mathbb{R}^d} \frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{\sigma_t^2} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x} - m_t \mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{x} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} \\ &= dC' \int_{I^d} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} \\ &= dC' \|f_N - p_0\|_{L^2(I^d)}^2 \\ &\leq C'' N^{-\frac{2s}{d}} \end{aligned} \quad (57)$$

by the choice of  $f_N$ . Similarly, we can prove that  $J_B^{(1)}$  and  $J_B^{(3)}$  have the same upper bounds of  $N^{-\frac{2s}{d}}$  order. This proves that there is  $C_{B,1} > 0$  such that

$$I_{B,1} \leq C_{B,1} \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\frac{2s}{d}}. \quad (58)$$

**(VI-b) case:**  $m_t \leq \|\mathbf{x}\|_\infty \leq m_t + C_4 \sigma_t \sqrt{\log N}$ .

Unlike case (i), we do not have a constant lower bound of  $p_t(\mathbf{x})$  in this region, and thus we resort to the bound  $p_t(\mathbf{x}) \geq N^{-(2s+\omega)/d}$ , that is  $1/p_t(\mathbf{x}) \leq N^{(2s+\omega)/d}$  and  $1/f_1(\mathbf{x}, t) \leq N^{(2s+\omega)/d}$ . We have

$$\begin{aligned} &\left\| \frac{\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t)}{f_1(\mathbf{x})} - \frac{\sigma'_t \mathbf{h}_2(\mathbf{x}, t) + m'_t \mathbf{h}_3(\mathbf{x}, t)}{p_t(\mathbf{x})} \right\| \\ &\leq \frac{1}{f_1(\mathbf{x}, t)} \|(\sigma'_t \mathbf{f}_2(\mathbf{x}, t) + m'_t \mathbf{f}_3(\mathbf{x}, t)) - (\sigma'_t \mathbf{h}_2(\mathbf{x}, t) + m'_t \mathbf{h}_3(\mathbf{x}, t))\| \\ &\quad + \|\mathbf{v}_t(\mathbf{x})\| \frac{1}{f_1(\mathbf{x}, t)} |f_1(\mathbf{x}) - p_t(\mathbf{x})| \\ &\leq N^{(2s+\omega)/d} \tilde{C} \left\{ (|\sigma'_t| \sqrt{\log N} + |m'_t|) |p_t(\mathbf{x}) - f_1(\mathbf{x}, t)| + |\sigma'_t| \|\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t)\| \right. \\ &\quad \left. + |m'_t| \|\mathbf{f}_3(\mathbf{x}, t) - \mathbf{h}_3(\mathbf{x}, t)\| \right\}, \end{aligned} \quad (59)$$

where in the last inequality we use Lemma 19 (i).



Let  $\Delta_{t,N} := \{\mathbf{x} \in \mathbb{R}^d \mid m_t \leq \|\mathbf{x}\|_\infty \leq m_t + C_4\sigma_t\sqrt{\log N}\}$ . By the same argument using Jensen's inequality as the derivation of (56), we can obtain

$$\begin{aligned} I_{B,2} &:= \int_{\Delta_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\frac{2s+\omega}{d}}] \|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ &\leq C''' N^{\frac{4s+2\omega}{d}} \left[ \{(\sigma'_t)^2 \log N + (m'_t)^2\} \int_{\Delta_{t,N}} \int_{I^d} \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \right. \\ &\quad + (\sigma'_t)^2 \int_{\Delta_{t,N}} \int_{I^d} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\ &\quad \left. + (m'_t)^2 \int_{\Delta_{t,N}} \int_{I^d} \|\mathbf{y}\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \right] \quad (60) \end{aligned}$$

Due to the factor  $N^{(4s+2\omega)/d}$ , the integrals must have orders smaller than in the case of (56) to derive the desired bound of  $I_{B,2}$ . We will make use of Assumption (A1) about the higher-order smoothness around the boundary of  $I^d$ .

Because the three integrals can be bounded in a similar manner, we focus only on the second one, denoted by  $K_B^{(2)}$ . Since  $\delta, \omega > 0$  can be taken arbitrarily small, we can assume  $\check{s} > 6s - 1 + \delta\kappa + 2\omega$ . From Lemma 15 with  $\varepsilon = N^{-\check{s}/d}$ , there is  $C_b > 0$ , which is independent of  $x, t$ , and sufficiently large  $N$ , such that

$$\begin{aligned} &\left| \int_{I^d} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} \right. \\ &\quad \left. - \int_{A_{\mathbf{x}}} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} \right| \leq N^{-\frac{\check{s}}{d}}, \quad (61) \end{aligned}$$

where  $A_{\mathbf{x}}$  is given by

$$A_{\mathbf{x}} := \left\{ \mathbf{y} \in I^d \mid \left\| \mathbf{y} - \frac{\mathbf{x}}{m_t} \right\|_\infty \leq C_b \frac{\sigma_t \sqrt{\log N}}{m_t} \right\}.$$

