

Almost Linear Convergence under Minimal Score Assumptions: Quantized Transition Diffusion

Xunpeng Huang*

University of California San Diego

HUANGXUNPENG746@GMAIL.COM

Yingyu Lin*

University of California San Diego

YIL208@UCSD.EDU

Nikki Lijing Kuang

University of California San Diego

NIKKI.KUANG@GMAIL.COM

Hanze Dong

Microsoft Research

HENDRYDONG@GMAIL.COM

Difan Zou[†]

University of Hong Kong

DZOU@CS.HKU.HK

Yian Ma[†]

University of California San Diego

YIANMA@UCSD.EDU

Tong Zhang

University of Illinois Urbana-Champaign

TONGZHANG@TONGZHANG-ML.ORG

Editors: Steve Hanneke and Tor Lattimore

Abstract

Continuous diffusion models have demonstrated remarkable generative performance across diverse domains but are often constrained by the computational cost of simulating reverse Ornstein–Uhlenbeck processes via SDE/ODE solvers. Existing theoretical results typically establish query complexities that scale polynomially with both the dimension d and the error tolerance ϵ (e.g., $\tilde{O}(d/\epsilon)$). This mirrors the limitations of unadjusted Langevin algorithm, where standard first-order score solvers lack access to zeroth-order density information, precluding natural error-correction mechanisms and thus preventing the fast $\ln(1/\epsilon)$ convergence attainable by Metropolis-adjusted methods. In this paper, we develop an improved generative modeling method by introducing Quantized Transition Diffusion (QTD), a framework that reformulates continuous diffusion into a discrete generation problem through spatial quantization and the parameterization of zeroth-order information (e.g., density ratios). To sample from this discrete target, we propose a truncated uniformization algorithm that simulates the underlying continuous-time Markov chain of the discrete diffusion process without discretization error, while eliminating the restrictive bounded-score assumption required by prior uniformization-based approaches. Consequently, QTD attains ϵ -accuracy in total variation distance with a query complexity of $\mathcal{O}(d \ln^2(d/\epsilon))$, yielding a notable improvement in ϵ -dependence compared to existing continuous diffusion samplers. Crucially, our analysis capitalizes on a novel proof technique based on the infinitesimal chain rule of KL divergence, providing a fresh perspective on unifying continuous and discrete diffusion paradigms.

Keywords: Quantized diffusion, Truncated uniformization, Linear convergence

*. Equal contribution.

†. Corresponding author.

1. Introduction

Diffusion models (Sohl-Dickstein et al., 2015; Song and Ermon, 2019; Ho et al., 2020) have emerged as a powerful and widely adopted class of generative models, achieving state-of-the-art performance across diverse domains. Their core mechanism is a noising–denoising procedure: the forward process incrementally corrupts training data by adding noise, thereby mapping an unknown and potentially complex distribution to a tractable prior (e.g., a standard Gaussian), while the reverse process progressively denoises samples back toward the original data distribution by estimating the logarithmic gradient (*score*) of the noised distributions (Vincent, 2011; Song and Ermon, 2019). Despite their remarkable empirical success, advancing our theoretical understanding of these models—particularly regarding how to generate high-quality samples in high-dimensional settings with reduced computational cost—remains a significant challenge.

