

# ERROR ANALYSIS OF DISCRETE FLOW WITH GENERATOR MATCHING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Discrete flow models offer a powerful framework for learning distributions over discrete state spaces and have demonstrated superior performance compared to the discrete diffusion model. However, their convergence properties and error analysis remain largely unexplored. In this work, we develop a unified framework grounded in stochastic calculus theory to systematically investigate the theoretical properties of discrete flow. Specifically, we derive the KL divergence of two path measures regarding two continuous-time Markov chains (CTMCs) with different transition rates by deriving a Girsanov-type theorem, and provide a comprehensive analysis that encompasses the error arising from transition rate estimation and early stopping, where the first type of error has rarely been analyzed by existing works. Unlike discrete diffusion models, discrete flow incurs no truncation error caused by truncating the time horizon in the noising process. Building on generator matching and uniformization, we establish non-asymptotic error bounds for distribution estimation. Our results provide the first error analysis for discrete flow models.

## 1 INTRODUCTION

Discrete diffusion models have achieved significant progress in large language models (Nie et al., 2025; Zhu et al., 2025; Zhao et al., 2025; Yang et al., 2025). By learning the time reversal of the noising process of a continuous-time Markov chain (CTMC), the models transform a simple distribution (e.g., uniform (Hoogeboom et al., 2021; Lou et al., 2024) and masked (Ou et al., 2025; Shi et al., 2024; Sahoo et al., 2024)) that is easy to sample to the data distribution that has discrete structures. Discrete flow models Campbell et al. (2024); Gat et al. (2024); Shaul et al. (2025) provide a flexible framework for learning generating transition rate analogous to continuous flow matching (Albergo & Vanden-Eijnden, 2023; Liu et al., 2023; Lipman et al., 2023), offering a more comprehensive family of probability paths.

Recent theoretical analysis for discrete diffusion models has emerged through numerous studies (Chen & Ying, 2024; Zhang et al., 2025; Ren et al., 2025a;b). To obtain the transition rate in the reversed process, the concrete scores in these analyses are obtained by minimizing the concrete score entropy introduced in Lou et al. (2024); Benton et al. (2024c). In those works, the distribution errors of discrete diffusion models are divided into three parts: (a) truncation error from truncating the time horizon in the noising process; (b) concrete score estimation error; (c) discretization error from sampling algorithms. In our paper, we aim to investigate the theoretical properties of the discrete flow-based models using the generator matching training objective (Holderrieth et al., 2025) and the uniformization sampling algorithm (Chen & Ying, 2024), which offers zero truncation error and discretization error. Our analysis takes transition rate estimation error into account instead of imposing a stringent condition on it in previous works, which is related to the early stopping parameter. We decompose the estimation error into stochastic error and approximation error, and then balance the stochastic error and early stopping error by choosing the early stopping parameter. Our stochastic error bound is aligned with the SOTA result in continuous flow (Gao et al., 2024b); the early stopping error bound matches the most recent result in discrete diffusion (Zhang et al., 2025). Furthermore, we present a comprehensive error analysis for the neural network class with the ReLU activation function, by controlling stochastic error, approximation error and early stopping error simultaneously.

The main contributions in this paper are summarized as follows.

1. We obtain a Girsanov-type theorem for CTMC with rigorous proofs. The KL divergence of two path measures of two CTMCs is derived in terms of the integral of the Bregman divergence of two transition rates, which motivates us to use the generator matching objective to analyze the distribution error. In the sampling stage, we use the uniformization technique to sample in an exact way.
2. We establish the non-asymptotic error bound for distribution error in discrete flow models, by taking estimation error into consideration, which is less explored in existing works. There are three sources of error in our framework: (a) stochastic error from estimation through empirical risk minimization; (b) approximation error of the selected function class; (c) early stopping error. We carefully analyze the stochastic error using empirical process theory, and discuss the choice of early stopping parameter to balance the stochastic error and early stopping error. We also provide a comprehensive analysis of our assumptions and the terminal time singularity of the transition rate. In addition, we present a thorough error analysis for the neural network class with the ReLU activation function, simultaneously analyzing these three sources of error. To the best of our knowledge, this is the first theoretical error analysis for discrete flow models.

All technical proofs are deferred to the Appendix.

NOTATION. Let  $[N] = \{1, 2, \dots, N\}$  for a positive integer  $N$ . We use  $\dot{\kappa}_t$  to denote the time derivative of a function  $\kappa_t$  of  $t$ . For a  $\mathcal{D}$ -dimensional vector  $z$ , let  $z^d$  and  $z^{\setminus d}$  denote the  $d$ -th element of the vector  $z$  and the  $(\mathcal{D} - 1)$ -dimensional vector  $(z^1, \dots, z^{d-1}, z^{d+1}, \dots, z^{\mathcal{D}})^\top$ . We use  $\mathbb{1}(\cdot)$  to denote an indicator function. We denote the Hamming distance of two vectors  $z, x$  by  $d^H(z, x)$ . For two quantities  $x, z$ , define the Kronecker delta  $\delta_x(z)$  satisfying  $\delta_x(z) = 1$  if  $x = z$  and  $\delta_x(z) = 0$  if  $x \neq z$ . For a random data  $\mathbb{D}_n = \{Z_i\}_{i \in [n]}$ , we denote  $\|f(Z)\|_{L^\infty(\mathbb{P}_n)} = \max_{i \in [n]} |f(Z_i)|$ , where  $\mathbb{P}_n$  is the empirical measure of  $\mathbb{D}_n$ . In this paper, some universal constants  $C, c > 0$  are allowed to vary from line to line.

## 2 RELATED WORKS

**Discrete Flow.** Discrete flow models (Campbell et al., 2024; Gat et al., 2024; Shaul et al., 2025) provide a more flexible framework to construct the transition rate and probability path than discrete diffusion models (Campbell et al., 2022; Lou et al., 2024), which also achieve superior performance in graph generation (Qin et al., 2025), visual generation and multimodal understanding (Wang et al., 2025). The training objective of discrete flow models can be cross-entropy (Campbell et al., 2024; Gat et al., 2024), negative ELBO (Shaul et al., 2025), or Bregman divergence (Holderrieth et al., 2025). In this work, we use the estimator through minimizing the empirical version of Bregman divergence to control the distribution error.

**Theoretical Analysis of Discrete Diffusion.** Discrete diffusion models (Austin et al., 2021; Campbell et al., 2022; Vignac et al., 2023; Sun et al., 2023; Lou et al., 2024; Benton et al., 2024c) have emerged as a CTMC-based framework for learning distribution on finite state spaces. Some recent works (Chen & Ying, 2024; Zhang et al., 2025; Ren et al., 2025a;b) investigate the theoretical property of discrete diffusion models. Chen & Ying (2024) proposed to use uniformization algorithm for sampling in an exact way, and derived the distribution error bound for the scenario where the state space is a hypercube. Zhang et al. (2025) studied the theoretical results of discrete diffusion based on Girsanov-based method in the general state space  $\mathcal{S}^{\mathcal{D}}$ . Ren et al. (2025a) derived the error bound for both  $\tau$ -leaping and uniformization algorithms using stochastic integrals. Ren et al. (2025b) proposed high-order solvers for sampling and derived error bound similar to Ren et al. (2025a), which enjoy more accuracy than  $\tau$ -leaping algorithm. The training objective in these literature is the concrete score entropy introduced in Lou et al. (2024); Benton et al. (2024c). Unfortunately, in existing works, it is typical to assume strong regularity conditions directly on the estimation error, which is related to the early stopping parameter. Furthermore, in discrete diffusion models, there is a non-zero truncation error arising from the truncation of the time horizon in the noising process.

**Theoretical Analysis of Continuous Flow.** Continuous flow matching (Liu et al., 2023; Albergo & Vanden-Eijnden, 2023; Lipman et al., 2023) is a powerful simulation-free method for learning continuous data distribution based on continuous normalizing flow (Chen et al., 2018). Error bounds on the Wasserstein distance between two flow ODEs have been extensively studied (Albergo & Vanden-Eijnden, 2023; Benton et al., 2024b; Gao et al., 2024a;b). Albergo & Vanden-Eijnden (2023) and

Benton et al. (2024b) derived the Wasserstein bound in terms of the spatial Lipschitz constants of the velocity fields by utilizing Grönwall’s inequality. Gao et al. (2024a) established the spatial Lipschitz regularity of the velocity field for a range of target distributions. Gao et al. (2024b) presented a comprehensive error analysis of continuous flow matching by rigorously bounding four types of errors, including early stopping error, discretization error, stochastic error and approximation error.

### 3 THEORETICAL BACKGROUND ON DISCRETE FLOW-BASED MODELS

In this section, we introduce some theoretical background about continuous-time Markov chains (Norris, 1998) and discrete flow-based models (Campbell et al., 2024; Gat et al., 2024; Shaul et al., 2025; Wang et al., 2025).

#### 3.1 CONTINUOUS-TIME MARKOV CHAIN AS STOCHASTIC INTEGRAL

In this subsection, we consider the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . We formally define continuous-time Markov chains (CTMCs) as follows.

**Definition 1** (CTMC). *Consider a  $\mathcal{D}$ -dimensional finite state space  $\mathcal{S}^{\mathcal{D}}$ , where  $\mathcal{S} = \{1, 2, \dots, |\mathcal{S}|\}$ . Let  $u_t(z, x)_{z, x \in \mathcal{S}^{\mathcal{D}}}$  be a (time-dependent) rate matrix that satisfying: (a)  $u_t(z, x)$  is continuous in  $t$ ; (b)  $u_t(z, x) \geq 0$  for any  $z \neq x$ ; (c)  $\sum_{z \in \mathcal{S}^{\mathcal{D}}} u_t(z, x) = 0$  for any  $x \in \mathcal{S}^{\mathcal{D}}$ . We call a process  $\{X(t)\}_{t \geq 0}$  a continuous time Markov chain with the transition rate matrix  $u_t(z, x)_{z, x \in \mathcal{S}^{\mathcal{D}}}$  and the natural filtration  $\mathcal{F}_t = \sigma(\{X(s) : 0 \leq s \leq t\})$  if it satisfies that for any  $z, x \in \mathcal{S}^{\mathcal{D}}$*

1. *the transition rate is the generator of CTMC:  $\mathbb{P}(X(t+h) = z | X(t) = x) = \delta_x(z) + u_t(z, x)h + o(h)$ ;*
2. *it has Markov property:  $\mathbb{P}(X(t+h) = z | \mathcal{F}_t) = \mathbb{P}(X(t+h) = z | X(t))$ .*

Given a rate matrix  $u_t$  with uniformly bounded entries, the following uniformization technique allows us to construct a CTMC with  $u_t$  through a Poisson process, which offers us a sampling algorithm without discretization error (see, e.g., Ren et al., 2025a; Chen & Ying, 2024).

**Proposition 1** (Uniformization). *Assume that  $u_t(z, x)_{z, x \in \mathcal{S}^{\mathcal{D}}}$  is a rate matrix satisfying that (a)  $-u_t(x, x) \leq M$  for any  $x \in \mathcal{S}^{\mathcal{D}}$ ; (b)  $u_t(z, x)$  is  $L$ -Lipschitz continuous in  $t$  for any pair  $(z, x) \in \mathcal{S}^{\mathcal{D}} \times \mathcal{S}^{\mathcal{D}}$ . Suppose that  $T_1, T_2, \dots$  are the arrival times of a Poisson process  $N(t)$  with rate  $M$ . Let  $X(t)$  be a jump process with initial distribution  $p_0$  and natural filtration  $\mathcal{F}_t$  such that at  $t = T_1, T_2, \dots$ , the process  $X(t)$  jumps to position  $z \neq X(t-)$  with probability  $u_t(z, X(t-))/M$ . Then  $X(t)$  is a Markov process satisfying for  $z \neq x$ ,*

$$\mathbb{P}(X(t+h) = z | X(t) = x) = u_t(z, x) + R_t,$$

where  $R_t \leq (M^2 + L)h^2 = O(h^2)$ . Thus,  $X(t)$  is a CTMC with the rate  $u_t$ .

Here, the remainder  $R_t$  is uniformly bounded for any  $t$  and  $z \neq x$  under the assumptions in Proposition 1, which is crucial for developing the theory of CTMC.

To rewrite a CTMC  $X(t)$  as a stochastic integral, we define a random measure  $N((t_1, t_2], A)$  associated with the CTMC  $X(t)$ .

**Definition 2** (Random Measure Associated with CTMC). *Suppose that  $X(t)$  is a CTMC with rate  $u_t$  and natural filtration  $\mathcal{F}_t$ . Define the random measure associated with  $X(t)$  and  $A \subseteq \mathcal{S}^{\mathcal{D}}$  as*

$$N((t_1, t_2], A) = \#\{t_1 < s \leq t_2; \Delta X(s) \neq 0, X(s) \in A\},$$

where  $\Delta X(t) = X(t) - X(t-)$ .

Suppose that  $N((t_1, t_2], A)$  is the random measure associated with a CTMC  $X(t)$ , where  $X(t)$  is constructed by uniformization with the Poisson process  $N(t)$ . Then  $N(t, A) \triangleq N((0, t], A) \leq N(t) < \infty$  a.s. for any  $t \geq 0$  and  $A \subseteq \mathcal{S}^{\mathcal{D}}$  (see Lemma 2.3.4 in Applebaum, 2009). Thus, similar to the Poisson random measure associated with a Lévy process (see Section 2.3.2 in Applebaum, 2009), there are some equivalent representations:

$$N((t_1, t_2], A) = \sum_{t_1 < s \leq t_2} \mathbb{1}_A((X(s))) (1 - \delta_{X(s-)}(X(s))) = \sum_{n \in \mathbb{N}} \delta_{(t_1, t_2]}(T_n^A),$$

where  $\{T_n^A\}_{n \in \mathbb{N}}$  are arrival times of the counting process  $N(t, A)$ . Consequently,  $X(t)$  can be written as the following stochastic integral:

$$\begin{aligned} X(t) &= X(0) + \sum_{n \in \mathbb{N}} \left[ X(T_n) - X(T_{n-1}) \right] \delta_{[0,t]}(T_n) \\ &= X(0) + \sum_{n \in \mathbb{N}} \sum_z \left[ z - X(T_{n-}) \right] \mathbb{1}_z(X(T_n)) \delta_{[0,t]}(T_n) \\ &= X(0) + \int_0^t \int_{\mathcal{S}^{\mathcal{D}}} (z - X(s-)) N(ds, dz), \end{aligned} \quad (1)$$

where  $\{T_n\}_{n \in \mathbb{N}}$  are arrival times of  $N(t, \mathcal{S}^{\mathcal{D}})$ . Here, we formally define the integrator  $N(t, A)$  in the CTMC representation as Definition 2, which is clearer and more formal than the relevant definition compared to Proposition 3.2 in Ren et al. (2025a).

**Remark 1** (Comparison to Lévy-Itô Decomposition). *Recall that a Lévy process  $X(t)$  has the following Lévy-Itô decomposition (see Theorem 2.4.16 in Applebaum, 2009):*

$$X(t) = X(0) + bt + B(t) + \int_{|x| < 1} x \tilde{N}(t, dx) + \int_{|x| \geq 1} x N(t, dx),$$

where  $B(t)$  is a Brownian motion,  $N(t, A)$  is the Poisson random measure associated with the Lévy process  $X(t)$ , and  $\tilde{N}(t, A)$  is the compensator of  $N(t, A)$ . Here, the second argument in  $N(t, A)$  is the jump size of the associated Lévy process at time  $t$ . However, the second argument of the random measure associated with a CTMC in Definition 2 is the position after jump at time  $t$ . Both the integrands in equation 1 and that in Lévy-Itô decomposition are jump size. Moreover, since CTMCs do not have independent increments, the random measure in Definition 2 is not independently scattered; that is,  $N(t, A_1)$  and  $N(t, A_2)$  are not necessarily independent for disjoint  $A_1, A_2 \subseteq \mathcal{S}^{\mathcal{D}}$ , which is not a Poisson random measure. There is no need to use the independently scattered property in our technical proofs.

Given a CTMC  $X(t)$  with transition rate  $u_t$  and marginal densities  $\{p_t\}_{t \geq 0}$  w.r.t. the counting measure, the following proposition demonstrates that the marginal densities  $p_t$  satisfy the following Kolmogorov forward equation (a.k.a. continuity equation).

**Proposition 2** (Kolmogorov Forward Equation). *The CTMC  $X(t)$  satisfies the following equation:*

$$\dot{p}_t(x) = \sum_{z \in \mathcal{S}^{\mathcal{D}}} u_t(x, z) p_t(z) = \underbrace{\sum_{z \neq x} u_t(x, z) p_t(z)}_{\text{Incoming Flux}} - \underbrace{\sum_{s \neq x} u_t(z, x) p_t(x)}_{\text{Outgoing Flux}}.$$

In our work, we say  $u_t$  can generate the probability path  $p_t$ , if it satisfies the above Kolmogorov equation.

### 3.2 DISCRETE FLOW-BASED MODELS

We aim to learn a transition rate  $u_t$  of a CTMC  $X(t)$  that can transport from a source distribution  $p_0$  to a target data distribution  $p_1$ . To obtain such a transition rate for sampling, a natural method is to learn the conditional expectation of the conditional transition rate  $u_t(z, x|x_1)$  that generates the conditional probability path  $p_{t|1}(\cdot|x_1)$ , since  $u_t(z, x) = \mathbb{E}[u_t(z, x|X(1))|X(\mathbf{t}) = x, \mathbf{t} = t]$  can generate the target probability path  $p_t$  (see Proposition 3.1 of Campbell et al., 2024), i.e., it satisfies the Kolmogorov forward equation, where  $X(\mathbf{t}) \sim p_{\mathbf{t}|1}(\cdot|X(1))$  given  $X(1)$  and  $\mathbf{t} \sim \mathcal{U}([0, 1])$ . Therefore, we are free to define the conditional probability path and the conditional transition rate.