Note that if  $\mathbf{x} \in \Delta_{t,N}$  and  $\mathbf{y} \in A_{\mathbf{x}}$ , then

$$-1 \leq y_j \leq -1 + \frac{C_b \sigma_t \sqrt{\log N}}{m_t} \quad \text{or} \quad 1 - \frac{C_b \sigma_t \sqrt{\log N}}{m_t} \leq y_j \leq 1$$

for each  $j = 1, \dots, d$ . Because we assume that  $t \leq 3N^{-\frac{\kappa-1-\delta}{d}}$  and  $\sigma_t \sim b_0 t^\kappa$ , we can assume that  $m_t \geq \sqrt{D_0}/2$  from Assumption (A3). Then, for sufficiently large  $N$ , we see that  $\mathbf{y} \in I^d \setminus I_N^d$ . This can be seen in  $\sigma_t \sim b_0 t^\kappa \leq b_0 3^\kappa N^{-\frac{1-\delta\kappa}{d}}$  and thus  $C_b \sigma_t \sqrt{\log N}/m_t \leq \frac{2C_b b_0 3^\kappa}{\sqrt{D_0}} N^{-\frac{1-\delta\kappa}{d}}$ . For  $\mathbf{y} \in I^d \setminus I_N^d$ , we can use the second bound in (25);  $\|f_N - p_0\|_{L^2(I^d \setminus I_N^d)} \leq N^{-\check{s}/d}$ . It follows from (61) that

$$\begin{aligned} K_B^{(2)} &= \int_{\Delta_{t,N}} \int_{I^d} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\ &\leq \int_{\Delta_{t,N}} \left\{ \int_{A_{\mathbf{x}}} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} + N^{-\check{s}/d} \right\} d\mathbf{x} \\ &\leq \int_{I^d \setminus I_N^d} \int_{\mathbb{R}^d} \left\| \frac{\mathbf{x} - m_t\mathbf{y}}{\sigma_t} \right\|^2 \frac{1}{(\sqrt{2\pi}\sigma_t)^d} e^{-\frac{\|\mathbf{x}-m_t\mathbf{y}\|^2}{2\sigma_t^2}} d\mathbf{x} (f_N(\mathbf{y}) - p_0(\mathbf{y}))^2 d\mathbf{y} + N^{-\check{s}/d} |\Delta_{t,N}|. \end{aligned}$$

Since the volume  $|\Delta_{t,N}|$  is upper bounded by  $D'\sigma_t\sqrt{\log N}$  with some constant  $D' > 0$ , we have

$$\begin{aligned} K_B^{(2)} &\leq d \|f_N - p_0\|_{L^2(I^d \setminus I_N^d)}^2 + C'(\sigma_t \sqrt{\log N}) N^{-\check{s}/d} \\ &\leq C'' \left( N^{-2\check{s}/d} + N^{-(\check{s}+1-\delta\kappa)/d} \log^{d/2} N \right) \\ &= O \left( N^{-\frac{\check{s}+1-\delta\kappa}{d}} \log^{d/2} N \right), \quad (62) \end{aligned}$$

where the last line uses  $\check{s} > 1$  (A1). The integrals  $K_B^{(1)}$  and  $K_B^{(3)}$  have an upper bound of the same order. As a result, there is  $C_{B,2} > 0$  which does not depend on  $n$  or  $t$  such that

$$I_{B,2} \leq C_{B,2} \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\frac{\check{s}+1-4s-2\omega-\delta\kappa}{d}}.$$

Since we have taken  $\check{s}$  so that  $\check{s} > 6s - 1 + \delta\kappa + 2\omega$ , we have

$$I_{B,2} \leq C_{B,2} \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\frac{2\check{s}}{d}}. \quad (63)$$

### (VI-c)

It follows from (58) and (63) that there is  $C_B > 0$  such that

$$I_B \leq C_B \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\frac{2\check{s}}{d}} \quad (64)$$

for sufficiently large  $N$ .

### (V) Concluding the proof.

Combining (41), (46), (50), and (64) leads to the upper bound of the statement of the theorem. The argument of the network size in part (III) proves the corresponding statement.  $\square$

## C.5 BOUNDS OF THE APPROXIMATION ERROR FOR LARGER $t$

This subsection gives a proof of Theorem 8. The proof is parallel to that of Theorem 7 in many places, while the smoothness of the target density is more helpful.

**(I)' Restriction of the integral.** In a similar manner to part (I) of Section C.4, there is  $C_8 > 0$ , which does not depend on  $t$  such that for any neural network  $\phi(\mathbf{x}, t)$  with  $\|\phi(\mathbf{x}, t)\| \leq C_3\{|\sigma'_t|\sqrt{\log N} + |m'_t|\}$  the bound

$$\int_{\|\mathbf{x}\| \geq m_t + C_8\sqrt{\log N}} p_t(\mathbf{x}) \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 d\mathbf{x} \lesssim \{|\sigma'_t|\sqrt{\log N} + |m'_t|\} N^{-\eta}$$

holds for any  $t \in [N^{-\frac{\kappa-1-\delta}{d}}, 1]$ . Also, in a similar way, we can restrict the integral to the region  $\{\mathbf{x} \mid p_t(\mathbf{x}) \geq N^{-\eta}\}$  up to a negligible difference. Consequently, we obtain

$$\begin{aligned} & \int_{\mathbb{R}^d} p_t(\mathbf{x}) \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 d\mathbf{x} \\ &= \int_{\|\mathbf{x}\| \leq m_t + C_8\sqrt{\log N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\eta}] p_t(\mathbf{x}) \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 d\mathbf{x} \\ &+ O(\{|\sigma'_t|\sqrt{\log N} + |m'_t|\} N^{-\eta}). \end{aligned} \quad (65)$$

**(II)' Decomposition of integral.** We consider a  $B$ -spline approximation of  $p_t(\mathbf{x})$ . Unlike Section C.4, we can regard  $p_{t_*}$  as the target distribution, which is of class  $C^\infty$ . This will cause a tighter bound than in Section C.4 and easier analysis. More precisely, it is easy to see that  $p_t$ , the convolution between  $p_0$  and the Gaussian distribution  $N_d(m_t \mathbf{y}, \sigma_t^2 I_d)$ , can be rewritten as the convolution between  $p_{t_*}$  and  $N_d(\tilde{m}_t \mathbf{y}, \tilde{\sigma}_t^2 I_d)$  for any  $t > t_*$ , that is,

$$p_t(\mathbf{x}) = \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x}-\tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} p_{t_*}(\mathbf{y}) d\mathbf{y},$$

where

$$\tilde{m}_t := \frac{m_t}{m_{t_*}}, \quad \tilde{\sigma}_t := \sqrt{\sigma_t^2 - \left(\frac{m_t}{m_{t_*}}\right)^2 \sigma_{t_*}^2}. \quad (66)$$

We can thus apply a similar argument to Section C.4.