Current theoretical analyses of continuous diffusion samplers in \mathbb{R}^d primarily address the simulation of reverse Ornstein–Uhlenbeck (OU) processes via SDE/ODE solvers (Chen et al., 2023b,a; Benton et al., 2024a; Li and Yan, 2024). These approaches yield complexities that are *polynomial* in the total variation error tolerance ϵ , ranging from $\tilde{O}(d/\epsilon^2)$ to $\tilde{O}(d^{1/3}/\epsilon^{2/3})$ (Chen et al., 2023a; Huang et al., 2024b; Li and Jiao, 2025). Due to the deep connections between diffusion-based inference and Monte Carlo methods (Huang et al., 2024a; He et al., 2024), the fundamental limitation of these approaches stems from the fact that standard diffusion samplers operate solely on first-order scores without access to zeroth-order density information; consequently, they cannot employ natural error-correction mechanisms, such as a Metropolis–Hastings filter, to eliminate the discretization bias introduced by Euler-type approximations. As a result, convergence is inherently constrained by the discretization step size, precluding the fast exponential rates attainable in high-precision MCMC. Although several recent works (Chen et al., 2026; Gatmiry et al., 2026) attempt to match the $\ln(1/\epsilon)$ TV-convergence rate of Metropolis-adjusted MCMC (see Table 2), they rely on strong assumptions—such as restricting the target distribution to the convolution of a bounded distribution and a Gaussian, or imposing smoothness conditions on the forward marginal scores—that limit direct comparison with existing methods requiring only minimal score assumptions. The challenge of reconciling fast convergence with mild score assumptions has motivated research into discrete diffusion on finite state spaces (e.g., $\{1, 2, \dots, K\}^d$), where density ratios (zeroth-order information) are explicitly parameterized. Chen and Ying (2024) demonstrate that the continuous-time Markov chain (CTMC) governing the inference process can be simulated via *uniformization*—a Poisson process simulation technique—using only $\mathcal{O}(d \ln(d/\epsilon))$ discrete score evaluations. Nonetheless, their analysis relies on the restrictive assumption that discrete score estimators remain bounded; allowing unbounded scores can lead to uncontrollable transition rates and redundant updates, thereby hindering convergence. Given these obstacles, two fundamental questions arise:

1. Can uniformization provide a *linear-convergence* alternative to continuous diffusion samplers without incurring significant computational overhead?
2. Is it possible to remove the bounded-score assumption for discrete estimators without compromising efficiency?

In this paper, we introduce QTD, a novel **Quantized Transition Diffusion** framework that directly addresses these challenges, achieving the convergence of the total variation (TV) distance in $\mathcal{O}(d \ln^2(d/\epsilon))$ expected score evaluations under minimal assumptions. To achieve linear convergence, we leverage zeroth-order information by adapting the density-ratio parameterization of

discrete diffusion to the continuous setting. Specifically, we partition high-probability regions of the target distribution into fine hypercubes and approximate the local density within each as uniform. By sufficiently refining these partitions, we ensure the induced histogram-like distribution approximates the true continuous distribution arbitrarily well in TV distance, effectively reducing the continuous sampling task to a discrete generation problem. Crucially, the availability of parameterized zeroth-order information allows for effective error correction; unlike Euler-type ODE/SDE solvers, where discretization error grows with step size, our method ensures that errors remain negligible even over long-distance state transitions. Furthermore, to eliminate the reliance on bounded discrete scores, we introduce *truncated uniformization*. This technique constrains the outgoing transition rates of each state, ensuring that the effective discrete scores remain naturally bounded. We prove that this modification preserves the *linear inference convergence* rate of standard uniformization. Our key contributions are summarized as follows:

- We propose the QTD framework and provably improve the inference rate from polynomial to logarithmic dependence on ϵ . Specifically, QTD generate d -dimensional samples to approximate the data distribution with $\Theta(\epsilon)$ -TV error with only $\mathcal{O}(d \ln^2(d/\epsilon))$ expected score evaluations.
- We provide a constructive framework that models continuous distributions via efficient discrete diffusion. Under this framework, employing an efficient uniformization-based sampler ensures the output distribution remains arbitrarily close to the continuous target distribution.
- We introduce the *truncated uniformization* technique for an unbiased and tractable CTMC simulation. This method removes the restricted bounded-score assumption imposed in prior discrete diffusion analyses (Chen and Ying, 2024; Zhang et al., 2024).
- We develop a novel proof technique which is simpler and more general than Girsanov-based analyses used in (Chen and Ying, 2024; Zhang et al., 2024) for the discrete diffusion inference. In particular, we leverage the chain rule of KL divergence over infinitesimal time intervals to derive convergence guarantees.