It is worth noting that in the sampling stage, considering a  $|\mathcal{S}|^{\mathcal{D}}$ -dimensional vector-valued function  $(u_t(z, x))_{z \in \mathcal{S}^{\mathcal{D}}}$  of current time  $t$  and state  $x$  is intractable when  $\mathcal{D}$  is relatively large. To handle such a high-dimensional scenario, a common approach is to construct a coordinate-wise conditional probability path and transition rate, that is,

$$p_{t|1}(x|x_1) = \prod_{d=1}^{\mathcal{D}} p_{t|1}^d(x^d|x_1^d), \text{ and } u_t(z, x|x_1) = \sum_{d=1}^{\mathcal{D}} \delta_{x \setminus d}(z \setminus d) u_t^d(z^d, x^d|x_1^d), \quad (2)$$

which means that the elements of the vector  $X(t)$  are independent conditional on  $X(1)$ . Here,  $u_t^d(z^d, x^d|x_1^d)$  is the conditional transition rate that generates the conditional probability path  $p_{t|1}^d$ . A popular choice of probability path and the associated conditional transition rate used in the previous works (Campbell et al., 2024; Gat et al., 2024) is

$$p_{t|1}^d(x^d|x_1^d) = (1 - \kappa_t)p_0^d(x^d) + \kappa_t\delta_{x_1^d}(x^d); u_t^d(z^d, x^d|x_1^d) = \frac{\dot{\kappa}_t}{1 - \kappa_t}(\delta_{x_1^d}(z^d) - \delta_{x^d}(z^d)), \quad (3)$$

where  $\kappa_t : [0, 1] \rightarrow [0, 1]$  is a non-decreasing function satisfying  $\kappa_0 = 0$  and  $\kappa_1 = 1$ . Note that the conditional transition rate will blow up as  $t \rightarrow 1^-$ . Thus, in the sampling stage, we employ the early stopping technique; that is, we only consider the time interval  $[0, 1 - \tau]$  for a sufficiently small positive parameter  $\tau$ . We will discuss this time singularity in Section B.1.

After defining the conditional path and rate, the unconditional transition rate is given by

$$\begin{aligned} u_t(z, x) &= \sum_{x_1} u_t(z, x|x_1)p_{1|t}(x_1|x) = \sum_{d=1}^{\mathcal{D}} \delta_{x \setminus d}(z \setminus d) \sum_{x_1^d} u_t^d(z^d, x^d|x_1^d)p_{1|t}^d(x_1^d|x) \\ &\triangleq \sum_{d=1}^{\mathcal{D}} \delta_{x \setminus d}(z \setminus d) u_t^d(z^d, x), \end{aligned} \quad (4)$$

where  $p_{1|t}^d(x_1^d|x) = \sum_{x \setminus d} p_{1|t}(x_1|x)$ . Therefore, it suffices to consider a sparse rate matrix  $u_t$  satisfying  $u_t(z, x) = 0$  for any  $z, x \in \mathcal{S}^{\mathcal{D}}$  with Hamming distance  $d^H(z, x) > 1$ .

**Remark 2.** *With a slight abuse of notation, throughout this article, the summation over the state space is restricted to states where the Hamming distance is not larger than 1; specifically, we use  $\sum_{z \in A} = \sum_{z \in A: d^H(z, X(t-)) \leq 1}$ .*

In our work, similar to Chen & Ying (2024); Zhang et al. (2025), we consider the *uniform* source distribution. There are three reasons: (a) a uniform distribution has full support, which is similar to the Gaussian distribution used in continuous flow-based models; (b) in any state, each coordinate can possibly be changed in the sampling stage, compared to an absorbing source distribution; (c) it has a good performance in practice (e.g., flow-based multimodal large language models Wang et al., 2025).

**Training via Bregman divergence.** Let  $\tau \in (0, 1/2)$  be an early stopping parameter. Suppose that we have i.i.d. samples  $\{\mathbf{t}_i, X_i(1)\}_{i \in [n]}$ , and  $X_i(\mathbf{t}_i) \sim p_{\mathbf{t}_i|1}(\cdot|X_i(1))$  for each  $i \in [n]$ , where  $\mathbf{t}_i$  samples from the uniform distribution  $\mathcal{U}([0, 1 - \tau])$ . Denote  $\mathbb{D}_n = \{Z_i\}_{i \in [n]} = \{(\mathbf{t}_i, X_i(\mathbf{t}_i), X_i(1))\}_{i \in [n]}$  and  $v(x, Z_i) = u_{\mathbf{t}_i}(x, X_i(\mathbf{t}_i)|X_i(1))$ , where  $u_t(z, x|x_1)$  is a conditional transition rate that generates the conditional probability path  $p_{t|1}(x|x_1)$ . Denote the Bregman divergence with a convex function  $F(\cdot)$  as  $D_F$ , which is defined as

$$D_F(a||b) = F(a) - F(b) - \langle \nabla F(b), a - b \rangle.$$

Inspired by generator matching Holderrieth et al. (2025), we consider the following rate estimator through empirical risk minimization (ERM) with a function class  $\mathcal{G}_n$ :

$$\begin{aligned} \hat{u} &= \arg \min_{u \in \mathcal{G}_n} \frac{1}{n} \sum_{i=1}^n \sum_{z \neq X_i(\mathbf{t}_i)} D_F(v(z, Z_i)||u_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i))) \\ &= \arg \min_{u \in \mathcal{G}_n} \frac{1}{n} \sum_{i=1}^n \sum_{z \neq X_i(\mathbf{t}_i)} \left\{ -u_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i)|X_i(1)) \log u_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i)) + u_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i)) \right\} \quad (5) \\ &\triangleq \arg \min_{u \in \mathcal{G}_n} \mathcal{L}_n(u), \end{aligned}$$

where we take  $F(x) = x \log x$  and the summation over  $\{z : z \neq X_i(\mathbf{t}_i)\}$  has only  $O(\mathcal{D}|S|)$  complexity due to the sparsity of the rate matrix (Remark 2). Here, the training objective is finite if we choose an appropriate function class satisfying Assumption 2 in Section 5.

**Sampling via uniformization.** For accurate sampling without discretization error, following the uniformization argument introduced in Proposition 1, we consider the uniformization Algorithm 1 in the Appendix (also see, e.g. Algorithm 1 in Chen & Ying, 2024). Here, the discretization used in Algorithm 1 (see Section A in the Appendix) can reduce the average sampling steps compared to using a uniform bound  $M \geq \sup_{t \in [0, 1 - \tau]} \sum_{z \neq x} \dot{u}_t(z, x)$  for a large time interval  $[0, 1 - \tau]$ .

## 4 CHANGE OF MEASURE AND BOUND FOR KL DIVERGENCE

In this section, we develop a Girsanov-type theorem for CTMCs. As a result, we establish the bound for the KL divergence of two marginal distributions regarding two CTMCs with different rates.

First, we derive the compensator of the random measure  $N(t, A)$  defined in Definition 2, where  $A \subseteq \mathcal{S}^{\mathcal{D}}$ .

**Proposition 3** (Compensator). *Under the assumptions in Proposition 1, the (martingale-valued) compensated random measure of  $N(t, A)$  is*

$$\tilde{N}(t, A) = N(t, A) - \int_0^t \sum_{z \in A} u_s(z, X(s-)) \mathbb{1}(z \neq X(s-)) ds.$$

### 4.1 CHANGE OF MEASURE

Suppose that  $X(t)$  is a CTMC with transition rate  $u^X$  and natural filtration  $\mathcal{F}_t$  under probability space  $(\Omega, \mathcal{F}, \mathbb{P}^X)$ . Our goal is to find a probability measure  $\mathbb{P}^Y \ll \mathbb{P}^X$  such that  $X(t)$  is a  $\mathbb{P}^Y$ -CTMC with a rate  $u^Y$ . Denote  $N(t, A)$  as the random measure associated with  $X(t)$ . Let  $\exp(W(t))$  be an exponential  $\mathbb{P}^X$ -martingale (see Lemma 3 in Appendix), where  $W(t)$  has the following form:

$$W(t) = W(0) + \int_0^t \sum_{x \neq X(s-)} F(s, x) ds + \int_0^t \int_{x \in \mathcal{S}^{\mathcal{D}}} K(s, x) N(ds, dx). \quad (6)$$

Since  $\exp(W(t))$  is a  $\mathbb{P}^X$ -martingale with mean one, we can define  $\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t = \frac{d\mathbb{P}_t^Y}{d\mathbb{P}_t^X} = \exp(W(t))$ , where  $\mathbb{P}_t$  is the restriction of the measure  $\mathbb{P}$  on  $(\Omega, \mathcal{F}_t)$ . If we take some specific predictable processes  $F$  and  $K$ , the following theorem shows that

$$\tilde{N}^Y(t, A) \triangleq N(t, A) - \int_0^t \sum_{z \in A} u_s^Y(z, X(s-)) \mathbb{1}(z \neq X(s-)) ds$$

is a  $\mathbb{P}^Y$ -martingale, and  $X(t)$  is a  $\mathbb{P}^Y$ -CTMC with rate  $u_t^Y$ .

**Theorem 1** (Change of Measure). *Assume that both transition rates  $u_t^X$  and  $u_t^Y$  satisfy the conditions in Proposition 1. Suppose that  $u_t^X(x, X(t-)) = 0$  implies  $u_t^Y(x, X(t-)) = 0$ . Choosing the predictable processes*

$$K(t, x) = \log \frac{u_t^Y(x, X(t-))}{u_t^X(x, X(t-))} \text{ and } F(t, x) = (u_t^X(x, X(t-)) - u_t^Y(x, X(t-))) \mathbb{1}(x \neq X(t-))$$

in Equation 6. The R-N derivative can be written as  $\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t = \exp(W(t))$ , where

$$W(t) = \int_0^t \sum_{x \neq X(s-)} [u_s^X(x, X(s-)) - u_s^Y(x, X(s-))] ds + \int_0^t \sum_{x \neq X(s-)} \log \frac{u_s^Y(x, X(s-))}{u_s^X(x, X(s-))} N(ds, x).$$

Then  $\tilde{N}^Y(t, A)$  is a  $\mathbb{P}^Y$ -martingale for any  $A \subseteq \mathcal{S}^{\mathcal{D}}$ . Moreover,  $X(t)$  is a  $\mathbb{P}^Y$ -CTMC with rate  $u_t^Y$ .

Theorem 1 is similar to Theorem F.12 in Pham et al. (2025), and can also be derived by applying Theorem F.12 in Pham et al. (2025) twice, following the argument in their proof of Theorem 2.3. Compared to the Girsanov's theorem used in Zhang et al. (2025), our result can be applied to two CTMCs with arbitrary transition rates. In our result, we provide an explicit expression of the RN derivative in terms of two transition rates. This form is closely related to Theorem 9 in Chen et al. (2023) and Theorem 5.2.12 in Applebaum (2009) for Brownian motion with drift.

### 4.2 BOUND FOR KL DIVERGENCE

Suppose that  $X(t)$  has marginal distribution  $p_t^X$  on the probability space  $(\Omega, \mathcal{F}, \mathbb{P}^X)$  and marginal distribution  $p_t^Y$  on the probability space  $(\Omega, \mathcal{F}, \mathbb{P}^Y)$ . According to Theorem 1, we can derive the bound for KL divergence of  $p_t^X$  and  $p_t^Y$ .

**Theorem 2** (Bound for KL Divergence). *Suppose that  $X(t)$  is a CTMC with natural filtration  $\mathcal{F}_t$  and rate  $u_t^X$  under measure  $\mathbb{P}^X$ ; with rate  $u_t^Y$  under measure  $\mathbb{P}^Y$ . Under the conditions in Theorem 1, the KL divergence between  $\mathbb{P}_t^X$  and  $\mathbb{P}_t^Y$  is*

$$D_{KL}(\mathbb{P}_t^X || \mathbb{P}_t^Y) = \mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s)} D_F(u_s^X(x, X(s)) || u_s^Y(x, X(s))) ds \right\},$$

where  $D_F$  is the Bregman divergence with function  $F(x) = x \log x$ . Consequently, it holds that

$$D_{KL}(p_t^X || p_t^Y) \leq \mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s)} D_F(u_s^X(x, X(s)) || u_s^Y(x, X(s))) ds \right\}.$$

Theorem 2 is *crucial* for developing the estimation error. Due to the sparsity of the transition rate matrix, intuitively, if  $D_F(u_t^X(x, z) || u_t^Y(x, z)) \leq \epsilon$  uniformly for  $t \in [0, 1 - \tau]$  and  $(x, z)$  with  $d^H(x, z) = 1$ , then  $D_{KL}(p_{1-\tau}^X || p_{1-\tau}^Y) \leq \mathcal{D}|\mathcal{S}|\epsilon$ , which is linear in  $\mathcal{D}$ . This matches the results in discrete diffusion models (Chen & Ying, 2024; Zhang et al., 2025; Ren et al., 2025a) and continuous diffusion models (Benton et al., 2024a). In the next section, we will systematically perform an error analysis for our distribution estimation based on discrete flow models.

## 5 MAIN RESULTS

In this section, we establish the error bound for discrete flow-based models. We will first present the additional notation and assumptions. Then we will present the error bounds for transition rate estimation and early stopping, and discuss the implications of the results.

### 5.1 ADDITIONAL NOTATIONS AND ASSUMPTIONS

We first introduce some additional notation. Similar to the training objective (Equation 5), we denote the oracle transition rate as

$$u^0 = \arg \min_u \mathbb{E} \left\{ \sum_{z \neq X(\mathbf{t})} \left[ -u_{\mathbf{t}}(z, X(\mathbf{t}) | X(1)) \log u_{\mathbf{t}}(z, X(\mathbf{t})) + u_{\mathbf{t}}(z, X(\mathbf{t})) \right] \right\} \triangleq \arg \min_u \mathcal{L}(u),$$

where  $Z = (\mathbf{t}, X(\mathbf{t}), X(1))$  is a test point independent of the data  $\mathbb{D}_n$ . Through marginalization trick (Lipman et al., 2024), we have  $u_t^0(z, x) = \mathbb{E}[u_{\mathbf{t}}(z, x | X(1)) | X(\mathbf{t}) = x, \mathbf{t} = t]$ . We also define the best approximation in the function class  $\mathcal{G}_n$  as

$$u^* = \arg \min_{u \in \mathcal{G}_n} \mathbb{E} \left\{ \sum_{z \neq X(\mathbf{t})} \left[ -u_{\mathbf{t}}(z, X(\mathbf{t}) | X(1)) \log u_{\mathbf{t}}(z, X(\mathbf{t})) + u_{\mathbf{t}}(z, X(\mathbf{t})) \right] \right\} \triangleq \arg \min_{u \in \mathcal{G}_n} \mathcal{L}(u).$$

Define

$$g(u, z, Z) \triangleq -u_{\mathbf{t}}(z, X(\mathbf{t}) | X(1)) \log \frac{u_{\mathbf{t}}(z, X(\mathbf{t}))}{u_{\mathbf{t}}^0(z, X(\mathbf{t}))} + u_{\mathbf{t}}(z, X(\mathbf{t})) - u_{\mathbf{t}}^0(z, X(\mathbf{t})).$$

Then, the approximation error of  $u^*$  is

$$\begin{aligned} \mathcal{L}(u^*) - \mathcal{L}(u^0) &= \mathbb{E} \left[ \sum_{z \neq X(\mathbf{t})} g(u^*, z, Z) \right] \\ &= \mathbb{E} \left\{ \sum_{z \neq X(\mathbf{t})} \left[ -u_{\mathbf{t}}(z, X(\mathbf{t}) | X(1)) \log \frac{u_{\mathbf{t}}^*(z, X(\mathbf{t}))}{u_{\mathbf{t}}^0(z, X(\mathbf{t}))} + u_{\mathbf{t}}^*(z, X(\mathbf{t})) - u_{\mathbf{t}}^0(z, X(\mathbf{t})) \right] \right\} \\ &= \mathbb{E} \left[ \sum_{z \neq X(\mathbf{t})} D_F(u_{\mathbf{t}}^0(z, X(\mathbf{t})) || u_{\mathbf{t}}^*(z, X(\mathbf{t}))) \right] \\ &= \inf_{u \in \mathcal{G}_n} \mathbb{E} \left[ \sum_{z \neq X(\mathbf{t})} D_F(u_{\mathbf{t}}^0(z, X(\mathbf{t})) || u_{\mathbf{t}}(z, X(\mathbf{t}))) \right], \end{aligned} \tag{7}$$

where the third equation we use the marginalization trick.

**Assumption 1** (Boundedness). *The conditional transition rate is uniformly bounded from above:  $u_{\mathbf{t}}(z, x | x_1) \leq \overline{M}_c$  uniformly for any  $t \in [0, 1 - \tau]$ ,  $z, x, x_1 \in \mathcal{S}^{\mathcal{D}}$ , where  $d^H(z, x) = 1$ . Moreover, the oracle rate is uniformly bounded from below:  $u_{\mathbf{t}}^0(z, x) > \underline{M}_c$  uniformly for any  $t \in [0, 1 - \tau]$ ,  $z, x \in \mathcal{S}^{\mathcal{D}}$ , where  $d^H(z, x) = 1$ .*

**Assumption 2** (Function Class). *The functions in  $\mathcal{G}_n$  (divided by the oracle transition rate) are uniformly bounded by  $\underline{M}$  and  $\overline{M}$  from below and above, respectively:  $u_t(z, x)/u_t^0(z, x) \in [\underline{M}, \overline{M}]$  for any  $t \in [0, 1 - \tau]$ ,  $(z, x) \in \{\mathcal{S}^{\mathcal{D}} \times \mathcal{S}^{\mathcal{D}} : d^H(z, x) = 1\}$  and  $u \in \mathcal{G}_n$ .*

**Remark 3.** *We only impose conditions on the state pairs with Hamming distance equal to one. Assumption 1 can be satisfied by choosing mixture path and uniform source distribution in Equation 3; see Section B.2 for further discussion. Assumption 2 provides the strong convexity of the Bregman divergence and is crucial to satisfying the condition in Theorem 1. The parameters in Assumption 2 can be fixed if the ranges of  $u$  and  $u^0$  are bounded by two positive constants from above and below; see Section B.3 for a specific example.*

## 5.2 ESTIMATION ERROR DECOMPOSITION

Suppose that  $p_t$  (resp.  $\hat{p}_t$ ) is the marginal distribution of a CTMC with rate  $u_t^0$  (resp.  $\hat{u}_t$ ) at time  $t$ , where  $p_0 = \hat{p}_0$ . For a random sample  $\mathbb{D}_n$ , according to Theorem 2 in the previous section, we have

$$\mathbb{E}_{\mathbb{D}_n}[D_{KL}(p_{1-\tau}||\hat{p}_{1-\tau})] \leq (1-\tau)\mathbb{E}_{\mathbb{D}_n}\mathbb{E}_Z\left\{\sum_{z \neq X(\mathbf{t})} g(\hat{u}, z, Z)\right\} = \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) - \mathcal{L}(u^0)].$$

The following proposition shows that this error bound can be decomposed to the stochastic error and the approximation error.