We use a  $B$ -spline approximation of  $p_{t_*}$ . For  $\eta > 0$ , take  $\alpha \in \mathbb{N}$  such that  $\alpha > \frac{3d\eta}{2\delta\kappa}$ . It is easy to see that the derivatives of  $p_{t_*}(\mathbf{x})$  satisfy

$$\left\| \frac{\partial^k}{\partial x_{i_1} \cdots \partial x_{i_k}} p_{t_*}(\mathbf{x}) \right\| \leq \frac{C_a}{\sigma_{t_*}^k}$$

for any  $k \leq \alpha$  and  $(i_1, \dots, i_k)$ . From Assumption (A3), there are  $t^\dagger \in [0, 1]$  and  $b^\dagger > 0$  such that  $\sigma_t \geq b^\dagger t^\kappa$  for any  $0 \leq t \leq t^\dagger$ . If we set  $c^\dagger := (t^\dagger)^\kappa > 0$ , then  $\sigma_t \geq b^\dagger t^\kappa \vee c^\dagger$  for any  $t \in [0, 1]$ . We can see that

$$\frac{p_{t_*}}{t_*^{-\alpha\kappa} \vee c^\dagger} \in B_{\infty, \infty}^\alpha(\mathbb{R}^d)$$

holds, because for any  $k \leq \alpha$  we have

$$\left\| \frac{\partial^k}{\partial x_{i_1} \cdots \partial x_{i_k}} \frac{p_{t_*}(\mathbf{x})}{t_*^{-\alpha\kappa} \vee c^\dagger} \right\| \leq \frac{C_\alpha (b^\dagger t_*^{-k\kappa} \wedge (c^\dagger)^{-k})}{t_*^{-\alpha\kappa} \vee c^\dagger} \leq C_\alpha (b^\dagger t_*^{(\alpha-k)\kappa} \wedge (c^\dagger)^{-(k+1)}) \leq C_\alpha ((b^\dagger \wedge (c^\dagger))^{-(k+1)}),$$

which implies  $\frac{p_{t_*}}{t_*^{-\alpha\kappa} \vee c^\dagger} \in W_\infty^\alpha(\mathbb{R}^d)$  and  $\left\| \frac{p_{t_*}}{t_*^{-\alpha\kappa} \vee c^\dagger} \right\|_{W_\infty^\alpha(\mathbb{R}^d)} \leq C_\alpha ((b^\dagger \wedge (c^\dagger))^{-(k+1)})$  (constant).

Notice that, by a similar argument as in the proof of Lemma 19 (ii), there is  $C_5 > 0$  such that

$$\int_{\|\mathbf{y}\|_\infty \geq C_5 \sqrt{\log N}} (\|\mathbf{y}\|^2 + 1) p_{t_*}(\mathbf{y}) d\mathbf{y} \leq N^{-3\eta} \quad (67)$$

holds. We therefore consider a  $B$ -spline approximation on  $[-C_5 \sqrt{\log N}, C_5 \sqrt{\log N}]^d$ . Letting

$$N_* := \lceil t_*^{-d\kappa} N^{\delta\kappa} \rceil$$

be the number of  $B$ -spline bases, from Theorem 12, we can find a function  $f_{N^*}$  of the form

$$f_{N^*}(\mathbf{x}) = (t_*^{-\alpha\kappa} \vee c^\dagger) \sum_{i=1}^{N^*} \alpha_i \mathbf{1}[\|\mathbf{x}\|_\infty \leq C_5 \sqrt{\log N}] M_{k_i, j_i}^d(\mathbf{x})$$

with  $|\alpha_i| \leq 1$  and  $C_9 > 0$  such that

$$\|p_{t_*} - f_{N^*}\|_{L^2([-C_5 \sqrt{\log N}, C_5 \sqrt{\log N}]^d)} \leq C_9 (\log N)^{\alpha/2} (N^*)^{-\frac{\alpha}{d}} (t_*^{-\alpha\kappa} \vee c^\dagger)$$

holds.

From  $N^* \geq t_*^{-d\kappa} N^{\delta\kappa}$  and  $\alpha > \frac{3d\eta}{2\delta\kappa}$ , the right-hand side is bounded by

$$C' (\log N)^{\alpha/2} N^{-\delta\alpha\kappa/d} \leq C' N^{-3\eta/2}$$

for sufficiently large  $N$ , which implies

$$\|p_{t_*} - f_{N^*}\|_{L^2([-C_5 \sqrt{\log N}, C_5 \sqrt{\log N}]^d)} \leq C_{10} N^{-3\eta/2} \quad (68)$$

for sufficiently large  $N$ . Note also that the coefficient  $\tilde{\alpha}_i := \alpha_i (t_*^{-\alpha\kappa} \vee c^\dagger)$  of the basis  $M_{k_i, j_i}^d$  in  $f_{N^*}$  is bounded by

$$|\tilde{\alpha}_i| \leq t_*^{-\alpha\kappa} \vee c^\dagger \leq N^{\frac{\kappa-1}{d}-\delta\alpha\kappa} = N^{\frac{\alpha(1-\delta\kappa)}{d}}.$$