2. Preliminaries

Given that our strategy rests on quantizing the continuous target distribution to obtain an approximation supported on a discrete space, we begin this section by formally defining the discrete forward and reverse Markov processes, parameterized by transition rate matrices. Subsequently, we detail the construction of noising and denoising dynamics on this discrete distribution for the purposes of training and inference. Furthermore, we introduce uniformization (van Dijk, 1992; van Dijk et al., 2018) as a technique for the efficient simulation of these processes. For reference, all notations introduced herein are summarized in Table 3 in Appendix A.

Problem setup. Without loss of generality, we focus on distributions that admit probability density functions in Euclidean space. These continuous density functions are represented by $p: \mathbb{R}^d \rightarrow \mathbb{R}^+$. Specifically, let the data distribution be $p_* \propto \exp(-f_*)$ for some potential function f_* . We consider the task of approximating p_* using some discrete distribution with probability mass function $q_*: \mathcal{Y} \rightarrow \mathbb{R}_0^+$, defined on a finite discrete space \mathcal{Y} . This discrete approximation is modeled via a forward Markov process $\{\mathbf{y}_t^{\rightarrow}\}_{t=0}^T$ and named as discrete diffusion model (Lou et al., 2024; Zhang et al., 2024; Chen and Ying, 2024), with initial distribution $q_0^{\rightarrow} = q_*$ that evolves toward

the uniform distribution. Then, the marginal distribution at time t is denoted by q_t^\rightarrow , the joint and conditional distributions over different time steps $t' > t$ are given by

$$(\mathbf{y}_{t'}^\rightarrow, \mathbf{y}_t^\rightarrow) \sim q_{t',t}^\rightarrow \quad \text{and} \quad q_{t'|t}^\rightarrow(\mathbf{y}'|\mathbf{y}) = q_{t',t}^\rightarrow(\mathbf{y}', \mathbf{y})/q_t^\rightarrow(\mathbf{y}).$$

For simplicity, we set the forward process to be a time-homogeneous CTMC constructed via the transition rate function $R^\rightarrow: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, which implies that both conditional and marginal distributions satisfy

$$\frac{dq_{t|s}^\rightarrow}{dt}(\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R^\rightarrow(\mathbf{y}, \mathbf{y}') \cdot q_{t|s}^\rightarrow(\mathbf{y}') \quad \text{and} \quad \frac{dq_t^\rightarrow}{dt}(\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R^\rightarrow(\mathbf{y}, \mathbf{y}') \cdot q_t^\rightarrow(\mathbf{y}'). \quad (1)$$

The transition rate function R^\rightarrow characterizes the instantaneous rate of transitioning from state \mathbf{y}' to \mathbf{y} and is formally defined as

$$R^\rightarrow(\mathbf{y}, \mathbf{y}') := \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \left(q_{\Delta t|0}^\rightarrow(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'}(\mathbf{y}) \right) \right], \quad (2)$$

where $\delta_{\mathbf{y}'}(\mathbf{y}) = 1$ if $\mathbf{y} = \mathbf{y}'$ and 0 otherwise.

Reverse process. Additional properties of R^\rightarrow are discussed in Appendix B.1. To sample from the target distribution $q_* = q_0^\rightarrow$ in practice, we simulate the reverse-time process \mathbf{y}_t^\leftarrow that starts from q_T^\rightarrow and moves backward.

$$\{\mathbf{y}_t^\leftarrow\}_{t=0}^T \quad \text{where} \quad \mathbf{y}_t^\leftarrow \sim q_t^\leftarrow = q_{T-t}^\rightarrow, \quad (\mathbf{y}_{t'}^\leftarrow, \mathbf{y}_t^\leftarrow) \sim q_{t',t}^\leftarrow, \quad \text{and} \quad q_{t',t}^\leftarrow = q_t^\leftarrow \cdot q_{t'|t}^\leftarrow,$$