**Proposition 4** (Estimation Error Decomposition). *For random sample  $\mathbb{D}_n$ , the excess risk of the estimator  $\hat{u}$  through ERM (Equation 5) satisfies*

$$\mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) - \mathcal{L}(u^0)] \leq \underbrace{\mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) + \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u})]}_{\text{Stochastic Error}} + 2 \underbrace{\inf_{u \in \mathcal{G}_n} \mathbb{E}\left[\sum_{z \neq X(\mathbf{t})} D_F(u_t(z, X(\mathbf{t}))||u_t^0(z, X(\mathbf{t})))\right]}_{\text{Approximation Error}}.$$

## 5.3 STOCHASTIC ERROR BOUND

Now, we establish the upper bound of the stochastic error in Proposition 4 using the empirical process theory. Before presenting our stochastic error bound, we first introduce the definition of uniform covering number.

**Definition 3** (Uniform Covering Number, Jiao et al. (2023)). *Let  $S$  be a subset of  $\mathbb{R}^n$ . Given a positive real number  $\epsilon$ , a subset  $\mathcal{C}$  of  $S$  is called an  $\epsilon$ -covering of  $S$  w.r.t. the infinity norm if for any  $x \in S$ , there is  $z \in \mathcal{C}$  such that  $\|z - x\|_{L^\infty} < \epsilon$ . The minimal cardinality  $\mathcal{N}_n(\epsilon, S, L^\infty)$  of all possible  $\mathcal{C}$  is called the covering number of  $S$ . For a given sequence  $\mathbf{s} = \{(t_i, z_i, x_i)\}_{i=1}^n$ , let  $\mathcal{F}|_{\mathbf{s}} = \left\{ \left( u_{t_1}(z_1, x_1), u_{t_2}(z_2, x_2), \dots, u_{t_n}(z_n, x_n) \right) : u \in \mathcal{F} \right\}$ . We define the covering number of  $\mathcal{F}$  constrained on  $\mathbf{s}$  as  $\mathcal{N}_n(\epsilon, \mathcal{F}|_{\mathbf{s}}, L^\infty)$ . The uniform covering number is defined as*

$$\mathcal{N}_n(\epsilon, \mathcal{F}, L^\infty) = \max \{ \mathcal{N}_n(\epsilon, \mathcal{F}|_{\mathbf{s}}, L^\infty) : \mathbf{s} \in ([0, 1 - \tau] \times \mathcal{S}^{\mathcal{D}} \times \mathcal{S}^{\mathcal{D}})^n \}.$$

**Theorem 3** (Stochastic Error). *Assume that Assumption 1 and Assumption 2 hold. Let  $\mathbb{P}_n$  be the empirical measure w.r.t.  $\mathbb{D}_n$ . If  $n \geq CK_1^2|\mathcal{S}|^2\mathcal{D}[\log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty) + \log \mathcal{D}]$ , then the stochastic error in Proposition 4 has the following upper bound:*

$$\mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) + \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u})] \leq \frac{CK_1^2|\mathcal{S}|^2\mathcal{D}[\log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty) + \log \mathcal{D}]}{n} \triangleq \gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n),$$

where  $\mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)$  is the uniform covering number of  $\mathcal{G}_n$ , and  $K_1 = \frac{2\overline{M}_c^2(\overline{M}+1)}{\min(1, \underline{M})\underline{M}_c}$ .

Theorem 3 shows that the stochastic error has the convergence rate

$$O(n^{-1}K_1^2|\mathcal{S}|^2\mathcal{D}[\log \mathcal{N}(1/(2n), \mathcal{G}_n, L^\infty(\mathbb{P}_n)) + \log \mathcal{D}]).$$

This convergence rate is nearly linear in  $\mathcal{D}$  (up to some logarithmic multiplier), which nearly aligns with the stochastic error bound of continuous flow matching (Gao et al., 2024b). We will further discuss this stochastic error bound in Section 5.5. By Pinsker's inequality and Jensen's inequality, we have the following upper bound for the total variation between  $p_{1-\tau}$  and  $\hat{p}_{1-\tau}$ .

$$\begin{aligned} \mathbb{E}_{\mathbb{D}_n}[\text{TV}(p_{1-\tau}, \hat{p}_{1-\tau})] &\leq \sqrt{\frac{1}{2}\mathbb{E}_{\mathbb{D}_n}[D_{KL}(p_{1-\tau}||\hat{p}_{1-\tau})]} \\ &\leq \sqrt{\frac{1}{2}\gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n)} + \sqrt{\inf_{u \in \mathcal{G}_n} \mathbb{E}\left[D_F\left(\sum_{z \neq X(\mathbf{t})} u_t^0(z, X(\mathbf{t})) \middle| \middle| \sum_{z \neq X(\mathbf{t})} u_t(z, X(\mathbf{t}))\right)\right]}. \end{aligned}$$

#### 5.4 EARLY STOPPING ERROR BOUND

We present the early stopping error bound for the total variation between  $p_1$  and  $p_{1-\tau}$  in the following theorem.

**Theorem 4** (Early Stopping Error). *Consider the conditional probability path in Equation 3. Suppose that the source distribution is uniform. Then the total variation between  $p_1$  and  $p_{1-\tau}$  has the following error bound*

$$TV(p_1, p_{1-\tau}) \leq 1 - \exp \left\{ \mathcal{D} \log \left( - (1 - \kappa_{1-\tau}) \frac{|\mathcal{S}| - 1}{|\mathcal{S}|} + 1 \right) \right\} \triangleq \varrho(\mathcal{D}, |\mathcal{S}|, \tau).$$

If the time schedule we used is the linear schedule  $\kappa_t = t$ , then the early stopping error of the discrete flow in Theorem 4 has the same convergence rate as that of the discrete diffusion derived in Theorem 1 of Zhang et al. (2025) as  $\tau \rightarrow 0^+$ .

#### 5.5 OVERALL DISTRIBUTION ERROR

Using triangle inequality, the total variation between  $p_1$  and  $\hat{p}_{1-\tau}$  can be bounded by the early stopping error and the estimation error. Combining the results of Proposition 4, Theorem 3 and Theorem 4, the overall bound for our distribution estimation is

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_n} \left[ TV(p_1, \hat{p}_{1-\tau}) \right] \\ & \leq \mathbb{E}_{\mathcal{D}_n} \left[ TV(p_{1-\tau}, \hat{p}_{1-\tau}) \right] + TV(p_1, p_{1-\tau}) \\ & \leq \underbrace{\sqrt{\frac{1}{2} \gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n)}}_{\text{Stochastic Error}} + \underbrace{\sqrt{\inf_{u \in \mathcal{G}_n} \mathbb{E} \left[ D_F \left( \sum_{z \neq X(\mathbf{t})} u_{\mathbf{t}}^0(z, X(\mathbf{t})) \parallel \sum_{z \neq X(\mathbf{t})} u_{\mathbf{t}}(z, X(\mathbf{t})) \right) \right]}}_{\text{Approximation Error}} + \underbrace{\varrho(\mathcal{D}, |\mathcal{S}|, \tau)}_{\text{Early Stopping Error}}, \end{aligned} \quad (8)$$

where  $\gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n)$  and  $\varrho(\mathcal{D}, |\mathcal{S}|, \tau)$  are defined in Theorem 3 and Theorem 4, respectively.

Given the time schedule  $\kappa_t$ , the sequence length  $\mathcal{D}$  and the vocabulary size  $|\mathcal{S}|$ , there are two quantities related to the above error bound. As  $\tau \rightarrow 0^+$ , the stochastic error increases and the early stopping error decreases; as  $\mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)$  increases, the stochastic error increases and the approximation error decreases. We can choose a suitable early stopping parameter  $\tau$  to balance the stochastic error and the early stopping error. We defer the error analysis for the specific neural network class with the ReLU activation function to Section B.3.

**Balancing the Stochastic Error and Early Stopping Error.** Given  $\overline{M}, \underline{M}$  in Assumption 2 and the vocabulary size  $|\mathcal{S}|$ , we analyze the convergence rate in Equation 8. Consider the linear schedule  $\kappa_t = t$ . Suppose that  $\mathcal{D} \leq \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)$  and the approximation error goes to zero when the size of  $\mathcal{G}_n$  goes to infinity. We mainly focus on the choice of the early stopping parameter  $\tau$  to balance the stochastic error and the early stopping error. Note that if  $\tau \mathcal{D} \rightarrow 0$ , the early stopping error has the convergence rate  $O(\tau \mathcal{D})$ . By taking  $\overline{M}_c, \underline{M}_c$  of Assumption 1 into consideration ( $K_1^2$  has the convergence rate  $\tau^{-4}$  by Section B.2), the stochastic error has the convergence rate  $O(n^{-1} \tau^{-4} \mathcal{D} \log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty))$ . Choosing  $\tau = [n^{-1} \mathcal{D}^{-1} \log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)]^{1/6}$ , if  $n \geq c \mathcal{D}^5 \log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)$ , then the summation of the early stopping error and the squared root of the stochastic error has the convergence rate  $O([n^{-1} \mathcal{D}^5 \log \mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)]^{1/6})$ .

## 6 SIMULATION

To study the empirical performance of discrete flow-based models with different dimensions, sample sizes, and hyperparameters, we conduct several simulation experiments, demonstrating the consistency of our theoretical analysis and empirical results. The data distribution and implementation details are described in Section F.1. Fig. 1, Fig. 2 and Fig. 3 present the simulation results for the empirical version of estimation error (up to an additive constant) with different sample size  $n$ , dimension  $\mathcal{D}$  and early stopping parameter  $\tau$ . We can see that the estimation error is nearly linear in  $\mathcal{D}$ , and the estimation error decreases as  $n$  and  $\tau$  increase. In addition, we also calculate the total variation (of the empirical joint distribution on the first 3 dimensions,  $8^3 = 512$  states in total) with different dimension  $\mathcal{D}$  and  $\tau$  (see Section F.2 for details). As presented in Fig. 4, the total variation decreases first and then increases as  $\tau$  increases, with the minimum achieved around  $\tau = 0.03$  or 0.05.

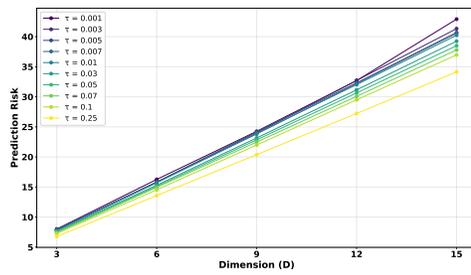


Figure 1: Prediction risk v.s. dimension.

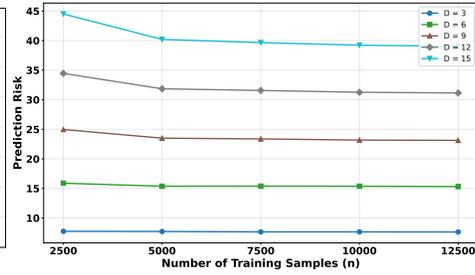


Figure 2: Prediction risk v.s. sample size.

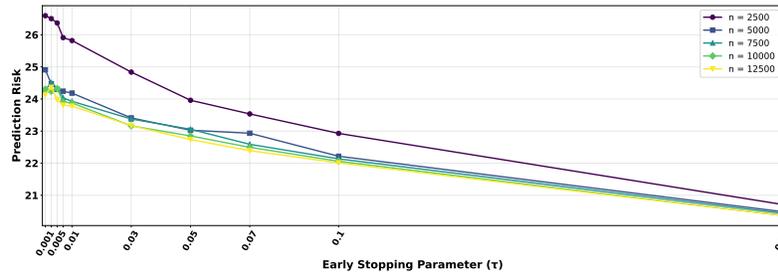


Figure 3: Prediction risk v.s. early stopping parameter.

## 7 CONCLUSION AND FUTURE WORKS

In this article, we have established the first non-asymptotic error bounds for the total variation between the data distribution and the distribution generated by the learned transition rate in discrete flow models. To derive the KL divergence between two CTMCs, we develop a Girsanov-type theorem based on stochastic calculus theory. Building on this result, the estimation error can be decomposed into stochastic error and approximation error. We control the stochastic error by using empirical process theory and balance the stochastic error and the early stopping error by carefully analyzing the boundedness condition. Our theoretical framework serves as an essential theoretical guarantee for studies grounded in CTMCs.

There are two directions deserving future research. Firstly, the error bound we derived has zero discretization error via uniformization technique; while, in practice, the uniformization algorithm does not enjoy sampling efficiency. It would be interesting to investigate the theoretical results of accuracy-efficiency trade-offs for discrete flow models. Secondly, the current work analyzes the approximation error for ReLU networks only; deriving the approximation error for networks with self-attention layers is a challenging problem in approximation theory, requiring greater effort and more careful analysis.

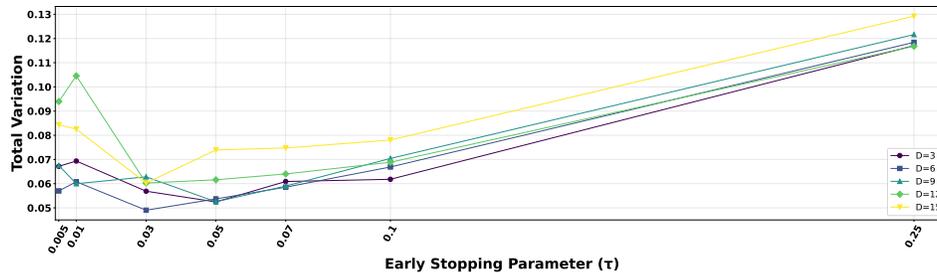


Figure 4: Total variation (of the empirical joint distribution on the first 3 dimensions) v.s. early stopping parameter with uniformization algorithm.

## 540 ETHICS STATEMENT

541

542 This work adheres to the ICLR Code of Ethics. Our research mainly focuses on the theoretical un-  
543 derstanding of the error of the discrete flow model and does not involve human subjects or sensitive  
544 personal data.

545

## 546 REPRODUCIBILITY STATEMENT

547

548 The main theoretical contributions and their assumptions are summarized in Section 4, Section 5  
549 and Section B. The complete proof can be found in Section C, Section D, and Section E.

550

## 551 REFERENCES

552

553 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic inter-  
554 polants. In *International Conference on Learning Representations*, 2023.

555

556 Martin Anthony and Peter L Bartlett. *Neural network learning: Theoretical foundations*. cambridge  
557 university press, 2009.

558 David Applebaum. *Lévy processes and stochastic calculus*. Cambridge University Press, 2009.

559

560 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Struc-  
561 tured denoising diffusion models in discrete state-spaces. In *Advances in Neural Information  
562 Processing Systems*, volume 34, pp. 17981–17993, 2021.

563 Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension  
564 and Pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning  
565 Research*, 20(63):1–17, 2019.

566

567 Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly  $d$ -linear con-  
568 vergence bounds for diffusion models via stochastic localization. In *The Twelfth International  
569 Conference on Learning Representations*, 2024a.

570 Joe Benton, George Deligiannidis, and Arnaud Doucet. Error bounds for flow matching methods.  
571 *Transactions on Machine Learning Research*, 2024b. ISSN 2835-8856.

572

573 Joe Benton, Yuyang Shi, Valentin De Bortoli, George Deligiannidis, and Arnaud Doucet. From  
574 denoising diffusions to denoising Markov models. *Journal of the Royal Statistical Society Series  
575 B: Statistical Methodology*, 86(2):286–301, 2024c.

576 Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and  
577 Arnaud Doucet. A continuous time framework for discrete denoising models. In *Advances in  
578 Neural Information Processing Systems*, volume 35, pp. 28266–28279, 2022.

579 Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative  
580 flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design.  
581 In *International Conference on Machine Learning*, 2024.

582

583 Hongrui Chen and Lexing Ying. Convergence analysis of discrete diffusion model: Exact imple-  
584 mentation through uniformization. *arXiv preprint arXiv:2402.08095*, 2024.

585 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary  
586 differential equations. In *Advances in Neural Information Processing Systems*, volume 31, 2018.

587

588 Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as  
589 learning the score: theory for diffusion models with minimal data assumptions. In *The Eleventh  
590 International Conference on Learning Representations*, 2023.

591 Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.

592

593 Yuan Gao, Jian Huang, and Yuling Jiao. Gaussian interpolation flows. *Journal of Machine Learning  
Research*, 25(253):1–52, 2024a.

- 594 Yuan Gao, Jian Huang, Yuling Jiao, and Shurong Zheng. Convergence of continuous normalizing  
595 flows for learning probability distributions. *arXiv preprint arXiv:2404.00551*, 2024b.  
596
- 597 Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and  
598 Yaron Lipman. Discrete flow matching. In *Advances in Neural Information Processing Systems*,  
599 volume 37, pp. 133345–133385, 2024.
- 600 László Györfi, Michael Kohler, Adam Krzyżak, and Harro Walk. *A distribution-free theory of*  
601 *nonparametric regression*. Springer, 2002.  
602
- 603 Peter Holderrieth, Marton Havasi, Jason Yim, Neta Shaul, Itai Gat, Tommi Jaakkola, Brian Karrer,  
604 Ricky TQ Chen, and Yaron Lipman. Generator matching: Generative modeling with arbitrary  
605 markov processes. In *International Conference on Learning Representations*, 2025.  
606
- 607 Emiel Hooeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows  
608 and multinomial diffusion: Learning categorical distributions. In *Advances in Neural Information*  
609 *Processing Systems*, volume 34, pp. 12454–12465, 2021.
- 610 Yuling Jiao, Guohao Shen, Yuanyuan Lin, and Jian Huang. Deep nonparametric regression on  
611 approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *The Annals of*  
612 *Statistics*, 51(2):691–716, 2023.
- 613 Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow  
614 matching for generative modeling. In *International Conference on Learning Representations*,  
615 2023.  
616
- 617 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ  
618 Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv*  
619 *preprint arXiv:2412.06264*, 2024.  
620
- 621 Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and  
622 transfer data with rectified flow. In *The Eleventh International Conference on Learning Repre-*  
623 *sentations*, 2023.
- 624 Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios  
625 of the data distribution. In *International Conference on Machine Learning*, 2024.  
626
- 627 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin,  
628 Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. In *International Conference*  
629 *on Learning Representations*, 2025.
- 630 James R Norris. *Markov chains*, volume 2. Cambridge University Press, 1998.  
631
- 632 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li.  
633 Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In  
634 *International Conference on Learning Representations*, 2025.
- 635 Le-Tuyet-Nhi Pham, Dario Shariatian, Antonio Ocello, Giovanni Conforti, and Alain Oliviero Dur-  
636 mus. Discrete markov probabilistic models: An improved discrete score-based framework with  
637 sharp convergence bounds under minimal assumptions. In *Forty-second International Conference*  
638 *on Machine Learning*, 2025.  
639
- 640 Yiming Qin, Manuel Madeira, Dorina Thanou, and Pascal Frossard. DeFoG: Discrete flow matching  
641 for graph generation. In *Forty-second International Conference on Machine Learning*, 2025.
- 642 Yinuo Ren, Haoxuan Chen, Grant M. Rotskoff, and Lexing Ying. How discrete and continuous  
643 diffusion meet: Comprehensive analysis of discrete diffusion models via a stochastic integral  
644 framework. In *The Thirteenth International Conference on Learning Representations*, 2025a.  
645
- 646 Yinuo Ren, Haoxuan Chen, Yuchen Zhu, Wei Guo, Yongxin Chen, Grant M. Rotskoff, Molei Tao,  
647 and Lexing Ying. Fast solvers for discrete diffusion models: Theory and applications of high-  
order algorithms. In *Frontiers in Probabilistic Inference: Learning meets Sampling*, 2025b.