In a similar manner to part (II) of Section C.4, define  $f_1, f_2, f_3$ , and  $f_4$  using  $f_{N^*}$  as

$$f_1(\mathbf{x}, t) := \tilde{f}_1(\mathbf{x}, t) \vee N^{-\eta} \quad \text{with} \quad \tilde{f}_1(\mathbf{x}, t) := \int \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x}-\tilde{m}_t\mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} f_{N^*}(\mathbf{y}) d\mathbf{y}, \quad (69)$$

$$f_2(\mathbf{x}, t) := \int \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x}-\tilde{m}_t\mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} f_{N^*}(\mathbf{y}) d\mathbf{y}.$$

$$f_3(\mathbf{x}, t) := \int \mathbf{y} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x}-\tilde{m}_t\mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} f_{N^*}(\mathbf{y}) d\mathbf{y}, \quad (70)$$

$$f_4(\mathbf{x}, t) := \frac{\tilde{\sigma}_t' f_2(\mathbf{x}, t) + \tilde{m}_t' f_3(\mathbf{x}, t)}{f_1(\mathbf{x}, t)} \mathbf{1} \left[ \left| \frac{f_2}{f_1} \right| \leq C_5 \sqrt{\log N} \right] \mathbf{1} \left[ \left| \frac{f_3}{f_1} \right| \leq C_5 \right].$$

We have a similar decomposition of the integral to (46):

$$\begin{aligned} & \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\eta}] \|\phi(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \leq \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\eta}] \|\phi(\mathbf{x}, t) - \mathbf{f}_4(\mathbf{x}, t)\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & \quad + \int_{D_{t,N}} \mathbf{1}[p_t(\mathbf{x}) \geq N^{-\eta}] \|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} \\ & =: \tilde{I}_A + \tilde{I}_B \end{aligned} \quad (71)$$

**(III)' Bound of  $\tilde{I}_A$  (neural network approximation of B-spline)**

Using exactly the same argument as in part (III) of Section C.4, we can show

$$\tilde{I}_A = O(\text{poly}(\log N)N^{-\eta'})$$

for arbitrary  $\eta' > 0$ , and thus it is negligible.

The size of the neural network is given by  $L = O(\log^4 N)$ ,  $\|W\|_\infty = O(N)$ ,  $S = O(N^*) = O(t_*^{-d\kappa}N^{\delta\kappa})$ , and  $B = \exp(O(\log N \log \log N))$ .

**(IV)' Bound of  $\tilde{I}_B$  (B-spline approximation of the true vector field)**

By replacing  $p_0$  with  $p_{t_*}$  in (52), define  $\mathbf{h}_2(\mathbf{x}, t)$  and  $\mathbf{h}_3(\mathbf{x}, t)$  with  $\tilde{m}_t$  and  $\tilde{\sigma}_t$  by

$$\begin{aligned}\mathbf{h}_2(\mathbf{x}, t) &:= \int_{\mathbb{R}^d} \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} p_{t_*}(\mathbf{y}) d\mathbf{y}, \\ \mathbf{h}_3(\mathbf{x}, t) &:= \int_{\mathbb{R}^d} \mathbf{y} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} p_{t_*}(\mathbf{y}) d\mathbf{y}.\end{aligned}$$

Then, by a similar argument to (19), we have

$$\begin{aligned}\|\mathbf{f}_4(\mathbf{x}, t) - \mathbf{v}_t(\mathbf{x})\| &\leq N^\eta \tilde{C} \left\{ (|\tilde{\sigma}'_t| \sqrt{\log N} + |\tilde{m}'_t|) |p_t(\mathbf{x}) - f_1(\mathbf{x}, t)| \right. \\ &\quad \left. + |\tilde{\sigma}'_t| \|\mathbf{f}_2(\mathbf{x}, t) - \mathbf{h}_2(\mathbf{x}, t)\| + |\tilde{m}'_t| \|\mathbf{f}_3(\mathbf{x}, t) - \mathbf{h}_3(\mathbf{x}, t)\| \right\}\end{aligned}$$

for some constant  $\tilde{C}$ , and thus

$$\begin{aligned}\tilde{I}_B &\leq C' N^{2\eta} \left[ \{(\tilde{\sigma}'_t)^2 \log N + (\tilde{m}'_t)^2\} \int_{D_{t,N}} \left\| \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 dx \right. \\ &\quad + (\tilde{\sigma}'_t)^2 \int_{D_{t,N}} \left\| \int_{\mathbb{R}^d} \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 dx \\ &\quad \left. + (\tilde{m}'_t)^2 \int_{D_{t,N}} \left\| \int_{\mathbb{R}^d} \mathbf{y} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 dx \right]. \quad (72)\end{aligned}$$

Here, we show a bound of

$$\tilde{J}_{B,2} := \int_{D_{t,N}} \left\| \int_{\mathbb{R}^d} \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 dx.$$

The other two integrals can be bounded similarly.

Let  $\rho := \frac{1}{\sqrt{2D_0}} > 0$ , where  $D_0$  is given in Assumption (A3). We derive a bound of  $\tilde{J}_{B,2}$  in the two cases of  $t$  separately: **(IV-a)'**  $\tilde{m}_t \geq \rho$ , and **(IV-b)'**  $\tilde{m}_t \leq \rho$ .