whose dynamic follows from

$$\frac{dq_t^\leftarrow}{dt}(\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R_t^\leftarrow(\mathbf{y}, \mathbf{y}') \cdot q_t^\leftarrow(\mathbf{y}') \quad \text{where} \quad R_t^\leftarrow(\mathbf{y}, \mathbf{y}') := R^\rightarrow(\mathbf{y}', \mathbf{y}) \cdot \frac{q_t^\leftarrow(\mathbf{y})}{q_t^\leftarrow(\mathbf{y}')}, \quad (3)$$

proven in Appendix B.2. Similar to the R^\rightarrow in the forward process, R_t^\leftarrow characterizes the transition rates for the time-inhomogeneous reverse process $\{\mathbf{y}_t^\leftarrow\}_{t=0}^T$, i.e.,

$$R_t^\leftarrow(\mathbf{y}, \mathbf{y}') := \lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'}(\mathbf{y})}{\Delta t^{-1}} \right], \quad (4)$$

as shown in Appendix B.3. In practice, the probability density ratio $q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$ will usually be approximated with neural networks due to its unknown closed form, which is presented as

$$\tilde{v}_{t,\mathbf{y}'}(\cdot) \approx v_{t,\mathbf{y}'}(\cdot) = q_t^\leftarrow(\cdot)/q_t(\mathbf{y}').$$

To simulate the reverse process in Eq. (3), we must estimate the time-varying rate matrix R_t^\leftarrow , which depends on the intractable ratio $q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$. We approximate this ratio using neural networks trained via score entropy minimization (Lou et al., 2024; Benton et al., 2024b),

$$L_{\text{SE}}(\hat{v}) = \int_0^T \mathbb{E}_{\mathbf{y}_t \sim q_t^\rightarrow} \left[\sum_{\mathbf{y} \neq \mathbf{y}_t} R^\rightarrow(\mathbf{y}_t, \mathbf{y}) \cdot D_\phi(v_{T-t,\mathbf{y}_t}(\mathbf{y}) \parallel \tilde{v}_{T-t,\mathbf{y}_t}(\mathbf{y})) \right] dt, \quad (5)$$

where $D_\phi(\cdot \parallel \cdot)$ denotes the Bregman divergence, and $\phi(c) = c \ln c$. Similar to the score estimation loss in continuous cases (Chen et al., 2023b), the loss L_{SE} is not directly estimable. Instead, implicit score entropy and denoising score entropy (Lou et al., 2024; Benton et al., 2024b) are introduced to enable an equivalent minimization.

Uniformization. Uniformization (van Dijk, 1992; van Dijk et al., 2018) is a sampling technique that enables the exact simulation of the discrete reverse process by a Poisson process (Chen and Ying, 2024). Consider the reverse process shown in Eq. (3), the conditional transition probability is

$$q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) = \begin{cases} \Delta t \cdot R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \Delta t \sum_{\tilde{\mathbf{y}} \neq \mathbf{y}} R_t^{\leftarrow}(\tilde{\mathbf{y}}, \mathbf{y}) & \mathbf{y}' = \mathbf{y} \end{cases} \quad (6)$$

in an infinitesimal time Δt due to Eq. (4), where the $o(\Delta t)$ term is omitted. Suppose the probability of transitioning to a different state is upper bounded by $\Delta t \cdot \beta$:

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) := R_t^{\leftarrow}(\mathbf{y}) \leq \beta, \quad \forall t. \quad (7)$$

We can then simulate Eq. (3) with the following uniformization procedure:

1. With probability $\Delta t \cdot \beta$, allow a state transition.
2. Conditioning on an allowed transition, move from \mathbf{y} to \mathbf{y}' with probability

$$M_t(\mathbf{y}'|\mathbf{y}) = \begin{cases} \beta^{-1} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \beta^{-1} R_t^{\leftarrow}(\mathbf{y}) & \text{otherwise} \end{cases}.$$

Under these two steps, the practical conditional probability satisfies

$$\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) = \begin{cases} \Delta t \cdot \beta \cdot R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \beta^{-1} = \Delta t \cdot R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \Delta t \cdot \beta + \Delta t \cdot \beta \cdot (1 - \beta^{-1} \cdot R_t^{\leftarrow}(\mathbf{y})) = 1 - \Delta t \cdot R_t^{\leftarrow}(\mathbf{y}) & \mathbf{y}' = \mathbf{y}, \end{cases} \quad (8)$$

which exactly matches Eq. (6). Under this condition, the number of transition events within a time interval $[s, t]$ follows a Poisson distribution (van Dijk, 1992; van Dijk et al., 2018) whose expectation is $\beta(t - s)$, which coincides with the number of required evaluations of the transition rate function R_t^{\leftarrow} . This implies choosing a tighter upper bound β directly leads to better complexity.