- 648 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu,  
649 Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language  
650 models. In *Advances in Neural Information Processing Systems*, volume 37, pp. 130136–130184,  
651 2024.
- 652 Neta Shaul, Itai Gat, Marton Havasi, Daniel Severo, Anuroop Sriram, Peter Holderrieth, Brian Kar-  
653 rer, Yaron Lipman, and Ricky T. Q. Chen. Flow matching with general discrete paths: a kinetic-  
654 optimal perspective. In *International Conference on Learning Representations*, 2025.
- 655 Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and general-  
656 ized masked diffusion for discrete data. In *Advances in Neural Information Processing Systems*,  
657 volume 37, pp. 103131–103167, 2024.
- 658 Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time  
659 discrete diffusion models. In *International Conference on Learning Representations*, 2023.
- 660 AW van der Vaart and Jon A Wellner. *Weak Convergence and Empirical Processes: With Applica-*  
661 *tions to Statistics*. Springer Nature, 2023.
- 662 Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal  
663 Frossard. DiGress: discrete denoising diffusion for graph generation. In *International Conference*  
664 *on Learning Representations*, 2023.
- 665 Zhengchao Wan, Qingsong Wang, Gal Mishne, and Yusu Wang. Elucidating flow matching ODE  
666 dynamics via data geometry and denoisers. In *Forty-second International Conference on Machine*  
667 *Learning*, 2025.
- 668 Jin Wang, Yao Lai, Aoxue Li, Shifeng Zhang, Jiacheng Sun, Ning Kang, Chengyue Wu, Zhenguo Li,  
669 and Ping Luo. FUDOKI: Discrete flow-based unified understanding and generation via kinetic-  
670 optimal velocities. *arXiv preprint arXiv:2505.20147*, 2025.
- 671 Shuntuo Xu, Zhou Yu, and Jian Huang. Estimating unbounded density ratios: Applications in error  
672 control under covariate shift. *arXiv preprint arXiv:2504.01031*, 2025.
- 673 Ling Yang, Ye Tian, Bowen Li, Xinchun Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang.  
674 MMaDA: Multimodal large diffusion language models. *arXiv preprint arXiv:2505.15809*, 2025.
- 675 Zikun Zhang, Zixiang Chen, and Quanquan Gu. Convergence of score-based discrete diffusion  
676 models: A discrete-time analysis. In *The Thirteenth International Conference on Learning Rep-*  
677 *resentations*, 2025.
- 678 Siyan Zhao, Devaansh Gupta, Qinqing Zheng, and Aditya Grover. d1: Scaling reasoning in diffusion  
679 large language models via reinforcement learning. *arXiv preprint arXiv:2504.12216*, 2025.
- 680 Fengqi Zhu, Rongzhen Wang, Shen Nie, Xiaolu Zhang, Chunwei Wu, Jun Hu, Jun Zhou, Jianfei  
681 Chen, Yankai Lin, Ji-Rong Wen, et al. LLaDA 1.5: Variance-reduced preference optimization for  
682 large language diffusion models. *arXiv preprint arXiv:2505.19223*, 2025.
- 683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## APPENDIX

The appendix is organized as follows. Section A gives the uniformization algorithm. Section B offers some discussions on time singularity, assumptions and approximation error analysis. Section C presents some useful lemmas. Section D gives the proof of CTMC theory. Section 5 provides the proof of main results. Section F presents additional simulation experiments and implementation details.

### A UNIFORMIZATION ALGORITHM

---

#### Algorithm 1 Sampling via Uniformization

---

**Require:** A learned transition rate  $\hat{u}$ , an early stopping parameter  $\tau > 0$ , time partition  $0 = t_0 < t_1 < \dots < t_N = 1 - \tau$ , parameters  $\lambda_1, \lambda_2, \dots, \lambda_N$  satisfying  $\sup_{t \in [t_k, t_{k+1}]} \sum_{z \neq x} \hat{u}_t(z, x) \leq \lambda_{k+1}$  for any  $k \in \{0, 1, \dots, N-1\}$ ,  $x \in \mathcal{S}^{\mathcal{D}}$ .

- 1: Draw  $Y_0 \sim \mathcal{U}(\mathcal{S}^{\mathcal{D}})$ .
- 2: **for**  $k = 0$  to  $N - 1$  **do**
- 3:   Draw  $M \sim \text{Poisson}(\lambda_{k+1}(t_{k+1} - t_k))$ .
- 4:   Sample  $M$  points i.i.d. from  $\mathcal{U}([t_k, t_{k+1}])$  and sort them as  $\tau_1 < \tau_2 < \dots < \tau_M$ .
- 5:   Set  $Z_0 = Y_k$ .
- 6:   **for**  $j = 0$  to  $M - 1$  **do**
- 7:     Set  $Z_{j+1} = \begin{cases} z, & \text{with probability } \hat{u}_{\tau_j}(z, Z_j)/\lambda_{k+1} \\ Z_j, & \text{with probability } \sum_{z \neq Z_j} \hat{u}_{\tau_j}(z, Z_j)/\lambda_{k+1} \end{cases}$ , where  $z \neq Z_j$ .
- 8:   **end for**
- 9:   Set  $Y_{k+1} = Z_M$ .
- 10: **end for**
- 11: **return**  $Y_N \sim \hat{p}_{1-\tau}$

---

### B DISCUSSIONS

#### B.1 DISCUSSION ON TIME SINGULARITY

We consider one token case and linear time schedule  $\kappa_t = t$  for simplicity. Note that the conditional transition rate is  $u_t(z, x|x_1) = \frac{\dot{\kappa}_t}{1-\kappa_t}(\delta_{x_1}(z) - \delta_x(z))$ , which can generate the conditional probability path  $p_{t|1}(x|x_1) = \kappa_t \delta_{x_1}(x) + (1-\kappa_t)p_0(x)$ , where  $\kappa_t \rightarrow 1$  as  $t \rightarrow 1$ . Then, when  $p_1$  has full support, the marginal transition rate has a bounded limit as  $t \rightarrow 1$ :

$$\begin{aligned} u_t(z, x) &= \frac{\dot{\kappa}_t}{1-\kappa_t}(p_{1|t}(z|x) - \delta_x(z)) \\ &= \frac{\dot{\kappa}_t}{1-\kappa_t} \left( \frac{p_1(z)p_{t|1}(x|z)}{\sum_z p_1(z)p_{t|1}(x|z)} - \delta_x(z) \right) \\ &= \frac{\dot{\kappa}_t p_0(x)(p_1(z) - \delta_x(z))}{\kappa_t p_1(x) + (1-\kappa_t)p_0(x)} \\ &\rightarrow \frac{\dot{\kappa}_1 p_0(x)(p_1(z) - \delta_x(z))}{p_1(x)}. \end{aligned}$$

When  $p_1$  is not fully supported, for  $x \in \{x \in \mathcal{S}^{\mathcal{D}} : p_1(x) = 0\}$ , the above limit goes to infinity if  $p_1(z) \in (0, 1)$  since  $u_t(z, x) = \frac{\dot{\kappa}_t(p_1(z) - \delta_x(z))}{1-\kappa_t}$  in this case. This result is analogous to the continuous flow matching counterparts (see Proposition C.4 in Wan et al., 2025).

For the transition rate estimator, however, since the conditional transition rate explodes as  $t \rightarrow 1$ , it is hard to control the estimation error of  $\hat{u}$  (see the discussion in Section B.2). Therefore, it is necessary to use the early stopping technique to balance the errors arising from estimation and early stopping.

## B.2 DISCUSSION ON ASSUMPTION 1

In this subsection, we discuss Assumption 1 in some specific scenario, which also provides a reference for time schedule selection.

First, we discuss the upper bound  $\overline{M}_c$  in Assumption 1.

1. Consider the polynomial schedule  $\kappa_t = t^s$  in Equation 3, where  $s \geq 1$ . Then, by mean value theorem, for any  $z^d \neq x^d$ , if  $t \in (0, 1)$ , we have

$$u_t^d(z^d, x^d | x_1^d) \leq \frac{\dot{\kappa}_t}{1 - \kappa_t} = \frac{st^{s-1}}{s(t + \theta(1-t))^{s-1}(1-t)} \leq \frac{1}{1-t},$$

which means that the conditional transition rate has the uniform upper bound  $\overline{M}_c = \frac{1}{\tau}$ .

2. Consider the cosine schedule  $\kappa_t = \cos^2(\frac{\pi}{2}(1-t))$  in Equation 3, which is kinetic optimal as mentioned in Shaul et al. (2025). Then, by  $\tan(a) > a$  ( $a \in (0, \pi/2)$ ), for any  $z^d \neq x^d$ , if  $t \in (0, 1)$ , we have

$$u_t^d(z^d, x^d | x_1^d) \leq \frac{\dot{\kappa}_t}{1 - \kappa_t} = \frac{\pi \cos(\frac{\pi}{2}(1-t)) \sin(\frac{\pi}{2}(1-t))}{\sin^2(\frac{\pi}{2}(1-t))} = \frac{\pi}{\tan(\frac{\pi}{2}(1-t))} \leq \frac{2}{1-t},$$

which means that the conditional transition rate has the uniform upper bound  $\overline{M}_c = \frac{2}{\tau}$ .

Given other quantities, if  $\tau \rightarrow 0$ , then the estimation error will go to infinity as presented in Theorem 3, since  $K_1$  will go to infinity as  $\overline{M}_c \rightarrow +\infty$ . By the proof of Theorem 3, under the above two scenarios, the factor  $K_1^2$  in the error bound of Theorem 3 grows at the order of  $O(\tau^{-4})$ .

Next, we discuss the lower bound  $\underline{M}_c$  for the oracle transition rate  $u_t^0(z, x)$ , where  $d^H(z, x) =$

1. We define the mixing coefficient  $\alpha \in (0, 1)$  such that  $\alpha \leq \frac{p_1(x_1^d | x_1^d)}{p_1(x_1^d)} \leq \alpha^{-1}$  for any  $x_1 \in \mathcal{S}^D$  and  $d \in [D]$  (for example, if  $p_1(x_1) = \prod_{d=1}^D p_1(x_1^d)$ , then  $\alpha = 1$ ). We denote  $\beta = \sup_{d, z^d, x^d} (p_1^d(z^d) / p_1^d(x^d))$ , where  $p_1^d$  is the marginal distribution of  $p_1$  (for example, if  $p_1^d$  is of uniform distribution, then  $\beta = 1$ ). Suppose that the source distribution is uniform; that is,  $p_{t|1}^d(x^d | x_1^d) = \frac{1-\kappa_t}{|\mathcal{S}|} + \kappa_t \delta_{x_1^d}(x^d)$ . Then, for any  $d \in [D]$  and  $(z, x) \in \mathcal{S}^D \times \mathcal{S}^D$  such that  $z^d \neq x^d, z^{\setminus d} = x^{\setminus d}$ , we have

$$\begin{aligned} u_t^0(z, x) &= \sum_{x_1} u_t^d(z^d, x^d | x_1^d) p_{1|t}(x_1 | x) \\ &= \frac{\sum_{x_1} u_t^d(z^d, x^d | x_1^d) p_{t|1}^d(x^d | x_1^d) \prod_{i \neq d} p_{t|1}^i(x^i | x_1^i) p_1(x_1)}{\sum_{x_1} p_{t|1}^d(x^d | x_1^d) \prod_{i \neq d} p_{t|1}^i(x^i | x_1^i) p_1(x_1)} \\ &= \frac{\frac{\dot{\kappa}_t}{|\mathcal{S}|} p_1^d(z^d) \sum_{x_1^{\setminus d}} \prod_{i \neq d} p_{t|1}^i(x^i | x_1^i) p_1(x_1^{\setminus d} | X(1)^d = z^d)}{\frac{1-\kappa_t}{|\mathcal{S}|} \sum_{x_1^{\setminus d}} \prod_{i \neq d} p_{t|1}^i(x^i | x_1^i) p_1(x_1^{\setminus d}) + \kappa_t p_1^d(x^d) \sum_{x_1^{\setminus d}} \prod_{i \neq d} p_{t|1}^i(x^i | x_1^i) p_1(x_1^{\setminus d} | X(1)^d = x^d)} \\ &\geq \frac{|\mathcal{S}|^{-1} \dot{\kappa}_t \alpha}{(1 - \kappa_t) \beta + \kappa_t \beta / \alpha} \\ &\geq |\mathcal{S}|^{-1} \dot{\kappa}_t \alpha^2 / \beta. \end{aligned}$$

To give a specific example, we consider an AR( $k$ ) autoregressive structure for the sequence; that is,  $p_1(x_1^d | x_1^{<d}) = p_1(x_1^d | x_1^{((l-k)\vee 1):(d-1)})$  for  $d > 2$ . Then we have

$$\begin{aligned} &\frac{p_1(x_1^{\setminus d} | x_1^d)}{p_1(x_1^{\setminus d})} \\ &= \frac{p_1(x_1^1, \dots, x_1^{d-1}) \prod_{l=d}^D p_1(x_1^l | x_1^{((l-k)\vee 1):(l-1)})}{p_1(x_1^1, \dots, x_1^{d-1}) p_1(x_1^d) \prod_{j=1}^k p(x_1^{d+j} | x_1^{((d-k)\vee 1):(d+j-1)}) \prod_{l=d+k+1}^D p_1(x_1^l | x_1^{((l-k)\vee 1):(l-1)})} \\ &= \frac{p_1(x_1^d | x_1^{((d-k)\vee 1):(d-1)}) \prod_{j=1}^k p_1(x_1^{d+j} | x_1^{((d+j-k)\vee 1):(d+j-1)})}{p_1(x_1^d) \prod_{j=1}^k p(x_1^{d+j} | x_1^{((d-k)\vee 1):(d+j-1)})}. \end{aligned}$$

If  $\log p_1(x_1^d)$  and  $\log p_1(x_1^d | x_1^{((d-2k)\vee 1):(d-1)}) \in [-c, 0]$  for each  $d \in [\mathcal{D}]$ ,  $x_1 \in \mathcal{S}^{\mathcal{D}}$ , then the right-hand side of the above equation satisfying  $\alpha < \frac{p_1(x_1^d | x_1^d)}{p_1(x_1^d)} < \alpha^{-1}$  with  $\alpha = \exp(-(k+1)c)$  and  $\beta = \exp(-c)$ , which are independent of  $\mathcal{D}$ .

Therefore, if  $\alpha$  and  $\beta$  are positive constants and the vocabulary size  $|\mathcal{S}|$  is fixed, then  $\underline{M}_c = (|\mathcal{S}|^{-1} \alpha^2 / \beta) \inf_{t \in [0, 1-\tau]} \kappa_t$  is the uniform lower bound for the oracle transition rate  $u_t^0(z, x)$ , where  $d^H(z, x) = 1$ . Empirically, one may adopt the linear schedule  $\kappa_t = t$ , which is commonly used in practice. More generally, a composite schedule can be employed:  $\kappa_t = \frac{t + \kappa_t^0}{2}$ , where  $\kappa_t^0$  is a schedule that is a non-decreasing function of  $t$ .

### B.3 ERROR ANALYSIS WITH RELU NETWORKS

In this subsection, we write  $u_t(z, x)$  to  $u(t, z, x)$  for any transition rate  $u$ . We consider using neural network (NN) functions with ReLU activation function to approximate the oracle transition rate  $u^0$ . To impose some regularity condition on the oracle rate, we first introduce the Hölder class. For a finite constant  $B_0 > 0$  and input dimension  $d \in \mathbb{N}^+$ , the Hölder class of functions  $\mathcal{H}^\beta([0, 1]^d, B_0)$  is defined as

$$\mathcal{H}^\beta([0, 1]^d, B_0) = \left\{ f : [0, 1]^d \rightarrow \mathbb{R}, \max_{\|\alpha\|_1 \leq s} \|\partial^\alpha f\|_\infty \leq B_0, \max_{\|\alpha\|_1 = s} \sup_{x \neq y} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{\|x - y\|_2^r} \leq B_0 \right\},$$

where  $\beta = r + s$ ,  $s = \lfloor \beta \rfloor \in \mathbb{N}$ ,  $\partial^\alpha = \partial^{\alpha_1} \dots \partial^{\alpha_d}$  with  $\alpha = (\alpha_1, \dots, \alpha_d)^\top \in \mathbb{N}^d$  and  $\|\alpha\|_1 = \sum_{i=1}^d \alpha_i$ .

**Assumption 3** (Hölder smoothness). *There exists a function  $\bar{u}^0$  belongs to the Hölder class  $\mathcal{H}^\beta([0, 1] \times [0, |\mathcal{S}|]^{2\mathcal{D}}, B_0)$  for a given  $\beta > 0$  and a finite constant  $B_0 > \underline{M}_c$  such that  $u^0 = \bar{u}^0$  for any  $t \in [0, 1]$  and  $x \in \mathcal{S}^{\mathcal{D}} \times \mathcal{S}^{\mathcal{D}}$  with  $d^H(x_{1:\mathcal{D}}, x_{(\mathcal{D}+1):(2\mathcal{D})}) = 1$ .*

Without loss of generality, we rewrite  $\bar{u}^0$  as  $u^0$  through this analysis. This assumption requires that the oracle transition rate  $u^0$  does not have terminal time singularity at time  $t = 1$ ; e.g.,  $p_1$  has full support as mentioned in Section B.1.