**Case (IV-a)'**:  $\tilde{m}_t \geq \rho$ .

By rewriting the inner integral on  $\mathbf{y}$  by a Gaussian integral, we have

$$\begin{aligned}
\tilde{J}_{B,2} &= \int_{D_{t,N}} \frac{1}{\tilde{m}_t^{2d}} \left\| \int_{\mathbb{R}^d} \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \left( \frac{\tilde{m}_t}{\sqrt{2\pi\tilde{\sigma}_t}} \right)^d e^{-\frac{\tilde{m}_t^2 \|\mathbf{y} - \mathbf{x}/\tilde{m}_t\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 d\mathbf{x} \\
&\leq \int_{D_{t,N}} \frac{1}{\tilde{m}_t^{2d}} \int_{\mathbb{R}^d} \left\| \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \right\|^2 \left( \frac{\tilde{m}_t}{\sqrt{2\pi\tilde{\sigma}_t}} \right)^d e^{-\frac{\tilde{m}_t^2 \|\mathbf{y} - \mathbf{x}/\tilde{m}_t\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\
&\leq \int_{D_{t,N}} \frac{1}{\tilde{m}_t^d} \int_{\mathbb{R}^d} \left\| \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \right\|^2 \frac{1}{(\sqrt{2\pi\tilde{\sigma}_t})^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\
&\leq (2D_0)^{d/2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \left\| \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \right\|^2 \frac{1}{(\sqrt{2\pi\tilde{\sigma}_t})^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y}))^2 d\mathbf{x} d\mathbf{y} \\
&\leq \rho^{-d/2} d \int_{\mathbb{R}^d} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y}))^2 d\mathbf{y} \\
&\leq \rho^{-d/2} d \left[ \int_{[-C_t\sqrt{\log N}, C_5\sqrt{\log N}]^d} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y}))^2 d\mathbf{y} + \int_{\|\mathbf{y}\| \geq C_5\sqrt{\log N}} p_{t_*}(\mathbf{y})^2 d\mathbf{y} \right] \\
&\leq \rho^{-d/2} d \left\{ \|f_{N^*} - p_{t_*}\|_{L^2([-C_5\sqrt{\log N}, C_5\sqrt{\log N}]^d)}^2 + N^{-3\eta} \right\} \\
&\leq C N^{-3\eta}
\end{aligned}$$

where the second line uses Jensen's inequality and the last two lines are based on (67) and (68). The constant  $C > 0$  does not depend on  $t$  or  $N$ .

**Case (IV-b)'  $\tilde{m}_t \leq \rho$ .**

In this case, we can show  $\tilde{\sigma}_t^2 \geq \frac{1}{2D_0}$ . In fact, from  $\tilde{m}_t \leq \rho$ , we have

$$\tilde{\sigma}_t^2 = \sigma_t^2 - \tilde{m}_t^2 \sigma_{t_*}^2 \geq \sigma_t^2 - \rho^2 \sigma_{t_*}^2.$$

From Assumption (A3),  $m_t^2 + \sigma_t^2 \geq D_0^{-1}$ . Since  $m_t^2 \leq \rho^2 m_{t_*}^2$  by assumption, we have

$$\sigma_t^2 \geq D_0^{-1} - m_t^2 \geq D_0^{-1} - \rho^2 m_{t_*}^2.$$

Combining these two inequalities, we obtain

$$\tilde{\sigma}_t^2 \geq D_0^{-1} - \rho^2 (m_{t_*}^2 + \rho_{t_*}^2) \geq D_0^{-1} - \rho^2 D_0 = \frac{1}{2D_0},$$

where the last equality holds from the definition  $\rho = \frac{1}{\sqrt{2D_0}}$ .

We divide the integral of  $\tilde{I}_{B,2}$  into the regions  $\{\mathbf{y} \mid \|\mathbf{y}\|_\infty \geq C_5\sqrt{\log N}\}$  and  $\{\mathbf{y} \mid \|\mathbf{y}\|_\infty \leq C_5\sqrt{\log N}\}$ . In the region  $\{\mathbf{y} \mid \|\mathbf{y}\|_\infty \geq C_5\sqrt{\log N}\}$ , using  $\tilde{\sigma}_t^2 \geq 1/(2D_0)$  and  $f_{N^*}(\mathbf{y}) = 0$ , we have a bound

$$\begin{aligned}
&\left\| \int_{\{\|\mathbf{y}\|_\infty \geq C_5\sqrt{\log N}\}} \frac{\mathbf{x} - \tilde{m}_t \mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi\tilde{\sigma}_t})^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t_*}(\mathbf{y})) d\mathbf{y} \right\|^2 \\
&\leq \left( \frac{D_0}{\pi} \right)^d (2D_0)^2 \int_{\{\|\mathbf{y}\|_\infty \geq C_5\sqrt{\log N}\}} \|\mathbf{x} - \tilde{m}_t \mathbf{y}\|^2 (p_{t_*}(\mathbf{y}))^2 d\mathbf{y} \\
&\leq C \left( \frac{D_0}{\pi} \right)^d (2D_0)^2 \int_{\{\|\mathbf{y}\|_\infty \geq C_5\sqrt{\log N}\}} (C' \log N + \rho^{-2} \|\mathbf{y}\|^2) p_{t_*}(\mathbf{y}) d\mathbf{y} \\
&\leq C'' N^{-3\eta} \log N,
\end{aligned}$$

where we use (67) and the fact  $\mathbf{x} \in D_{t,N}$  implies  $\|\mathbf{x}\|^2 \leq C' \log N$  for some  $C'$ .