3. Quantized Transition Diffusion

In this section, we introduce Quantized Transition Diffusion (QTD), a framework designed for efficient sampling from continuous data distributions. At a high level, QTD approximates the continuous target via a histogram-like distribution, thereby transforming the continuous generation task into a sampling problem over a discrete CTMC defined on a hypercube partition. Then, we construct the noising dynamics based on the adjacency structure induced by a binary encoding of this partition. This scheme effectively balances the trade-off between the out-degree of each state and the global connectivity of the transition graph, ensuring rapid forward mixing while maintaining low iteration complexity for the inference updates. Finally, we propose a *truncated uniformization* scheme to simulate the reverse-time CTMC. Crucially, this approach enables efficient simulation without relying on the restrictive boundedness assumption on discrete scores, a limitation prevalent in prior uniformization-based inference algorithms (Chen and Ying, 2024; Liang et al., 2025).

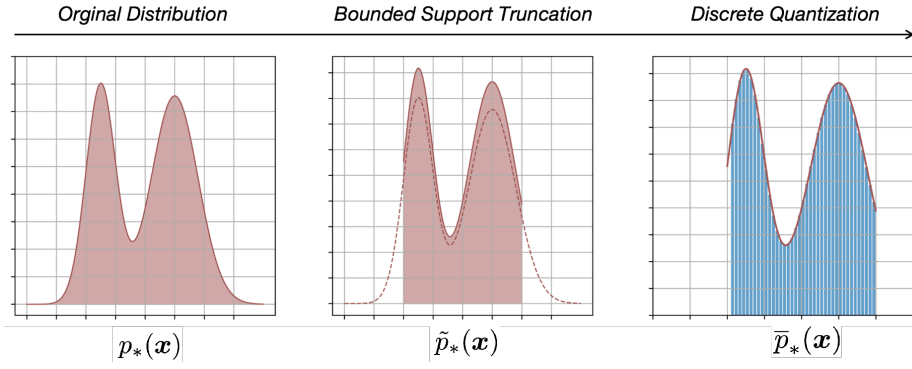


Figure 1: Visualization of the histogram approximation. The first step regularizes the original distribution in some bounded sets but controls the TV gap by Lemma 10. The second step quantizes the probability density to a histogram-like distribution but controls the TV gap by Lemma 11.

3.1. Histogram Approximation

To approximate the target distribution p_* defined in the Euclidean space \mathbb{R}^d with a histogram-like distribution, we first restrict its support to a bounded region, which can be represented by a cube of side length L as follows:

$$\text{Cube}(L) := \{\mathbf{x} \mid -L \leq \mathbf{x}_j \leq L, \forall j \in \{1, 2, \dots, d\}\}.$$

Given that $\text{Cube}(L)$ covers most probability mass of p_* , we construct a probability density restricted to this region to approximate p_* :

$$\tilde{p}_*(\mathbf{x}) := \frac{p_*(\mathbf{x})}{\int_{\mathbf{x} \in \text{Cube}(L)} p_*(\mathbf{x}) d\mathbf{x}} \quad \forall \mathbf{x} \in \text{Cube}(L). \quad (9)$$

Standard concentration arguments allow us to control the TV distance between p_* and \tilde{p}_* . Next, we quantize \tilde{p}_* over $\text{Cube}(L)$ by discretizing each dimension into $K := 2L/l$ intervals of width l , with partition points defined by:

$$l_i = -L + i \cdot l \quad i \in \{0, 1, \dots, K\} \quad \text{and} \quad -L \leq l_i \leq L.$$