#### B.3.1 BOUNDING APPROXIMATION ERROR

Following the proof of Lemma B.6 in Xu et al. (2025), we prove that there exists a sequence of neural networks with the ReLU activation function that can control the approximation error. We focus on the region  $(t, x) \in [0, 1] \times [0, |\mathcal{S}|]^{2\mathcal{D}}$ . Let  $\tilde{u}^0(t, x') = u^0(t, |\mathcal{S}|x')$  for  $x' \in [0, 1]^{2\mathcal{D}}$ . Then  $\tilde{u}^0 \in \mathcal{H}^\beta([0, 1]^{2\mathcal{D}+1}, |\mathcal{S}|^\beta B_0)$ . Lemma 5 implies that for any  $S_1, S_2 \in \mathbb{N}^+$ , there exists a function  $f^*$  implemented by a ReLU network with depth  $D = 21(\lfloor \beta \rfloor + 1)^2 S_1 \lceil \log_2(8S_1) \rceil$  and width  $W = 38(\lfloor \beta \rfloor + 1)^2 (2\mathcal{D} + 1)^{\lfloor \beta \rfloor + 1} S_2 \lceil \log_2(8S_2) \rceil$  such that

$$|f^*(t, x') - \tilde{u}^0(t, x')| \leq 18|\mathcal{S}|^\beta B_0 (\lfloor \beta \rfloor + 1)^2 (2\mathcal{D} + 1)^{\lfloor \beta \rfloor + (\beta \vee 1)/2} (S_1 S_2)^{-2\beta/(2\mathcal{D}+1)},$$

for any  $(t, x') \in [0, 1]^{2\mathcal{D}+1} \setminus \Omega([0, 1]^{2\mathcal{D}+1}, K, \delta)$ , where  $K = \lceil (S_1 S_2)^{2/(2\mathcal{D}+1)} \rceil$  and  $\delta$  is an arbitrary number in  $(0, 1/(3K))$ . Let  $f^{**}(t, x) = f^*(t, x/|\mathcal{S}|)$  for  $x \in [0, |\mathcal{S}|]^{2\mathcal{D}}$  be a network with depth  $D + 1$ . Note that a clip function can be expressed as a two-layer ReLU network, then we define the NN  $f^{***}(t, x) = \underline{M}_c + \text{ReLU}(B_0 - \text{ReLU}(B_0 - f^{**}(t, x)) - \underline{M}_c)$  with depth  $D + 3$ , whose range is  $[\underline{M}_c, B_0]$ .

Consequently, we have the following bound:

$$|f^{***}(t, x) - u^0(t, x)| \leq 18|\mathcal{S}|^\beta B_0 (\lfloor \beta \rfloor + 1)^2 (2\mathcal{D} + 1)^{\lfloor \beta \rfloor + (\beta \vee 1)/2} (S_1 S_2)^{-2\beta/(2\mathcal{D}+1)},$$

where  $(t, x) \in [0, 1] \times [|\mathcal{S}|]^{2\mathcal{D}}$  such that  $(t, x/|\mathcal{S}|) \in [0, 1] \times [|\mathcal{S}|]^{2\mathcal{D}} \setminus \Omega([0, 1]^{2\mathcal{D}+1}, K, \delta)$  and

$$\Omega([0, 1]^{2\mathcal{D}+1}, K, \delta) = \cup_{i=1}^{2\mathcal{D}+1} \{x = [x_1, x_2, \dots, x_{2\mathcal{D}+1}]^\top : x_i \in \cup_{k=1}^{K-1} (k/K - \delta, k/K)\}.$$

Note that  $\mathbf{t} \sim \mathcal{U}([0, 1 - \tau])$ , whose probability measure is absolutely continuous w.r.t. the Lebesgue measure. Moreover, when  $\delta$  is sufficiently small, the set  $\{1/|\mathcal{S}|, 2/|\mathcal{S}|, \dots, 1\}$  and the set  $\cup_{k=1}^{K-1} (k/K - \delta, k/K)$  are disjoint. Then, the event

$$\{(\mathbf{t}, z/|\mathcal{S}|, x/|\mathcal{S}|) \in [0, 1]^{2\mathcal{D}+1} \setminus \Omega([0, 1]^{2\mathcal{D}+1}, K, \delta) \text{ for any } (z, x) \in \mathcal{S}^{\mathcal{D}} \times \mathcal{S}^{\mathcal{D}}\}$$

holds with probability greater than  $1 - 2K\delta$ . Thus, by the smoothness of the Bregman divergence, following the proof of Theorem 4.2 in Jiao et al. (2023), since  $\delta$  is an arbitrary number in  $(0, 1/(3K))$ , the approximation error has the following bound

$$\begin{aligned} & \inf_{u \in \mathcal{G}_n} \mathbb{E} \left[ D_F \left( \sum_{z \neq X(\mathbf{t})} u^0(\mathbf{t}, z, X(\mathbf{t})) \parallel \sum_{z \neq X(\mathbf{t})} u(\mathbf{t}, z, X(\mathbf{t})) \right) \right] \\ & \leq \frac{1}{2\underline{M}_c} \mathbb{E} \left| \sum_{z \neq X(\mathbf{t})} f^{***}(\mathbf{t}, z, X(\mathbf{t})) - \sum_{z \neq X(\mathbf{t})} u^0(\mathbf{t}, z, X(\mathbf{t})) \right|^2 \\ & \leq C \underline{M}_c^{-1} |\mathcal{S}|^{2\lfloor \beta \rfloor} B_0^2 (\lfloor \beta \rfloor + 1)^4 (2\mathcal{D} + 1)^{2\lfloor \beta \rfloor + \beta \vee 1} (S_1 S_2)^{-4\beta/(2\mathcal{D}+1)}, \end{aligned} \quad (9)$$

where  $\mathcal{G}_n$  is the class of ReLU networks with range  $[\underline{M}_c, B_0]$ , depth  $\mathcal{D}^* = 21(\lfloor \beta \rfloor + 1)^2 S_1 \lceil \log_2(8S_1) \rceil + 3$ , width  $W^* = 38(\lfloor \beta \rfloor + 1)^2 (2\mathcal{D} + 1)^{\lfloor \beta \rfloor + 1} S_2 \lceil \log_2(8S_2) \rceil$ ,  $S_1, S_2 \in \mathbb{N}^+$ . If  $\underline{M}_c$  is fixed, then the parameters in Assumption 2 are fixed constants.

### B.3.2 CHOOSING HYPERPARAMETERS

Before the end of this section, the convergence rate might omit a logarithmic multiplier of  $n\tau^4$ . Similar to Section 5.5, given the vocabulary size  $|\mathcal{S}|$ , by Lemma 6 and Lemma 7, if  $n \geq C\tau^{-4} |\mathcal{S}|^2 \mathcal{D} S^* \mathcal{D}^* \log S^* \log n$ , the stochastic error is

$$\gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n) = O\left(\frac{|\mathcal{S}|^2 \mathcal{D} S^* \mathcal{D}^* \log S^* \log n}{n\tau^4}\right),$$

where  $S^*$  is the number of parameters of the ReLU networks in  $\mathcal{G}_n$ . Note that for a ReLU network with depth  $\mathcal{D}^*$  and width  $W^*$  and input dimension  $2\mathcal{D} + 1$ , we have (assume that  $\mathcal{D} \lesssim W^* \mathcal{D}^*$ )

$$S \leq \underbrace{W^*(2\mathcal{D} + 1) + W^*}_{\text{input layer}} + \underbrace{((W^*)^2 + W^*)(\mathcal{D}^* - 1)}_{\text{hidden layer}} + \underbrace{W^* + 1}_{\text{output layer}} = O((W^*)^2 \mathcal{D}^*).$$

Therefore, by choosing  $S_1 = O((n\tau^4)^{\frac{2\mathcal{D}+1}{(4\mathcal{D}+4\beta+2)}})$  and  $S_2 = O(1)$ , we have

$$\begin{aligned} W^* &= O\left((2\mathcal{D} + 1)^{\lfloor \beta \rfloor + 1}\right); \mathcal{D}^* = O\left((n\tau^4)^{\frac{2\mathcal{D}+1}{(4\mathcal{D}+4\beta+2)}} \log(n\tau^4)\right); \\ S^* &= O\left(\mathcal{D}^{2\lfloor \beta \rfloor + 2} n^{\frac{2\mathcal{D}+1}{(4\mathcal{D}+4\beta+2)}} \log(n\tau^4)\right), \end{aligned}$$

yielding that

$$\gamma_n(\mathcal{D}, |\mathcal{S}|, \mathcal{G}_n) = O\left(|\mathcal{S}|^2 \mathcal{D}^{2\lfloor \beta \rfloor + 2} (\log(n\tau^4))^4 (n\tau^4)^{-\frac{2\beta}{(2\mathcal{D}+2\beta+1)}}\right).$$

By Equation 9, the approximation error is

$$\begin{aligned} & \inf_{u \in \mathcal{G}_n} \mathbb{E} \left[ D_F \left( \sum_{z \neq X(\mathbf{t})} u^0(\mathbf{t}, z, X(\mathbf{t})) \parallel \sum_{z \neq X(\mathbf{t})} u(\mathbf{t}, z, X(\mathbf{t})) \right) \right] \\ & = O\left(|\mathcal{S}|^{2\beta} (2\mathcal{D} + 1)^{2\lfloor \beta \rfloor + \beta \vee 1} (n\tau^4)^{-\frac{2\beta}{(2\mathcal{D}+2\beta+1)}}\right). \end{aligned}$$

If the vocabulary size  $|\mathcal{S}|$  is fixed, then the summation of the approximation error and the stochastic error has the convergence rate (up to some logarithmic multiplier)

$$\mathbb{E}_{\mathbb{D}_n} [D_{KL}(p_{1-\tau} \parallel \hat{p}_{1-\tau})] = O\left(\mathcal{D}^{2\lfloor \beta \rfloor + 2\vee \beta} (n\tau^4)^{-\frac{2\beta}{(2\mathcal{D}+2\beta+1)}}\right). \quad (10)$$

### B.3.3 DISTRIBUTION CONSISTENCY

Combining Equation 8, Theorem 4 and Equation 10, as  $\tau\mathcal{D} \rightarrow 0$ , we have

$$\mathbb{E}_{\mathbb{D}_n} [\text{TV}(p_1, \hat{p}_{1-\tau})] = O\left(\mathcal{D}^{\lfloor\beta\rfloor+1\nu\beta}(n\tau^4)^{-\frac{\beta}{(2\mathcal{D}+2\beta+1)}} + \tau\mathcal{D}\right).$$

When  $\mathcal{D}^{\lfloor\beta\rfloor+1\nu\beta}(n\tau^4)^{-\frac{\beta}{(2\mathcal{D}+2\beta+1)}} + \tau\mathcal{D}$  goes to zero (e.g., the sequence dimension  $\mathcal{D}$  is fixed and the early stopping parameter  $\tau = o(n^{-1/4})$ ), we can find that the distribution estimation is consistent in the sense of  $\mathbb{E}_{\mathbb{D}_n} [\text{TV}(p_1, \hat{p}_{1-\tau})] \rightarrow 0$ .

## C SOME USEFUL LEMMAS

**Lemma 1** (Exit Time). *Under the assumptions in Proposition 1, considering the stopping time  $T_t = \min\{h > 0 : \Delta X(t+h) \neq 0\}$ , we have*

$$\mathbb{P}(T_t > h | \mathcal{F}_t) = \exp\left(\int_t^{t+h} u_s(X(t), X(t)) ds\right).$$

*Proof.* Consider the event  $E_n = \{X(t) = X(t + \frac{h}{2^n}) = \dots = X(t+h)\}$ . Since  $E_{n+1} \subseteq E_n$ , by dominated convergence theorem and Markov property, we have

$$\begin{aligned} & \mathbb{P}(T_t > h | \mathcal{F}_t) \\ &= \mathbb{P}(\cap_n E_n | \mathcal{F}_t) = \lim_{n \rightarrow \infty} \mathbb{P}(E_n | \mathcal{F}_t) \\ &= \lim_{n \rightarrow \infty} \prod_{i=1}^{2^n} \left( \mathbb{P}\left(X\left(t + \frac{i}{2^n}h\right) = X(t) \middle| X\left(t + \frac{i-1}{2^n}h\right) = X(t), \mathcal{F}_t\right) \right) \\ &= \lim_{n \rightarrow \infty} \exp\left(\sum_{i=1}^{2^n} \frac{h}{2^n} \frac{\log \mathbb{P}\left(X\left(t + \frac{i}{2^n}h\right) = X(t) \middle| X\left(t + \frac{i-1}{2^n}h\right) = X(t), \mathcal{F}_t\right) - \log 1}{h/(2^n)}\right) \\ &= \lim_{n \rightarrow \infty} \exp\left(\sum_{i=1}^{2^n} \frac{h}{2^n} \frac{u_{t+\frac{i-1}{2^n}h}(X(t), X(t)) \frac{h}{2^n} + R_i}{h/(2^n)}\right) \\ &= \exp\left(\int_t^{t+h} u_s(X(t), X(t)) ds\right), \end{aligned}$$

where the remainder  $R_i \leq C(h/2^n)^2$  for some constant  $C$  not depending on  $i$  by Proposition 1.  $\square$

**Lemma 2** (Itô's formula). *Consider the random measure  $N(t, A)$  associated with a CTMC  $X(t)$ , which is defined in Definition 2. Define a stochastic integral  $W(t) = W(0) + \int_0^t \int_A K(s, x) N(ds, dx)$ , where  $K$  is a predictable process and  $A \subseteq \mathcal{S}^{\mathcal{D}}$ . For each  $f \in C(\mathbb{R})$ ,  $t \geq 0$ , we have*

$$f(W(t)) - f(W(0)) = \int_0^t \int_A [f(W(s-) + K(s, x)) - f(W(s-))] N(ds, dx).$$

*Proof.* This proof is similar to Lemma 4.4.5 in Applebaum (2009). Let  $T_0^A = 0$  and  $T_n^A = \inf\{t > T_{n-1}^A : \Delta N(t, A) \neq 0\}$ . Note that

$$\begin{aligned} f(W(t)) - f(W(0)) &= \sum_{n=1}^{\infty} f(W(t \wedge T_n^A)) - f(W(t \wedge T_n^A -)) \\ &= \sum_{n=1}^{\infty} [f(W(t \wedge T_n^A -) + K(t \wedge T_n^A, X(t \wedge T_n^A))) - f(W(t \wedge T_n^A -))] \\ &= \int_0^t \int_A [f(W(s-) + K(s, x)) - f(W(s-))] N(ds, dx), \end{aligned}$$

which completes the proof.  $\square$

**Lemma 3** (Exponential martingale). *Suppose that  $N(t, A)$  is the random measure associated with the CTMC  $X(t)$ . Consider the following stochastic integral*

$$W(t) = W(0) + \int_0^t \sum_{x \neq X(s-)} F(s, x) ds + \int_0^t \int_{x \in \mathcal{S}^D} K(s, x) N(ds, dx).$$

Then  $\exp(W(t))$  is an exponential martingale if for each  $x \in \mathcal{S}^D$ ,

$$F(t, x) = -(e^{K(t, x)} - 1)u_t(x, X(t-))\mathbb{1}(x \neq X(t-)).$$

*Proof.* By Itô's product formula (e.g., Theorem 4.4.13 in Applebaum, 2009) and Itô's formula (Lemma 2), we have

$$\begin{aligned} d[\exp(W(t))] &= \exp(W(t)) \sum_{x \neq X(t-)} F(t, x) dt + \exp(W(t-)) \int_{x \in \mathcal{S}^D} (e^{K(t, x)} - 1) N(dt, dx) \\ &= \exp(W(t-)) \sum_{x \neq X(t-)} \left\{ F(t, x) dt + (e^{K(t, x)} - 1) N(dt, x) \right\} \\ &= \exp(W(t-)) \sum_{x \neq X(t-)} \left\{ (e^{K(t, x)} - 1) \tilde{N}(dt, x) \right\}, \end{aligned}$$

which completes the proof by Proposition 3.  $\square$

**Lemma 4.** *For any random variable  $Z \in \mathcal{F}_{t+h}$ , we have*

$$\mathbb{E}^Y(Z|\mathcal{F}_t) = \frac{\mathbb{E}^X \left[ \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_{t+h} Z | \mathcal{F}_t \right]}{\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t}.$$

*Proof.* For any set  $A \in \mathcal{F}_t$ , we have

$$\begin{aligned} \mathbb{E}^Y(Z\mathbb{1}_A) &= \mathbb{E}^X \left[ \mathbb{1}_A \mathbb{E}^X \left( \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_{t+h} Z | \mathcal{F}_t \right) \right] \\ &= \mathbb{E}^Y \left[ \frac{\mathbb{1}_A \mathbb{E}^X \left( \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_{t+h} Z | \mathcal{F}_t \right)}{\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t} \right], \end{aligned}$$

which completes the proof by the definition of the R-N derivative.  $\square$

**Lemma 5** (Theorem 3.3 in Jiao et al. (2023)). *Assume that  $f \in \mathcal{H}^\beta([0, 1]^d, B_0)$  with  $\beta = s + r$ ,  $s \in \mathbb{N}$  and  $r \in (0, 1]$ . For any  $S_1, S_2 \in \mathbb{N}^+$ , there exists a function  $\phi_0$  implemented by a ReLU network with width  $W = 38(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1} S_2 \lceil \log_2(8S_2) \rceil$  and depth  $D = 21(\lfloor \beta \rfloor + 1)^2 S_1 \lceil \log_2(8S_1) \rceil$  such that*

$$|f(x) - \phi_0(x)| \leq 18B_0(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + (\beta \vee 1)/2} (S_1 S_2)^{-2\beta/d},$$

for all  $x \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$ , where

$$\Omega([0, 1]^d, K, \delta) = \cup_{i=1}^d \{x = [x_1, x_2, \dots, x_d]^\top : x_i \in \cup_{k=1}^{K-1} (k/K - \delta, k/K)\},$$

with  $K = \lfloor (S_1 S_2)^{2/d} \rfloor$  and  $\delta$  an arbitrary number in  $(0, 1/(3K))$