For the other region  $\{\mathbf{y} \mid \|\mathbf{y}\|_\infty \leq C_5\sqrt{\log N}\}$ , first notice that for  $\mathbf{x} \in D_{t,N}$ , we have  $\|\mathbf{x} - \tilde{m}_t\mathbf{y}\|/\tilde{\sigma}_t \leq D'\sqrt{\log N}$  for some  $D' > 0$ . Application of Cauchy-Schwarz inequality derives

$$\begin{aligned} & \left\| \int_{\{\|\mathbf{y}\|_\infty \leq C_5\sqrt{\log N}\}} \frac{\mathbf{x} - \tilde{m}_t\mathbf{y}}{\tilde{\sigma}_t} \frac{1}{(\sqrt{2\pi}\tilde{\sigma}_t)^d} e^{-\frac{\|\mathbf{x} - \tilde{m}_t\mathbf{y}\|^2}{2\tilde{\sigma}_t^2}} (f_{N^*}(\mathbf{y}) - p_{t^*}(\mathbf{y})) d\mathbf{y} \right\|^2 \\ & \leq D'^2 \log N \left( \frac{D_0}{\pi} \right) \int_{\{\|\mathbf{y}\|_\infty \leq C_5\sqrt{\log N}\}} d\mathbf{y} \int_{\{\|\mathbf{y}\|_\infty \leq C_5\sqrt{\log N}\}} (f_{N^*}(\mathbf{y}) - p_{t^*}(\mathbf{y}))^2 d\mathbf{y} \\ & \leq D'' (\log N)^{\frac{d}{2}+1} \|f_{N^*} - p_{t^*}\|_{L^2([-C_5\sqrt{\log N}^2, C_5\sqrt{\log N}]^d)} \\ & \leq D'' (\log N)^{\frac{d}{2}+1} N^{-3\eta}. \end{aligned}$$

From the above two cases (IV-a)' and (IV-b)', we have for any  $t \in [t_*, 1]$

$$\tilde{J}_{B,2} \leq \text{poly}(\log N) N^{-\eta}.$$

As a consequence, there is a constant  $C''$  that does not depend on  $t$  and  $m \in \mathbb{N}$  such that

$$\tilde{I}_B \leq C'' \{(\sigma'_t)^2 \log N + (m'_t)^2\} N^{-\eta} \text{poly}(\log N).$$

The factor  $\text{poly}(\log N)$  can be erased if we take a larger  $\eta$  in the proof.  $\square$

## D APPROXIMATION OF FUNCTIONAL OPERATIONS BY NEURAL NETWORKS

This section reviews the accuracy of the approximation and the increase in complexity when we approximate functional operations by neural networks. The following results are shown in Oko et al. (2023, Section F) as well and in more original literature Nakada and Imaizumi (2020), Petersen and Voigtlaender (2018), Schmidt-Hieber (2019).

The operations used directly in this article are `recip`, `mult`, `clip`, and `sw`. The usage in Section C.4 (III) is explained as examples.

### D.1 CLIPPING FUNCTION

First, we consider the realization of the component-wise clipping function.

**Lemma 20.** *For any  $a, b \in \mathbb{R}^d$  with  $a_i \leq b_i$  ( $i = 1, 2, \dots, d$ ), there exists a neural network  $\text{clip}(\mathbf{x}; a, b) \in \mathcal{M}(L, W, S, B)$  with  $L = 2$ ,  $W = (d, 2d, d)^T$ ,  $S = 7d$ , and  $B = \max_{1 \leq i \leq d} \max\{|a_i|, b_i\}$  such that*

$$\text{clip}(\mathbf{x}; a, b)_i = \min\{b_i, \max\{x_i, a_i\}\} \quad (i = 1, 2, \dots, d)$$

*holds. When  $a_i = c_{\min}$  and  $b_i = c_{\max}$  for all  $i$ , we also use the notation  $\text{clip}(\mathbf{x}; c_{\min}, c_{\max}) := \text{clip}(\mathbf{x}; a, b)$ .*

### D.2 RECIPROCAL FUNCTION

Second, the reciprocal function  $x \mapsto 1/x$  is approximated by neural networks as follows.

**Lemma 21.** *For any  $0 < \varepsilon < 1$ , there is  $\text{recip}(\mathbf{x}') \in \mathcal{M}(L, W, S, B)$  such that*

$$\left| \text{recip}(\mathbf{x}') - \frac{1}{x} \right| \leq \varepsilon + \frac{|\mathbf{x} - \mathbf{x}'|}{\varepsilon^2} \quad (73)$$

*holds for any  $x \in [\varepsilon, \varepsilon^{-1}]$  and  $\mathbf{x}' \in \mathbb{R}$  with  $L = O(\log^2(\varepsilon^{-1}))$ ,  $\|W\|_\infty = O(\log^3(\varepsilon^{-1}))$ ,  $S = O(\log^4(\varepsilon^{-1}))$ , and  $B = O(\varepsilon^{-2})$ .*

### D.3 MULTIPLICATION

**Lemma 22.** *Let  $d \geq 2$ ,  $C \geq 1$ ,  $0 < \epsilon_{\text{err}} \leq 1$ . For any  $\varepsilon > 0$ , there exists a neural network  $\text{mult}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d) \in \mathcal{M}(L, W, S, B)$  with  $L = O(d \log \varepsilon^{-1} + d \log C)$ ,  $\|W\|_\infty = 48d$ ,  $S =$*

1620  $O(d \log \varepsilon^{-1} + d \log C)$ ,  $B = C^d$  such that (i)

$$1621 \left| \text{mult}(x'_1, \dots, x'_d) - \prod_{d'=1}^d x'_{d'} \right| \leq \varepsilon + dC^{d-1}\epsilon_{err},$$

1622 holds for all  $x \in [-C, C]^d$  and  $\mathbf{x}' \in \mathbb{R}^d$  with  $\|\mathbf{x} - \mathbf{x}'\|_\infty \leq \epsilon_{err}$ , (ii)  $|\text{mult}(\mathbf{x})| \leq C^d$  for all  $x \in \mathbb{R}^d$ ,  
1623 and (iii)  $\text{mult}(x'_1, \dots, x'_d) = 0$  if at least one of  $x'_i$  is 0.