That means the high-dimensional cube $\text{Cube}(L)$ will be decomposed into K^d cells (subsets), and each cell will cover a small region shown as follows

$$\text{Cell}(i_0, i_1, \dots, i_{d-1}) := \{\mathbf{x} \mid l_{i_j} < \mathbf{x}_j \leq l_{i_j+1}, \forall j \in \{0, 1, \dots, d-1\}\}. \quad (10)$$

We construct the piecewise constant distribution $\bar{p}_*(\mathbf{x})$ by averaging the original density \tilde{p}_* over each quantization cell. Specifically, for each cell $\text{Cell}(i_0, i_1, \dots, i_{d-1})$, we assign a constant density to all points \mathbf{x} in the cell:

$$\bar{p}_*(\mathbf{x}) = l^{-d} \cdot \int_{\mathbf{u} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} \tilde{p}_*(\mathbf{u}) d\mathbf{u}, \quad \mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1}). \quad (11)$$

By construction, the quantized distribution satisfies the normalization condition $\int_{\text{Cube}(L)} \bar{p}_*(\mathbf{x}) d\mathbf{x} = 1$. The following lemma establishes a bound on the TV distance between the approximation \bar{p}_* and

the ground truth p_* . This error analysis relies on two key properties: first, the sub-Gaussianity of the target distribution ensures that the bounded domain $\text{Cube}(L)$ captures almost all the probability mass; second, the smoothness of p_* ensures that the histogram-like approximation incurs a controllable discretization error. Consequently, with an appropriate selection of the domain size L and grid resolution l , \bar{p}_* converges to the continuous target distribution p_* .

Lemma 1 *Suppose the target distribution $p_* \propto \exp(-f_*)$ is σ sub-Gaussian and f_* is H -smooth, we can construct \bar{p}_* defined on a finite cube $\text{Cube}(L)$ with length*

$$L = \sigma \cdot \sqrt{2 \ln(2d/\epsilon)} \quad \text{and} \quad l = \left[2H \cdot \left(\sigma \sqrt{2d \ln(2d/\epsilon)} + d + \sqrt{dm_0} \right) \right]^{-1} \cdot \epsilon,$$

to satisfy $\text{TV}(p_*, \bar{p}_*) \leq 3\epsilon$.

We defer the proof to Appendix C. Under this condition, we have constructed a histogram-like distribution to approximate p_* , which can be visualized by Fig. 1.

Under this transformation, although distribution \bar{p}_* remains defined on \mathbb{R}^d , due to the histogram shape, it can be sampled by introducing a discrete distribution \bar{q}_* :

$$\bar{q}_*(\mathbf{y}) \propto \bar{p}_*(-L \cdot \mathbf{1} + l \cdot (\mathbf{y} - 0.5 \cdot \mathbf{1})), \quad \text{where} \quad \mathbf{y} \in \bar{\mathcal{Y}} := \{0, 1, \dots, K-1\}^d. \quad (12)$$

Here, each cell $\text{Cell}(\mathbf{y})$ in the partition is treated as a discrete state whose probability mass is proportional to the probability density of \bar{p}_* at the midpoint of that cell. Consequently, sampling from \bar{p} reduces to the following two-stage procedure:

1. Sample from the discrete distribution \bar{q}_* defined on $\bar{\mathcal{Y}}$;
2. Uniformly draw a sample from the cell, i.e., $\text{Cell}(\mathbf{y})$.

Thus far, we have reformulated the continuous generation task as a sampling problem defined over a hypercube partition. To address this within the framework of discrete diffusion, the remaining objective is to construct appropriate noising and denoising dynamics—modeled as a CTMC—for the quantized distribution \bar{q}_* .