**Lemma 6** (Theorem 12.2 in Anthony & Bartlett (2009)). *Let  $\mathcal{G}$  be a set of real functions that map a domain  $\mathcal{X}$  to a bounded interval  $[0, B]$ . The pseudo-dimension  $Pdim(\mathcal{G})$  of  $\mathcal{G}$  is defined as the largest integer  $m$  for which there exists  $(x_1, \dots, x_m, y_1, \dots, y_m) \in \mathcal{X} \times \mathbb{R}^m$  such that for any  $(b_1, \dots, b_m) \in \{0, 1\}^m$  there exists  $f \in \mathcal{G}$  such that  $f(x_i) > y_i$  if and only if  $b_i = 1$  for any  $i \in [m]$ . Then, for  $n \geq Pdim(\mathcal{G})$  and  $B \geq \epsilon$ , we have*

$$\mathcal{N}_n(\epsilon, \mathcal{G}, L^\infty) \leq \left( \frac{eBn}{\epsilon Pdim(\mathcal{G})} \right)^{Pdim(\mathcal{G})}.$$

**Lemma 7** (Theorem 7 in Bartlett et al. (2019)). *Let  $\mathcal{G}$  be a ReLU neural network function class with depth  $L$  and number of parameters  $S$ . Then, there exists a universal constant  $C$  such that*

$$Pdim(\mathcal{G}) \leq CSL \log S.$$

## D PROOF OF RESULTS FOR CTMC

### D.1 PROOF OF PROPOSITION 1

*Proof.* By construction,  $X(t)$  has Markov property since the Poisson process  $N(t)$  has independent increments. By conditioning argument (e.g., Theorem 3.7.9 in Durrett, 2019) and the definition of Poisson processes, for  $z \neq x$ , we have

$$\begin{aligned} \mathbb{P}(X(t+h) = z | X(t) = x) &= Mh \exp(-Mh) \mathbb{P}(X(t+h) = z | X(t) = x, N(t+h) - N(t) = 1) + O(h^2) \\ &= Mh \exp(-Mh) \times \left( \int_0^h \underbrace{\frac{1}{h}}_{\text{jump location: uniform distribution}} \frac{u_{t+\theta}(z, x)}{M} d\theta \right) + O(h^2) \\ &= u_t(z, x)h + R_t, \end{aligned}$$

where the remainder  $R_t$  satisfying

$$\begin{aligned} R_t &\leq Lh^2 \exp(-Mh) + \mathbb{P}(N(t+h) - N(t) > 1, X(t+h) = z | X(t) = x) \\ &\leq \exp(-Mh)[Lh^2 + \exp(Mh) - 1 - Mh] \leq (M^2 + L)h^2 = O(h^2). \end{aligned}$$

Consequently, since the state space is finite, we have

$$\mathbb{P}(X(t+h) = x | X(t) = x) = 1 - \sum_{z \neq x} \mathbb{P}(X(t+h) = z | X(t) = x) = 1 + u_t(x, x)h + O(h^2),$$

which completes the proof.  $\square$

### D.2 PROOF OF PROPOSITION 2

*Proof.* By Definition 1, we have

$$p_{t+h}(x) - p_t(x) = \sum_{z \in \mathcal{S}^D} (p_{t+h|t}(x|z) - \delta_z(x))p_t(z) = \sum_{z \in \mathcal{S}^D} (u_t(x, z)h + o(h))p_t(z).$$

Then,

$$\dot{p}_t(x) = \lim_{h \rightarrow 0^+} \frac{p_{t+h}(x) - p_t(x)}{h} = \sum_{z \in \mathcal{S}^D} u_t(x, z)p_t(z),$$

which completes the proof.  $\square$

### D.3 PROOF OF PROPOSITION 3

*Proof.* In this proof, we consider the event  $\{X(t) = x\}$  if we are given  $\mathcal{F}_t$ . Following the proof of Proposition 1, conditioning on  $\mathcal{F}_t$ , since the Poisson process  $N(t)$  has independent increments, we have

$$\begin{aligned} &\mathbb{P}(N(t+h, A) - N(t, A) = 1 | \mathcal{F}_t) \\ &\leq \mathbb{P}(N(t+h) - N(t) = 1) \sum_{z \in A} \mathbb{P}(X(t+h) = z, X(t) \neq z | N(t+h) - N(t) = 1, \mathcal{F}_t) \\ &\quad + \mathbb{P}(N(t+h) - N(t) > 1) \\ &= \sum_{z \in A} u_t(z, X(t)) \mathbb{1}(z \neq X(t))h + R_t, \end{aligned}$$

where  $R_t$  is bounded by  $|A|(M^2 + L)h^2$ . Since the rate of  $X(t)$  is bounded above by  $M$ , by uniformization, we have

$$\begin{aligned} &\mathbb{E} \left[ \left\{ N(t+h, A) - N(t, A) \right\} \mathbb{1}(N(t+h, A) - N(t, A) \geq 2) \middle| \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \left\{ N(t+h) - N(t) \right\} \mathbb{1}(N(t+h) - N(t) \geq 2) \middle| \mathcal{F}_t \right] \\ &= Mh - (Mh \exp(-Mh)) \leq M^2 h^2. \end{aligned}$$

Then,

$$\mathbb{E}\left[N(t+h, A) - N(t, A) \middle| \mathcal{F}_t\right] = \sum_{z \in A} u_t(z, X(t)) \mathbb{1}(z \neq X(t))h + R.$$

Here, the remainder  $R$  is bounded by  $M^2 h^2 + |A|(M^2 + L)h^2$ . Then, since  $X(t)$  has only finite jumps in  $[s, t]$  almost surely, by dominated convergence theorem, we have

$$\begin{aligned} & \mathbb{E}\left[N(t, A) - N(s, A) \middle| \mathcal{F}_s\right] \\ &= \sum_{i=1}^n \mathbb{E}\left[\mathbb{E}\left[N\left(s + \frac{i}{n}(t-s), A\right) - N\left(s + \frac{i-1}{n}(t-s), A\right) \middle| \mathcal{F}_{s+\frac{i-1}{n}(t-s)}\right] \middle| \mathcal{F}_s\right] \\ &= \sum_{i=1}^n \mathbb{E}\left[\sum_{z \in A} u_{s+\frac{i-1}{n}(t-s)}(z, X\left(s + \frac{i-1}{n}(t-s)\right)) \mathbb{1}(z \neq X\left(s + \frac{i-1}{n}(t-s)\right)) \frac{1}{n}(t-s) \middle| \mathcal{F}_s\right] \\ &\quad + O\left(\frac{1}{n}\right) \\ &\rightarrow \mathbb{E}\left[\int_s^t \sum_{z \in A} u_r(z, X(r-)) \mathbb{1}(z \neq X(r-)) dr \middle| \mathcal{F}_s\right], \end{aligned}$$

as  $n \rightarrow \infty$ . Thus, the process

$$\tilde{N}(t, A) = N(t, A) - \underbrace{\int_0^t \sum_{z \in A} u_s(z, X(s-)) \mathbb{1}(z \neq X(s-)) ds}_{\text{predictable process}}$$

is a (martingale-valued) compensated random measure of the random measure  $N(t, A)$ . We are done.  $\square$

#### D.4 PROOF OF THEOREM 1

*Proof.* By computing directly, we can obtain that

$$\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t = \exp \left\{ \int_0^t \sum_{x \neq X(s-)} [u_s^X(x, X(s-)) - u_s^Y(x, X(s-))] ds + \int_0^t \sum_{x \neq X(s-)} \log \frac{u_s^Y(x, X(s-))}{u_s^X(x, X(s-))} N(ds, x) \right\}.$$

Next, we divide the proof into two parts: first we prove that  $\tilde{N}^Y(t, A)$  is a  $\mathbb{P}^Y$ -martingale, and then we conclude that  $X(t)$  is a  $\mathbb{P}^Y$ -CTMC with rate  $u_t^Y$ .

$\tilde{N}^Y(t, A)$  is a  $\mathbb{P}^Y$ -martingale. By Itô's product formula (see Theorem 4.4.13 in Applebaum, 2009) and Lemma 3, we have

$$\begin{aligned} d[\tilde{N}^Y(t, A) \exp(W(t))] &= \tilde{N}^Y(t-, A) \exp(W(t-)) \sum_{x \neq X(t-)} \left\{ (e^{K(t,x)} - 1) \tilde{N}^X(dt, x) \right\} \\ &\quad + \exp(W(t-)) \tilde{N}^Y(dt, A) \\ &\quad + \underbrace{\exp(W(t-)) \sum_{x \neq X(t-): x \in A} \left\{ (e^{K(t,x)} - 1) N(dt, x) \right\}}_{\text{quadratic variation}} \\ &= \tilde{N}^Y(t-, A) \exp(W(t-)) \sum_{x \neq X(t-)} \left\{ (e^{K(t,x)} - 1) \tilde{N}^X(dt, x) \right\} \\ &\quad + \exp(W(t-)) \tilde{N}^X(dt, A) \\ &\quad + \exp(W(t-)) \sum_{x \neq X(t-): x \in A} \left\{ (e^{K(t,x)} - 1) N(dt, x) \right\} \end{aligned}$$

$$\begin{aligned}
& - \exp(W(t-)) \sum_{x \neq X(t-): x \in A} \left[ \frac{u_t^Y(x, X(t-))}{u_t^X(x, X(t-))} - 1 \right] u_t^X(x, X(t-)) dt \\
& = \tilde{N}^Y(t-, A) \exp(W(t-)) \sum_{x \neq X(t-)} \left\{ \frac{u_t^Y(x, X(t-))}{u_t^X(x, X(t-))} - 1 \right\} \tilde{N}^X(dt, x) \\
& \quad + \exp(W(t-)) \sum_{x \neq X(t-): x \in A} \frac{u_t^Y(x, X(t-))}{u_t^X(x, X(t-))} \tilde{N}^X(dt, x),
\end{aligned}$$

which implies that  $\tilde{N}^Y(t, A) \exp(W(t))$  is a  $\mathbb{P}^X$ -martingale. By Lemma 5.2.11 in Applebaum (2009),  $\tilde{N}^Y(t, A)$  is a  $\mathbb{P}^Y$ -martingale.

$X(t)$  is a  $\mathbb{P}^Y$ -CTMC with rate  $u_t^Y$ . Considering  $Z = \mathbb{1}(X(t+h) = z)$  in Lemma 4, we can obtain that

$$\mathbb{P}^Y(X(t+h) = z | \mathcal{F}_t) = \mathbb{E}^X[\exp(W(t+h) - W(t)) \mathbb{1}(X(t+h) = z) | \mathcal{F}_t].$$

Since the right-hand side in above equation only depends on  $X(t)$ , thus  $X(t)$  has Markov property under  $\mathbb{P}^Y$ . It suffices to show that  $\mathbb{P}^Y(X(t+h) = z | \mathcal{F}_t) = u_t(z, X(t))h + o(h)$  for any  $z \neq X(t)$ .

Since  $\exp(W(t+h) - W(t)) \leq \exp(Ch)$  for some constant  $C$ , by Lemma 4, we have

$$\begin{aligned}
\mathbb{P}^Y(N(t+h, A) - N(t, A) > 1 | \mathcal{F}_t) &= \frac{\mathbb{E}^X \left[ \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_{t+h} \mathbb{1}(N(t+h, A) - N(t, A) > 1) | \mathcal{F}_t \right]}{\frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \Big|_t} \\
&\leq \exp(Ch) \mathbb{P}^X(N(t+h, A) - N(t, A) > 1 | \mathcal{F}_t) \\
&= O(h^2).
\end{aligned}$$

Note that, for  $\mathbb{P}^Y$ -submartingale  $(\tilde{N}^Y)^2$ , by Itô product formula, we have Doob-Meyer decomposition:

$$\begin{aligned}
d(\tilde{N}^Y)^2(t, A) &= 2\tilde{N}^Y(t-) d\tilde{N}^Y(t, A) + d[\tilde{N}^Y, \tilde{N}^Y](t, A) \\
&= 2\tilde{N}^Y(t-, A) d\tilde{N}^Y(t, A) + dN(t, A).
\end{aligned}$$

Then, by Lemma 4 again,  $\mathbb{E}^Y[(\tilde{N}^Y(t+h, A) - \tilde{N}^Y(t, A))^2 | \mathcal{F}_t] = \mathbb{E}^Y[N(t+h, A) - N(t, A) | \mathcal{F}_t] = O(h)$ . By Cauchy-Schwarz inequality, we can obtain

$$\begin{aligned}
& \mathbb{E}^Y[(N(t+h, A) - N(t, A)) \mathbb{1}(N(t+h, A) - N(t, A) \geq 2) | \mathcal{F}_t] \\
& \leq \sqrt{\mathbb{E}^Y[(N(t+h, A) - N(t, A))^2 | \mathcal{F}_t]} \mathbb{P}^Y(N(t+h, A) - N(t, A) > 1 | \mathcal{F}_t) \\
& \leq \sqrt{\left( 2\mathbb{E}^Y[(\tilde{N}^Y(t+h, A) - \tilde{N}^Y(t, A))^2 | \mathcal{F}_t] + 2\mathbb{E}^Y \left[ \left( \int_t^{t+h} \sum_{z \in A} u_s^Y(z, X(s-)) \mathbb{1}(z \neq X(s-)) ds \right)^2 \Big| \mathcal{F}_t \right] \right)} \\
& \quad \times \sqrt{\mathbb{P}^Y(N(t+h, A) - N(t, A) > 1 | \mathcal{F}_t)} \\
& = O(h^{3/2})
\end{aligned}$$

Then, conditioning on  $\mathcal{F}_t$ , for  $z \neq X(t)$ , we have

$$\begin{aligned}
\mathbb{P}^Y(X(t+h) = z | \mathcal{F}_t) &\leq \mathbb{P}^Y(N(t+h, z) - N(t, z) = 1 | \mathcal{F}_t) + O(h^2) \\
&= \mathbb{E}^Y(N(t+h, z) - N(t, z) | \mathcal{F}_t) + O(h^{3/2}) \\
&= \int_t^{t+h} u_t^Y(z, X(t)) ds + \int_t^{t+h} u_s^Y(z, X(t)) - u_t^Y(z, X(t)) ds \\
& \quad + \mathbb{E}^Y \left( \int_t^{t+h} u_s^Y(z, X(s-)) - u_s^Y(z, X(t)) ds \Big| \mathcal{F}_t \right) + O(h^{3/2}) \\
&= u_t(z, X(t))h + O(h^{3/2}),
\end{aligned}$$

1188 since

$$\begin{aligned}
1189 \mathbb{E}^Y \left( \int_t^{t+h} u_s^Y(z, X(s-)) - u_s^Y(z, X(t)) ds \middle| \mathcal{F}_t \right) &\leq 2Mh \mathbb{P}^Y(N(t+h, \mathcal{S}^{\mathcal{D}}) - N(t, \mathcal{S}^{\mathcal{D}}) \geq 1 | \mathcal{F}_t) \\
1190 &\leq 2Mh \mathbb{E}^Y(N(t+h, \mathcal{S}^{\mathcal{D}}) - N(t, \mathcal{S}^{\mathcal{D}}) | \mathcal{F}_t) \\
1191 &\leq 2M^2 h^2 = O(h^2).
\end{aligned}$$

1195 This completes the proof. □

## 1199 D.5 PROOF OF THEOREM 2

1200 *Proof.* We compute directly:

$$\begin{aligned}
1202 &D_{\text{KL}}(\mathbb{P}_t^X || \mathbb{P}_t^Y) \\
1203 &= \mathbb{E}^X \left[ -\log \left( \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \middle|_t \right) \right] \\
1204 &= -\mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s-)} [u_s^X(x, X(s-)) - u_s^Y(x, X(s-))] + u_s^X(x, X(s-)) \log \frac{u_s^Y(x, X(s-))}{u_s^X(x, X(s-))} ds \right\} \\
1205 &\quad - \mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s-)} \log \frac{u_s^Y(x, X(s-))}{u_s^X(x, X(s-))} \tilde{N}^X(ds, x) \right\} \\
1206 &= \mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s)} D_F(u_s^X(x, X(s)) || u_s^Y(x, X(s))) ds \right\}, \\
1207 & \\
1208 & \\
1209 & \\
1210 & \\
1211 & \\
1212 & \\
1213 & \\
1214 &
\end{aligned}$$

1215 where the last equation holds since the integrator is the Lebesgue measure.

1216 By Jensen's inequality, we have

$$\begin{aligned}
1217 & \\
1218 &D_{\text{KL}}(p_t^X || p_t^Y) = \mathbb{E}^X \left[ -\log \left( \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \middle|_{\sigma(X_t)} \right) \right] \\
1219 &\leq \mathbb{E}^X \left[ -\log \left( \frac{d\mathbb{P}^Y}{d\mathbb{P}^X} \middle|_{\mathcal{F}_t} \right) \right] \\
1220 &= \mathbb{E}^X \left\{ \int_0^t \sum_{x \neq X(s)} D_F(u_s^X(x, X(s)) || u_s^Y(x, X(s))) ds \right\}. \\
1221 & \\
1222 & \\
1223 & \\
1224 & \\
1225 &
\end{aligned}$$

1226 This completes the proof. □

## 1228 E PROOF OF MAIN RESULTS

### 1230 E.1 PROOF OF PROPOSITION 4

1231 *Proof.* This proof is similar to the proof of Lemma 3.1 in Jiao et al. (2023). Note that

$$1233 \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}_n(\hat{u}) - \mathcal{L}_n(u^0)] \leq \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}_n(u^*) - \mathcal{L}_n(u^0)],$$

1234 which yields that

$$1235 -\mathcal{L}(u^0) \leq 2\mathcal{L}(u^*) - \mathcal{L}(u^0) - \mathbb{E}_{\mathbb{D}_n}[2\mathcal{L}_n(\hat{u})].$$

1236 Then we have

$$1237 \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) - \mathcal{L}(u^0)] \leq \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) + \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u})] + 2[\mathcal{L}(u^*) - \mathcal{L}(u^0)],$$

1238 which completes the proof by using Equation 7. □

## E.2 PROOF OF THEOREM 3

*Sketch of proof.* We first decompose the objective into  $\mathcal{D}$  terms

$$\mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) + \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u})] = \mathbb{E}_{\mathbb{D}_n} \left[ \mathbb{E}_Z \left[ \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) \right],$$

where we use  $\sum_{d:z \neq X(\mathbf{t})} = \sum_{z^d \neq X(\mathbf{t})^d; z \setminus d = X(\mathbf{t}) \setminus d}$  for notational simplicity. To bound the expectation, we focus on the tail probability

$$\mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) - \frac{2}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) > \frac{t}{\mathcal{D}} \right).$$

Next, we replace the expectation by an empirical mean of a "ghost" sample  $\mathbb{D}'_n$  independent of  $\mathbb{D}_n$ . Additionally, to use concentration inequality, we have to consider the nonnegative empirical mean  $\frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2$  instead of  $\frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i)$ . By symmetrization argument and introducing Rademacher variables  $\{\epsilon_i\}_{i=1}^n$ , the above probability can be bounded by

$$\begin{aligned} & \mathbb{P} \left\{ \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}], \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 > \left[ \frac{ct}{\mathcal{D}} + \frac{c}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right\} \\ & + \mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) > \frac{ct}{\mathcal{D}} + \frac{c}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right). \end{aligned}$$

Finally, conditioning on  $\mathbb{D}_n$  and introducing a covering, we can apply some concentration inequalities to bound these two terms.