1624 We note that some of  $x_i, x_j$  ( $i \neq j$ ) can be shared; for  $\prod_{i=1}^I x_{\alpha_i}$  with  $\alpha_i \in \mathbb{Z}_+$  ( $i = 1, \dots, I$ ) and  
1625  $\sum_{i=1}^I \alpha_i = d$ , there exists a neural network satisfying the same bounds as above; the network is  
1626 denoted by  $\text{mult}(\mathbf{x}; \alpha)$ .

#### 1631 D.4 SWITCHING

1632 **Lemma 23.** Let  $t_1 < t_2 < s_1 < s_2$ , and  $f(\mathbf{x}, t)$  be a scalar-valued function. Assume that  
1633  $|\varphi_1(\mathbf{x}, t) - f(\mathbf{x}, t)| \leq \varepsilon$  on  $[t_1, s_1]$  and  $|\varphi_2(\mathbf{x}, t) - f(\mathbf{x}, t)| \leq \varepsilon$  on  $[t_2, s_2]$ . Then, there exist neural  
1634 networks  $\text{sw}_1(t; t_2, s_1)$  and  $\text{sw}_2(t; t_2, s_1)$  in  $\mathcal{M}(L, W, S, B)$  with  $L = 3$ ,  $W = (1, 2, 1, 1)^T$ ,  $S = 8$ ,  
1635 and  $B = \max\{t_1, (s_1 - t_2)^{-1}\}$  such that

$$1636 |\text{sw}_1(t; t_2, s_1)\varphi_1(\mathbf{x}, t) + \text{sw}_2(t; t_2, s_1)\varphi_2(\mathbf{x}, t) - f(\mathbf{x}, t)| \leq \varepsilon$$

1637 holds for any  $t \in [t_1, s_2]$ .

#### 1641 D.5 CONSTRUCTION OF NETWORK IN SECTION C.4 (III)

1642 We detail the network size required for the approximation procedure presented in Section C.4 (III).

1643 **Clipping to  $\zeta_1$ :** From (39) and assumption (A3),  $\phi_5$  can be upper bounded by  $N^{K_0+1}$  for sufficiently  
1644 large  $N$ . Then, from Lemma 20, we can see that in the clipping of  $\phi_5$ , the increase in model sizes is  
1645 constant depending on  $d$  except  $B$ , which is multiplied by the upper bound  $N^{K_0+1}$ .

1646 **Approximating  $f_1^{-1}$  by  $\zeta_2$ :** Since we assume  $f_1 \geq N^{-(2s+\omega)/d}$ , clipping  $\phi_5$  from below by  
1647  $N^{-(2s+\omega)/d}$  does not increase the difference; thus  $|\zeta_1 - f_1| \leq D_5 N^{-\eta}$ . In Lemma 21, substituting  
1648  $x' = \zeta_1$ ,  $x = f_1$ , and  $\varepsilon = N^{-\chi}$  for  $\chi > \chi_0 + (2s + \omega)/d$  with an arbitrary  $\chi_0 > 0$ , we have

$$1649 \left| \text{recip}(\zeta_1(\mathbf{x}, t)) - \frac{1}{f_1(\mathbf{x}, t)} \right| \leq N^{-\chi} + N^{2\chi} |\zeta_1(\mathbf{x}, t) - f_1(\mathbf{x}, t)|.$$

1650 Since  $\eta > 0$  is arbitrary, by setting  $\eta$  and  $\chi$  so that  $\eta > 3\chi$ , we have

$$1651 \left| \text{recip}(\zeta_1(\mathbf{x}, t)) - \frac{1}{f_1(\mathbf{x}, t)} \right| \leq (D_5 + 1)N^{-\chi}. \quad (74)$$

1652 This is achieved by a neural network with  $L = O(\log^2 N)$ ,  $S = O(\log^4 N)$ ,  $\|W\|_\infty = O(\log \log N)$   
1653 and  $B = O(N^{2\chi})$ .

1654  **$\zeta_3 = \text{mult}(\zeta_2, \phi_6)$ :** Note that  $|\zeta_2| \leq N^{(2s+\omega)/d}$ ,  $\|\mathbf{f}_2 - \phi_6\| = O(N^{-\eta})$ , and  $|\zeta_2 - f_1^{-1}| \leq$   
1655  $O(N^{-\chi})$  from (74). We have also taken  $\eta$  so that  $\eta > 3\chi$ . In applying Lemma 22, we can set  
1656  $C = N^{(2s+\omega)/d}$  because  $|\zeta_2| \leq N^{(2s+\omega)/d}$ . Also,  $\epsilon_{err} := \max\{|\zeta_2 - 1/f_1|, \|\phi_6 - \mathbf{f}_2\|\} =$   
1657  $\max\{O(N^{-\chi}), O(N^{-\eta})\} = O(N^{-\chi})$ . With  $d = 2$  and  $\varepsilon = N^{-\chi_0}$ , we have

$$1658 |\zeta_3(\mathbf{x}, t) - \zeta_2 \cdot \phi_6(\mathbf{x}, t)| = O(N^{-\chi_0} + 2N^{(2s+\omega)/d}N^{-\chi}) = O(N^{-\chi_0} + 2N^{-\chi_0}) = O(N^{-\chi_0}),$$

1659 where we use the fact that  $\chi$  is taken to satisfy  $\chi > \chi_0 + (2s + \omega)/d$ .

1660 We then obtain

$$1661 \left\| \zeta_3 - \frac{\mathbf{f}_2}{f_1} \right\| \leq \left\| \zeta_3 - \phi_6 \zeta_2 \right\| + \left\| \phi_6 \zeta_2 - \frac{\mathbf{f}_2}{f_1} \right\|$$

$$1662 \leq \left\| \zeta_3 - \phi_6 \zeta_2 \right\| + \left\| \phi_6 \zeta_2 - \mathbf{f}_2 \zeta_2 \right\| + \left\| \mathbf{f}_2 \zeta_2 - \frac{\mathbf{f}_2}{f_1} \right\|$$

$$1663 = O(N^{-\chi_0}) + O(N^{(2s+\omega)/d}N^{-\eta}) + O(\sqrt{\log N}N^{-\chi})$$

$$1664 = O(N^{-\chi_0}), \quad (75)$$

1674 where in the second last line we use  $\|\mathbf{f}_2/f_1\| = O(\sqrt{\log N})$  and thus  $\|\mathbf{f}_2\| = O(\sqrt{\log N})$ .

1675 A similar argument derives

$$1676 \left\| \zeta_5 - \frac{\mathbf{f}_3}{f_1} \right\| = O(N^{-\chi_0}). \quad (76)$$

1677 For the neural network architecture in this multiplication,  $L$  and  $S$  are added by the order of  
1678  $O(\log \varepsilon^{-1} + \log C) = O(\log N)$ . The width  $W$  is of constant order and thus negligible. The width  
1679  $W$  is  $O(N^{2(2s+\omega)/d})$ .

1680 **Clipping to obtain  $\zeta_4$  and  $\zeta_6$ :** For these clipping procedures, the approximating networks have  $L$ ,  
1681  $W$  and  $S$  of constant order, while the weight values  $B$  are of  $O(\sqrt{\log N})$  and  $O(1)$ , respectively, and  
1682 thus they are negligible. The approximation errors are kept as  $N^{-\chi_0}$ .

1683 **Multiplication to obtain  $\zeta_7$  and  $\zeta_8$ :** As in the previous procedures, clipping by  $O(\sqrt{\log N})$  and  
1684  $O(1)$  does not increase the approximation error, while the increase in the size of the network is  
1685 negligible.

1686 In a similar manner to Oko et al. (Lemma B.1 2023), we can approximate  $\sigma'_t$  and  $m'_t$  by neural  
1687 networks so that  $|\sigma'_t - \widehat{\sigma}'_t| = O(N^{-\eta})$  and  $|m'_t - \widehat{m}'_t| = O(N^{-\eta})$ . The network sizes are  $L =$   
1688  $O(\log^2 N)$ ,  $\|W\|_\infty = O(\log^2 N)$ ,  $S = O(\log^3 N)$ , and  $B = O(\log N)$ . With arguments similar to  
1689 those of the previous procedures, we can show

$$1690 \left\| \zeta_7 - (\sigma'_t) \frac{\mathbf{f}_2}{f_1} \mathbf{1} \left[ \left\| \frac{\mathbf{f}_2}{f_1} \right\| \leq C_5 \sqrt{\log N} \right] \right\| = O(N^{-\chi_0}),$$

1691 and

$$1692 \left\| \zeta_8 - |m'_t| \frac{\mathbf{f}_3}{f_1} \mathbf{1} \left[ \left\| \frac{\mathbf{f}_3}{f_1} \right\| \leq C_5 \right] \right\| = O(N^{-\chi_0}),$$

1693 In total, we can find a neural network to approximate  $\mathbf{v}_t(\mathbf{x})$  with the approximation error of  $O(N^{-\chi_0})$   
1694 so that the network has the size of most polynomial orders for  $B$  and  $\|W\|_\infty$ , while  $O(\text{poly}(\log N))$   
1695 for  $S$  and  $L$ . As a result, the contributions to the log cover number are only  $O(\text{poly}(\log N))$ .

## 1700 E IDEA OF TIME DIVISION

1701 The basic idea of the time division used to derive the almost optimal minimax convergence rate is  
1702 depicted in Figure E.1.



1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

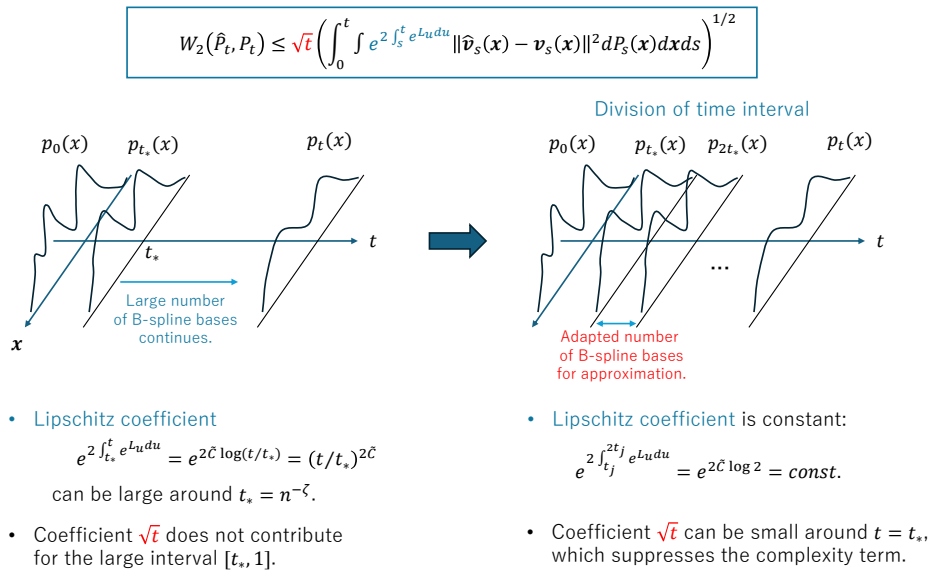


Figure E.1: The idea of time division for deriving the convergence rate