*Proof.* This proof is similar to the proof of Theorem 11.4 in Györfi et al. (2002), where they only consider the squared loss. Note that

$$\begin{aligned} & \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) + \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u})] \\ & = \mathbb{E}_{\mathbb{D}_n}[\mathcal{L}(\hat{u}) - \mathcal{L}(u^0) - 2\mathcal{L}_n(\hat{u}) + 2\mathcal{L}(u^0)] \\ & = \mathbb{E}_{\mathbb{D}_n} \left[ \mathbb{E}_Z \left[ \sum_{z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) \right] \\ & = \mathbb{E}_{\mathbb{D}_n} \left[ \mathbb{E}_Z \left[ \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) \right], \end{aligned}$$

Notice that

$$\begin{aligned} & \mathbb{E}_{\mathbb{D}} \left[ \mathbb{E}_Z \left[ \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) \right] \\ & \leq a_n + \int_{a_n}^{\infty} \mathbb{P} \left( \mathbb{E}_Z \left[ \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) > t \right) dt, \end{aligned}$$

where  $a_n$  is a quantity depending on  $n$ , which we will choose later. Thus, we can focus on the following probability first.

$$\mathbb{P} \left( \mathbb{E}_Z \left[ \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(\hat{u}, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\hat{u}, z, Z_i) > t \right)$$

Following the proof of symmetrization in probability (e.g., Lemma 2.3.7 in van der Vaart & Wellner (2023)), consider the following event

$$\mathcal{B}_1 = \{\hat{\mathcal{A}}(t) \text{ is a non-empty set}\},$$

where  $t \leq 1$  and

$$\begin{aligned} \hat{A}(t) &= \left\{ (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right] - \frac{2}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) > \frac{t}{\mathcal{D}} \right\} \\ &= \left\{ (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right] - \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right. \\ &\quad \left. > \frac{1}{3} \left( \frac{2t}{\mathcal{D}} + \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) + \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right) \right\}. \end{aligned} \quad (11)$$

Let  $(\bar{u}, \bar{d})$  be a (random) function such that  $(\bar{u}, \bar{d}) \in \hat{A}(t)$  if  $\mathcal{B}_1$  holds, and let  $(\bar{u}, \bar{d}) = (1, 1)$  if  $\mathcal{B}_1$  does not hold, where  $(\bar{u}, \bar{d})$  depends on  $\mathbb{D}_n$ . By assumptions, note that  $|g(u, z, Z)| \leq \frac{2\bar{M}_c^2(\bar{M}+1)}{\min(1, \bar{M})\bar{M}_c} \triangleq K_1$  for any  $u \in \mathcal{G}_n$  and  $z \in \mathcal{S}^{\mathcal{D}}$  such that  $d^H(z, X(\mathbf{t})) = 1$ . Let  $\mathbb{D}'_n = \{Z'_i\}_{i=1}^n$  be an independent copy of  $\mathbb{D}_n$ . By Markov's inequality, we have

$$\begin{aligned} &\mathbb{P} \left\{ \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) | \mathbb{D}_n \right] - \frac{1}{n} \sum_{i=1}^n \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(\bar{u}, z, Z'_i) > \frac{1}{4} \left[ \frac{t}{\mathcal{D}} + \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) | \mathbb{D}_n \right] \right] | \mathbb{D}_n \right\} \\ &\leq \frac{16\mathbb{E}[(\sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z))^2 | \mathbb{D}_n]}{n \left( \frac{t}{\mathcal{D}} + \mathbb{E}[\sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) | \mathbb{D}_n] \right)^2} \leq \frac{8K_1\mathcal{D}|\mathcal{S}|}{nt} \triangleq \eta_1. \end{aligned} \quad (12)$$

Since  $\mathcal{B}_1 \in \sigma(\mathbb{D}_n)$ , we have

$$\begin{aligned} &(1 - \eta_1) \mathbb{P} \left( \exists u \in \mathcal{G}_n : \mathbb{E} \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) - \frac{2}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) > t \right) \\ &\leq (1 - \eta_1) \mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) - \frac{2}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) > \frac{t}{\mathcal{D}} \right) \\ &\leq \mathbb{E} \left( \mathbb{1}(\mathcal{B}_1) \mathbb{P} \left\{ \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) | \mathbb{D}_n \right] - \frac{1}{n} \sum_{i=1}^n \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(\bar{u}, z, Z'_i) \right. \right. \\ &\quad \left. \left. \leq \frac{1}{4} \left[ \frac{t}{\mathcal{D}} + \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) | \mathbb{D}_n \right] \right] | \mathbb{D}_n \right\} \right) \\ &\leq \mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) - \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right. \\ &\quad \left. > \frac{1}{4} \left[ \frac{t}{\mathcal{D}} + \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right] \right), \end{aligned} \quad (13)$$

where we use equation 12 and the definition of  $\hat{A}(t)$  in the second inequality.

Since  $\sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i)$  might be negative, we want to focus on the nonnegative quantities  $\{\sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i)\}^2$ . By a union bound, we can obtain that

$$\begin{aligned} &\mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) - \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right. \\ &\quad \left. > \frac{1}{4} \left[ \frac{t}{\mathcal{D}} + \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right] \right) \\ &\leq 2\mathbb{P} \left( \exists (u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right)^2 - \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \right. \\ &\quad \left. > \frac{1}{2} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right)^2 + \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \right] \right) \end{aligned}$$

$$\begin{aligned}
& + \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) - \frac{1}{n} \sum_{i=1}^n \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \\
& > \frac{1}{4} \left[ \frac{t}{\mathcal{D}} + \mathbb{E} \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right], \\
& \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right)^2 - \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \\
& \leq \frac{1}{2} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right)^2 + \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \right], \\
& \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right)^2 - \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \\
& \leq \frac{1}{2} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right)^2 + \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \right] \\
& \triangleq 2\mathbb{P}(\mathcal{B}_2) + \mathbb{P}(\mathcal{B}_3)
\end{aligned}$$

**Bounding  $\mathbb{P}(\mathcal{B}_2)$ .** Let  $(\bar{u}, \bar{d}) \in \mathcal{G}_n \times [\mathcal{D}]$  be a function such that

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \left( \sum_{\bar{d}:z \neq X_i(\mathbf{t}_i)} g(\bar{u}, z, Z_i) \right)^2 - \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right]^2 \\
& > \frac{1}{2} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left( \sum_{\bar{d}:z \neq X_i(\mathbf{t}_i)} g(\bar{u}, z, Z_i) \right)^2 + \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right]^2 \right],
\end{aligned}$$

if  $\mathcal{B}_2$  holds; and  $(\bar{u}, \bar{d}) = (1, 1)$  otherwise, which depends on  $\mathbb{D}_n$ . Conditioning on  $\mathcal{B}_2$  and  $\mathbb{D}_n$ , by Markov's inequality, we have

$$\begin{aligned}
& \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \left( \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(\bar{u}, z, Z'_i) \right)^2 - \mathbb{E} \left[ \left( \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \right. \\
& > \frac{1}{8} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left( \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(\bar{u}, z, Z'_i) \right)^2 + \mathbb{E} \left[ \sum_{\bar{d}:z \neq X(\mathbf{t})} g^2(\bar{u}, z, Z) \middle| \mathbb{D}_n \right] \right] \middle| \mathbb{D}_n \Big\} \\
& = \mathbb{P} \left\{ \frac{7}{n} \sum_{i=1}^n \left( \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(\bar{u}, z, Z'_i) \right)^2 - 7 \mathbb{E} \left[ \left( \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \right. \\
& > \frac{t}{\mathcal{D}} + 2 \mathbb{E} \left[ \left( \sum_{z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \middle| \mathbb{D}_n \Big\} \\
& \leq \frac{49 \mathbb{E} \left[ \left( \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^4 \middle| \mathbb{D}_n \right]}{n \left( \frac{t}{\mathcal{D}} + 2 \mathbb{E} \left[ \left( \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \right)^2} \leq \frac{49 K_1^2 |\mathcal{S}|^2 \mathcal{D}}{4nt} \triangleq \eta_2.
\end{aligned} \tag{14}$$

Then, by using equation 14 and a symmetrization argument, we have

$$\begin{aligned}
& (1 - \eta_2) \mathbb{P}(\mathcal{B}_2) \\
& \leq \mathbb{E} \left( \mathbb{1}(\mathcal{B}_2) \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 - \mathbb{E} \left[ \left( \sum_{d:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \right. \right. \\
& \leq \frac{1}{8} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{\bar{d}:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 + \mathbb{E} \left[ \left( \sum_{\bar{d}:z \neq X(\mathbf{t})} g(\bar{u}, z, Z) \right)^2 \middle| \mathbb{D}_n \right] \right] \middle| \mathbb{D}_n \Big\} \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [D], \frac{4}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 - \frac{7}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 \right. \\
&\quad \left. > 3 \left[ \frac{t}{D} + \mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 \right] \right) \\
&= \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [D], \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 - \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right. \\
&\quad \left. > \frac{3}{11} \left[ \frac{2t}{D} + 2\mathbb{E} \left[ \sum_{d:z \neq X(\mathbf{t})} g(u, z, Z) \right]^2 + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 \right. \right. \\
&\quad \left. \left. + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right) \\
&\leq \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [D], \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 - \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right. \\
&\quad \left. > \frac{3}{11} \left[ \frac{2t}{D} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right) \\
&= \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [D], \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 - \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right. \\
&\quad \left. > \frac{3}{11} \left[ \frac{2t}{D} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right) \\
&\leq 2\mathbb{P}\left\{\exists(u, d) \in \mathcal{G}_n \times [D], \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right. \\
&\quad \left. > \frac{3}{11} \left[ \frac{t}{D} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right\},
\end{aligned}$$

where  $\{\epsilon_i\}_{i=1}^n$  is a sequence of i.i.d. Rademacher random variables, independent of  $\mathbb{D}_n$ .

Given  $\mathbb{D}_n$  and  $d \in [D]$ , we consider the following sequence with length  $n|\mathcal{S}|$ :

$$\mathbf{s}_{|\mathcal{S}|n} = \left\{ (\mathbf{t}_i, X_i(\mathbf{t}_i)|_{X_i(\mathbf{t}_i)^{d=s}}, X_i(\mathbf{t}_i)) \right\}_{i \in [n], s \in \mathcal{S}},$$

where we denote  $X_i(\mathbf{t}_i)|_{X_i(\mathbf{t}_i)^{d=s}} = (X_i(\mathbf{t}_i)^1, \dots, X_i(\mathbf{t}_i)^{d-1}, s, X_i(\mathbf{t}_i)^{d+1}, \dots, X_i(\mathbf{t}_i)^D)$ . Let  $\mathcal{H}_\delta^d(\mathbb{D}_n)$  be a  $L^\infty$   $\delta$ -covering of  $\mathcal{G}_n|_{\mathbf{s}_{|\mathcal{S}|n}}$  with minimal size. That is to say, given  $\mathbb{D}_n$  and  $d \in [D]$ , for any  $u \in \mathcal{G}_n$ , there exists  $\tilde{u} \in \mathcal{H}_\delta^d(\mathbb{D}_n)$  such that for any  $i \in [n]$  and  $z \in \mathcal{S}^D$  satisfying  $z^d \neq X_i(\mathbf{t}_i)^d$  and  $z^{\setminus d} = X_i(\mathbf{t}_i)^{\setminus d}$

$$|u_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i)) - \tilde{u}_{\mathbf{t}_i}(z, X_i(\mathbf{t}_i))| \leq \delta.$$

Note that, for  $d^H(z, X(\mathbf{t})) = 1$ ,

$$|g(u_1, z, Z) - g(u_2, z, Z)| \leq \left( \frac{\overline{M}_c}{\underline{M}_c} + 1 \right) |u_1(z, Z) - u_2(z, Z)| \stackrel{\Delta}{=} K_2 |u_1(z, Z) - u_2(z, Z)|. \quad (15)$$

Then, given  $\mathbb{D}_n$  and  $d \in [D]$ , for any  $u \in \mathcal{G}_n$ , there exists  $\tilde{u} \in \mathcal{H}_\delta^d(\mathbb{D}_n)$  such that for any  $i \in [n]$ ,

$$\left| \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right)^2 - \left( \sum_{d:z \neq X_i(\mathbf{t}_i)} g(\tilde{u}, z, Z_i) \right)^2 \right| \leq 2K_1 K_2 |\mathcal{S}|^2 \delta \stackrel{\Delta}{=} \eta_3 \delta.$$

Thus, taking  $\delta_1 = \frac{t}{7\eta_3\mathcal{D}}$ , by union bound, we have

$$\begin{aligned}
\mathbb{P}(\mathcal{B}_2) &\leq \frac{2}{1-\eta_2} \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [D], \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right. \\
&> \left. \frac{3}{11} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right] \right) \\
&\leq \frac{2}{1-\eta_2} \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{H}_{\delta_1}^d(\mathbb{D}_n), \eta_3 \delta_1 + \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right. \\
&> \left. \frac{3}{11} \left[ \frac{t}{\mathcal{D}} + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 - \eta_3 \delta_1 \right] \right) \\
&= \frac{2}{1-\eta_2} \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{H}_{\delta_1}^d(\mathbb{D}_n), \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right. \\
&> \left. \frac{3t}{11\mathcal{D}} + \frac{3}{11n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 - \frac{14}{11} \eta_3 \delta_1 \right) \\
&= \frac{2}{1-\eta_2} \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{H}_{\delta_1}^d(\mathbb{D}_n), \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right. \\
&> \left. \frac{t}{11\mathcal{D}} + \frac{3}{11n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right).
\end{aligned}$$

Conditioning on  $\mathbb{D}_n$ , by Hoeffding's inequality, we have

$$\begin{aligned}
&\mathbb{E} \left\{ \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{H}_{\delta_1}^d(\mathbb{D}_n), \frac{1}{n} \sum_{i=1}^n \epsilon_i \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right. \right. \\
&> \left. \left. \frac{t}{11\mathcal{D}} + \frac{3}{11n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \middle| \mathbb{D}_n \right) \right\} \\
&\leq \mathbb{E} \left\{ 2 \sum_{d=1}^{\mathcal{D}} \mathcal{N}_{|\mathcal{S}|n}(\delta_1, \mathcal{H}_{\delta_1}^d(\mathbb{D}_n), L^\infty) \right. \\
&\quad \left. \times \max_{u \in \mathcal{H}_{\delta_1}^d(\mathbb{D}_n)} \exp \left( - \frac{cn \left( \frac{t}{11\mathcal{D}} + \frac{3}{11n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \right)^2}{\frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z \neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^4} \right) \right\} \\
&\leq 2\mathcal{D} \mathcal{N}_{|\mathcal{S}|n} \left( \frac{t}{7\eta_3\mathcal{D}}, \mathcal{G}_n, L^\infty \right) \exp \left( - \frac{cnt}{K_1^2 |\mathcal{S}|^2 \mathcal{D}} \right).
\end{aligned}$$

Consequently, we have

$$\mathbb{P}(\mathcal{B}_2) \leq \left( \frac{4\mathcal{D}}{1-\eta_2} \right) \mathcal{N}_{|\mathcal{S}|n} \left( \frac{t}{7\eta_3\mathcal{D}}, \mathcal{G}_n, L^\infty \right) \exp \left( - \frac{cnt}{K_1^2 |\mathcal{S}|^2 \mathcal{D}} \right). \quad (16)$$

**Bounding  $\mathbb{P}(\mathcal{B}_3)$ .** For  $F(x) = x \log x$ , if  $a, b \in [c, C]$ , then we have the strong convexity  $D_F(a||b) \geq \frac{1}{2C}(a-b)^2$ . Thus, for any  $(u, d) \in \mathcal{G}_n \times [D]$ , by equation 15 and QM-AM inequality,

we have

$$\begin{aligned}
\mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} g(u, z, Z)\right] &= \mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} D_F(u_{\mathbf{t}}^0(z, X(\mathbf{t}))|u_{\mathbf{t}}(z, X(\mathbf{t})))\right] \\
&\geq \frac{C}{MM_c}\mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} (u_{\mathbf{t}}^0(z, X(\mathbf{t})) - u_{\mathbf{t}}(z, X(\mathbf{t})))^2\right] \\
&\geq \frac{C}{MM_c K_2^2}\mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} g^2(u, z, Z)\right] \\
&\geq \frac{C}{MM_c K_2^2 |\mathcal{S}|}\mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} g(u, z, Z)\right]^2 \\
&\geq C\alpha\mathbb{E}\left[\sum_{d:z\neq X(\mathbf{t})} g(u, z, Z)\right]^2,
\end{aligned} \tag{17}$$

where  $\alpha = (\overline{MM_c K_2^2} |\mathcal{S}|)^{-1} \wedge (3C^{-1}) \wedge (1/2)$ . Then, by the definition of  $\mathcal{B}_3$ , equation 17 and introducing random signs, we have

$$\begin{aligned}
&\mathbb{P}(\mathcal{B}_3) \\
&\leq \mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \sum_{d:z\neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) - \frac{1}{n} \sum_{i=1}^n \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \\
&\quad \left. > \left(\frac{1}{2} - \frac{C\alpha}{12}\right) \frac{t}{\mathcal{D}} + \frac{C\alpha}{24} \left(\frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X'_i(\mathbf{t}'_i)} g(u, z, Z'_i) \right\}^2 + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right)\right) \\
&\leq 2\mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \\
&\quad \left. > \left(\frac{1}{4} - \frac{C\alpha}{24}\right) \frac{t}{\mathcal{D}} + \frac{C\alpha}{24n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right),
\end{aligned}$$

where  $\{\epsilon_i\}_{i=1}^n$  is a sequence of i.i.d. Rademacher random variables, independent of  $\mathbb{D}_n$ . Note that, given  $\mathbb{D}_n$  and  $d \in [\mathcal{D}]$ , for any  $u \in \mathcal{G}_n$ , there exists  $\tilde{u} \in \mathcal{H}_\delta^d(\mathbb{D}_n)$  such that for any  $i \in [n]$ ,

$$\left| \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) - \sum_{d:z\neq X_i(\mathbf{t}_i)} g(\tilde{u}, z, Z_i) \right| \leq K_2 |\mathcal{S}| \delta \triangleq \eta_4 \delta.$$

Thus, if  $\delta_2 = \frac{t}{(16\eta_4 + 2\eta_3)\mathcal{D}}$ , conditioning on  $\mathbb{D}_n$ , by union bound and Bernstein's inequality (e.g., Lemma A.2 of Györfi et al., 2002), we have

$$\begin{aligned}
&2\mathbb{P}\left(\exists(u, d) \in \mathcal{G}_n \times [\mathcal{D}] : \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \\
&\quad \left. > \left(\frac{1}{4} - \frac{C\alpha}{24}\right) \frac{t}{\mathcal{D}} + \frac{C\alpha}{24n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2\right) \\
&= 2\mathbb{E}\left\{ \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{G}_n : \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \right. \\
&\quad \left. \left. > \left(\frac{1}{4} - \frac{C\alpha}{24}\right) \frac{t}{\mathcal{D}} + \frac{C\alpha}{24n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 \middle| \mathbb{D}_n \right)\right\} \\
&\leq 2\mathbb{E}\left\{ \sum_{d=1}^{\mathcal{D}} \mathbb{P}\left(\exists u \in \mathcal{H}_\delta^d(\mathbb{D}_n) : \eta_4 \delta_2 + \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i)\right. \right. \\
&\quad \left. \left. > \left(\frac{1}{4} - \frac{C\alpha}{24}\right) \frac{t}{\mathcal{D}} + \frac{C\alpha}{24n} \sum_{i=1}^n \left\{ \sum_{d:z\neq X_i(\mathbf{t}_i)} g(u, z, Z_i) \right\}^2 - \frac{C\alpha\eta_3\delta_2}{24} \middle| \mathbb{D}_n \right)\right\}
\end{aligned}$$

$$\begin{aligned}
&\leq 2\mathbb{E}\left\{\sum_{d=1}^{\mathcal{D}}\mathcal{N}_{|\mathcal{S}|n}(\delta_2, \mathcal{H}_\delta^d(\mathbb{D}_n), L^\infty)\max_{u\in\mathcal{H}_\delta^d(\mathbb{D}_n)}\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\epsilon_i\sum_{d:z\neq X_i(\mathbf{t}_i)}g(u, z, Z_i)\right.\right. \\
&\quad \left.\left.>\left(\frac{3}{16}-\frac{C\alpha}{24}\right)\frac{t}{\mathcal{D}}+\frac{C\alpha}{24n}\sum_{i=1}^n\left\{\sum_{d:z\neq X_i(\mathbf{t}_i)}g(u, z, Z_i)\right\}^2\middle|\mathbb{D}_n\right)\right\} \\
&\leq 4\mathbb{E}\left\{\sum_{d=1}^{\mathcal{D}}\mathcal{N}_{|\mathcal{S}|n}(\delta_2, \mathcal{H}_\delta^d(\mathbb{D}_n), L^\infty)\right. \\
&\quad \left.\times\max_{u\in\mathcal{H}_\delta^d(\mathbb{D}_n)}\exp\left(-\frac{cn\left[\left(\frac{3}{16}-\frac{C\alpha}{24}\right)\frac{t}{\mathcal{D}}+\frac{C\alpha}{24n}\sum_{i=1}^n\left\{\sum_{d:z\neq X_i(\mathbf{t}_i)}g(u, z, Z_i)\right\}^2\right]^2}{\frac{2}{n}\sum_{i=1}^n\left\{\sum_{d:z\neq X_i(\mathbf{t}_i)}g(u, z, Z_i)\right\}^2+\frac{2t}{\mathcal{D}}}\right)\right\} \\
&\leq 4\mathcal{D}\mathcal{N}_{|\mathcal{S}|n}\left(\frac{t}{(16\eta_4+2\eta_3)\mathcal{D}}, \mathcal{G}_n, L^\infty\right)\exp\left(-\frac{c\alpha nt}{\mathcal{D}}\right),
\end{aligned}$$

where the last inequality we use the inequality

$$\frac{(a+u)^2}{a+bu}\geq\frac{4a}{b^2}[(b-1)\vee 0]\text{ for any }a, b, u > 0.$$

Consequently, we have

$$\mathbb{P}(\mathcal{B}_3)\leq 4\mathcal{D}\mathcal{N}_{|\mathcal{S}|n}\left(\frac{t}{(16\eta_4+2\eta_3)\mathcal{D}}, \mathcal{G}_n, L^\infty\right)\exp\left(-\frac{c\alpha nt}{\mathcal{D}}\right). \quad (18)$$

**Deriving the final bound:** Combining Equation 16 and Equation 18 derived in previous parts, if  $t\geq CK_1^2|\mathcal{S}|^2\mathcal{D}/n$  (note that  $\alpha^{-1}\leq CK_1^2|\mathcal{S}|$  and  $K_2\leq CK_1$ ), we have

$$\begin{aligned}
&\mathbb{P}\left(\exists u\in\mathcal{G}_n:\mathbb{E}\sum_{d=1}^{\mathcal{D}}\sum_{d:z\neq X(\mathbf{t})}g(u, z, Z)-\frac{2}{n}\sum_{i=1}^n\sum_{d=1}^{\mathcal{D}}\sum_{d:z\neq X_i(\mathbf{t}_i)}g(u, z, Z_i)>t\right) \\
&\leq\frac{1}{1-\eta_1}(2\mathbb{P}(\mathcal{B}_2)+\mathbb{P}(\mathcal{B}_3)) \\
&\leq C\mathcal{D}\mathcal{N}_{|\mathcal{S}|n}\left(\frac{t}{CK_1^2|\mathcal{S}|^2\mathcal{D}}, \mathcal{G}_n, L^\infty\right)\exp\left(-\frac{c\alpha nt}{K_1^2|\mathcal{S}|^2\mathcal{D}}\right).
\end{aligned}$$

Thus, by taking  $a_n=\frac{CK_1^2|\mathcal{S}|^2\mathcal{D}[\log\mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)+\log\mathcal{D}]}{n}$ , we can obtain that

$$\begin{aligned}
&\mathbb{E}_{\mathbb{D}}\left[\mathbb{E}_Z\left[\sum_{z\neq X(\mathbf{t})}g(\hat{u}, z, Z)\right]-\frac{2}{n}\sum_{i=1}^n\sum_{z\neq X_i(\mathbf{t}_i)}g(\hat{u}, z, Z_i)\right] \\
&\leq a_n+\int_{a_n}^{\infty}\mathbb{P}\left(\mathbb{E}_Z\left[\sum_{z\neq X(\mathbf{t})}g(\hat{u}, z, Z)\right]-\frac{2}{n}\sum_{i=1}^n\sum_{z\neq X_i(\mathbf{t}_i)}g(\hat{u}, z, Z_i)>t\right)dt \\
&\leq a_n+\int_{a_n}^{\infty}\exp\left(-\frac{c\alpha nt}{K_1^2|\mathcal{S}|^2\mathcal{D}}+\log\mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)+\log\mathcal{D}\right)dt \\
&\leq a_n+\frac{CK_1^2|\mathcal{S}|^2\mathcal{D}}{n} \\
&\leq\frac{CK_1^2|\mathcal{S}|^2\mathcal{D}[\log\mathcal{N}_{|\mathcal{S}|n}(1/(2n), \mathcal{G}_n, L^\infty)+\log\mathcal{D}]}{n},
\end{aligned}$$

which completes the proof.  $\square$

### E.3 PROOF OF THEOREM 4

*Proof.* We follow the proof of Theorem 6 (2) in Chen & Ying (2024) and Theorem 1 in Zhang et al. (2025). We have

$$\text{TV}(p_1, p_{1-\tau})\leq\mathbb{P}(X(1)\neq X(1-\tau))$$

$$\begin{aligned}
1620 & \\
1621 & = 1 - \sum_{x_1} p(x_1) \prod_{d=1}^{\mathcal{D}} p_{1|1}^d(x_1^d|x_1^d) \\
1622 & \\
1623 & = 1 - \left( \kappa_t + \frac{1 - \kappa_t}{|\mathcal{S}|} \right)^{\mathcal{D}} \\
1624 & \\
1625 & = 1 - \exp \left\{ \mathcal{D} \log \left( - (1 - \kappa_{1-\tau}) \frac{|\mathcal{S}| - 1}{|\mathcal{S}|} + 1 \right) \right\}, \\
1626 & \\
1627 & 
\end{aligned}$$

1628 which completes the proof.  $\square$

## 1633 F IMPLEMENTATION DETAILS AND ADDITIONAL EXPERIMENTS

### 1636 F.1 IMPLEMENTATION DETAILS

1637 **Data.** We consider the data distribution with a blockwise AR(1) structure; that is, we first sample  
1638 dimension  $d = 1$  from  $X(1)^1 \sim \mathcal{U}(\mathcal{S})$  and then we sample dimension  $d = 2, 3$  from

$$1640 X(1)^d | X(1)^{d-1} \sim \begin{cases} 0.9 \mathcal{U}(X(1)^{d-1} + \{-2, -1, \dots, 2\}) + 0.1 \mathcal{U}(\mathcal{S}), & \text{if } X(1)^{d-1} \in [3, |\mathcal{S}| - 2] \\ \mathcal{U}(\mathcal{S}), & \text{otherwise} \end{cases},$$

1643 and finally we sample  $(X(1)^{3j-2}, X(1)^{3j-1}, X(1)^{3j})$  from distribution same as  
1644  $(X(1)^1, X(1)^2, X(1)^3)$  for  $j = 2, \dots, \mathcal{D}/3$  (if  $\mathcal{D} > 3$ ).

#### 1646 Experimental Setup

- 1647 • sample size:  $n \in \{2500, 5000, 7500, 10000, 12500\}$ ;
- 1648 • dimension:  $\mathcal{D} \in \{3, 6, 9, 12, 15\}$ ;
- 1649 • early stopping parameter:  $\tau \in \{0.001, 0.003, 0.005, 0.007, 0.01, 0.03, 0.05, 0.07, 0.1\}$ .

1652 **Training and Evaluation.** In our experiments, we consider the linear time scheduler  $\kappa_t = t$ .  
1653 Note that  $u_t^d(z^d, x) = \frac{1}{1-t} (p_{1|t}^d(z^d|x) - \delta_{x^d}(z^d))$  by equation 3 and equation 4. Thus, we can  
1654 parameterize the posterior  $p_{1|t}^{d,\theta}(z^d|x)$ . By equation 4, the estimated transition rate is  $\hat{u}_t(z, x) =$   
1655  $\sum_{d=1}^{\mathcal{D}} \frac{1}{1-t} \delta_{x^d}(z^d) (p_{1|t}^{d,\theta}(z^d|x) - \delta_{x^d}(z^d))$ , where the parameters can be obtained by minimizing  
1656 the following empirical risk (equivalent to the empirical version of ELBO derived in Equation 37  
1657 of Shaul et al. (2025)) based on data  $\mathbb{D}_n$ , which is equivalent to equation 5 (up to a constant not  
1658 depending on  $\theta$ ):

$$1660 -\frac{1}{n} \sum_{i=1}^n \sum_{d=1}^{\mathcal{D}} \frac{1}{1-t_i} \left\{ \left( 1 - \delta_{X_i(1)^d}(X_i(\mathbf{t}_i)^d) \right) \log p_{1|t}^{d,\theta}(X_i(1)^d | X_i(\mathbf{t}_i)) - \delta_{X_i(1)^d}(X_i(\mathbf{t}_i)^d) + p_{1|t}^{d,\theta}(X_i(\mathbf{t}_i)^d | X_i(\mathbf{t}_i)) \right\}.$$

1663 Let  $\{(\mathbf{t}'_i, X'_i(\mathbf{t}'_i), X'_i(1))\}_{i=1}^{n_{\text{test}}}$  be a test dataset independent of  $\mathbb{D}_n$ . To evaluate the estimation error  
1664 of the estimated transition rate  $\hat{u}$ , we use the following empirical prediction risk

$$1666 -\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \sum_{d=1}^{\mathcal{D}} \frac{1}{1-t'_i} \left\{ \left( 1 - \delta_{X'_i(1)^d}(X'_i(\mathbf{t}'_i)^d) \right) \log p_{1|t}^{d,\theta}(X'_i(1)^d | X_i(\mathbf{t}'_i)) - \delta_{X'_i(1)^d}(X'_i(\mathbf{t}'_i)^d) + p_{1|t}^{d,\theta}(X'_i(\mathbf{t}'_i)^d | X'_i(\mathbf{t}'_i)) \right\}.$$

1669 We set  $n_{\text{test}} = 100,000$  in our simulation.

1670 **Models** All our logit models use ReLU networks with 4 hidden layers with 256 dimensions. The  
1671 optimizer is Adam with learning rate 1e-3. We train on the  $\mathcal{D}$ -dimensional dataset for  $2000\mathcal{D}/3$   
1672 epochs with batch size 512.

1673

**Algorithm 2** Sampling via Tau-leaping (Algorithm 1 in Campbell et al., 2022)

---

**Require:** A learned transition rate  $\hat{u}$ , an early stopping parameter  $\tau > 0$ , time partition  $0 = t_0 < t_1 < \dots < t_N = 1 - \tau$ .

- 1: Draw  $Y_0 \sim \mathcal{U}(\mathcal{S}^{\mathcal{D}})$ .
- 2: **for**  $k = 0$  to  $N - 1$  **do**
- 3:   **for**  $d = 1$  to  $\mathcal{D}$  **do**
- 4:     **for**  $s \in \mathcal{S} \setminus Y_k^d$  **do**
- 5:       Draw  $P_{ds} \sim \text{Poisson}((t_{k+1} - t_k)\hat{u}_{t_k}^d(s, Y_k))$ .
- 6:     **end for**
- 7:   **end for**
- 8:   **for**  $d = 1$  to  $\mathcal{D}$  **do**
- 9:     **if**  $\sum_{s \in \mathcal{S} \setminus Y_k^d} P_{ds} > 1$  **then**
- 10:        $Y_{k+1}^d = Y_k^d$
- 11:     **else**
- 12:        $Y_{k+1}^d = Y_k^d + \sum_{s \in \mathcal{S} \setminus Y_k^d} P_{ds} \times (s - Y_k^d)$
- 13:     **end if**
- 14:   **end for**
- 15: **end for**
- 16: **return**  $Y_N \sim \hat{p}_{1-\tau}$

---

## F.2 OVERALL PERFORMANCE EVALUATION

**Tau-leaping algorithm.** We present the tau-leaping algorithm (Algorithm 2) described in Campbell et al. (2022).

**Implementation Details.** We train our models with sample size 100,000. To assess the performance of uniformization and tau-leaping sampling algorithms in practice, we calculate the total variation of the empirical joint distribution of the first 3 dimensions ( $8^3 = 512$  states in total) between the samples from the true data distribution and the generated samples. We choose  $N = 100$  and the time partition  $t_i = (1 - \tau) \times i/N$  for both algorithms, and  $\lambda_{k+1} = \mathcal{D}/(1 - t_{k+1})$  for uniformization algorithm. For evaluation, we generate 500,000 samples from the true data distribution and 100,000 samples from discrete flow-based models using different sampling algorithms. We also record the runtime of each algorithm for sampling 100,000 samples.

**Overall Performance and Comparison between Uniformization and Tau-leaping.** The simulation results are presented in Table 1. From the simulation results, we can obtain the following conclusions.

- As the early stopping parameter  $\tau$  increases, the total variation decreases first and then increases. The minimum is achieved between  $\tau = 0.01$  and  $\tau = 0.07$ .
- As  $\tau$  decreases, the runtime of uniformization increases and that of tau-leaping is almost fixed. This is because in each time interval  $[t_k, t_{k+1}]$ , the number of function calls depends on  $t_{k+1}$  for uniformization algorithm and is fixed for tau-leaping algorithm.
- The uniformization sampling algorithm performs well for moderately small  $\tau$ , and is sometimes worse than tau-leaping algorithm especially for extremely small  $\tau$ . One possible explanation is that the tau-leaping algorithm might reduce the large estimation error caused by extremely small  $\tau$ .

## LLM USAGE

We only use LLMs for writing refinement. No ideas or scientific contributions are generated by LLMs.

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

Table 1: Total variation (on the joint distribution of the first 3 dimensions) and runtime with uniformization and tau-leaping algorithms.

$\mathcal{D}$	$\tau$	Total Variation		Runtime (s)	
		Uniformization	Tau-leaping	Uniformization	Tau-leaping
3	0.01	<b>0.0670</b>	0.0679	31.3777	10.6464
	0.03	<b>0.0551</b>	0.0561	30.3401	10.6581
	0.05	0.0547	<b>0.0534</b>	29.1027	10.9488
	0.07	<b>0.0590</b>	0.0608	28.6259	10.6501
	0.1	<b>0.0613</b>	0.0684	27.7543	10.6748
6	0.01	<b>0.0588</b>	0.0612	44.1496	13.3098
	0.03	<b>0.0488</b>	0.0516	40.6444	13.6291
	0.05	<b>0.0509</b>	0.0522	40.3871	13.8540
	0.07	<b>0.0594</b>	0.0639	37.7639	13.4209
	0.1	<b>0.0647</b>	0.0683	37.7448	13.3129
9	0.01	0.0643	<b>0.0628</b>	56.9279	16.0244
	0.03	<b>0.0653</b>	0.0666	53.7150	16.4715
	0.05	<b>0.0539</b>	0.0581	56.2539	15.9569
	0.07	<b>0.0581</b>	0.0612	50.1394	17.3870
	0.1	<b>0.0689</b>	0.0719	50.7919	16.6256
12	0.01	0.1061	<b>0.0980</b>	73.9635	18.8918
	0.03	<b>0.0584</b>	0.0607	64.0545	19.3838
	0.05	<b>0.0610</b>	0.0640	64.8737	18.3714
	0.07	<b>0.0647</b>	0.0672	62.9912	18.3160
	0.1	<b>0.0680</b>	0.0723	57.3696	19.2553
15	0.01	<b>0.0816</b>	<b>0.0816</b>	82.3939	21.9227
	0.03	0.0627	<b>0.0621</b>	73.7929	21.7484
	0.05	<b>0.0718</b>	0.0729	70.4029	20.3761
	0.07	<b>0.0757</b>	0.0775	68.3231	20.5265
	0.1	<b>0.0784</b>	0.0803	65.7229	21.0171