3.2. The Discrete Forward Process and Binary Encoding

As established in Section 2, the forward process on the discrete state space $\bar{\mathcal{Y}}$ is governed by the transition rate matrix R^\rightarrow . The design of R^\rightarrow necessitates a fundamental trade-off between *reverse iteration complexity* and *forward mixing time*. On one hand, the computational cost of the reverse update scales with the number of non-zero entries in R_i^\leftarrow . By Eq. (3), R_i^\leftarrow inherits the sparsity structure of R^\rightarrow ; thus, a sparse R^\rightarrow is essential for efficient inference. On the other hand, the mixing time quantifies the rate of state space traversal; generally, a denser R^\rightarrow facilitates faster global exploration, standing in direct tension with the sparsity requirement.

Figure 2 illustrates this tension. Consider two extreme topologies for a d -dimensional space with K bins per dimension. Restricting transitions to immediate neighbors yields a favorable $\mathcal{O}(d)$ iteration complexity but a prohibitive $\mathcal{O}(Kd)$ mixing time, as the process must traverse the support incrementally. Conversely, a fully connected topology ensures rapid mixing but incurs an intractable $\mathcal{O}(K^d)$ computational cost per iteration. We therefore seek a middle ground that balances these two factors, targeting both iteration complexity and mixing time of $\mathcal{O}(d \log K)$.

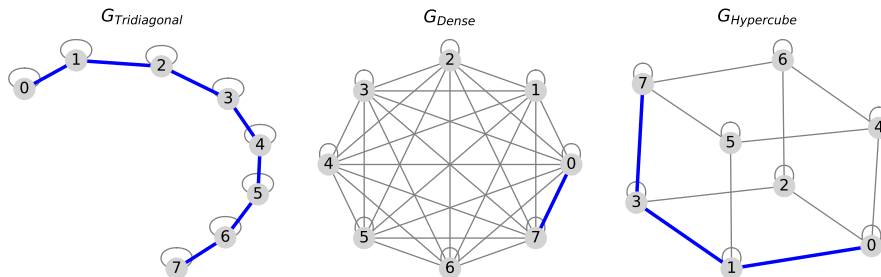


Figure 2: Comparison of transition graph topologies for discrete diffusion. We illustrate the connectivity of 8 discrete states ($d = 1, K = 8$) under different graph priors. Left ($G_{\text{Tridiagonal}}$): A simple neighbor-based chain yields low iteration complexity (degree 2) but suffers from a large diameter, leading to slow mixing. Center (G_{Dense}): A fully connected graph ensures one-step mixing but incurs a prohibitive $\mathcal{O}(K^d)$ cost. Right ($G_{\text{Hypercube}}$): Our proposed binary encoding (Section 3.2) maps states to a Boolean hypercube. This achieves an optimal balance, where both the out-degree (sparsity) and graph diameter (mixing) scale logarithmically with the number of bins, specifically $\mathcal{O}(d \log K)$. The bold blue edges highlight a diameter path—a shortest path between the two most distant vertices in each graph.

Specifically, we achieve this balance via Alg. 1 based on binary encoding. Assuming $\log_2 K$ is an integer, we define a bijection $\text{vBin}: \bar{\mathcal{Y}} \rightarrow \mathcal{Y} := \{0, 1\}^{d \log_2 K}$ that maps the original state $\bar{\mathbf{y}} = [\bar{y}_0, \dots, \bar{y}_{d-1}] \in \bar{\mathcal{Y}}$ to a binary vector

$$\mathbf{y} = [\mathbf{y}_0, \dots, \mathbf{y}_{d \log_2 K - 1}], \quad \text{where } \mathbf{y}_i = \lfloor \bar{y}_{\lfloor i / \log_2 K \rfloor} / 2^{(i - \lfloor i / \log_2 K \rfloor)} \rfloor \bmod 2.$$

This embedding transforms the geometry into a Boolean hypercube. By restricting transitions to states with a Hamming distance of 1, we limit connectivity to immediate bit flips. Consequently, the neighbor count scales as $\mathcal{O}(d \log K)$, ensuring efficient reverse updates, while the graph diameter is similarly bounded by $\mathcal{O}(d \log K)$, facilitating rapid mixing. Formally, the forward transition rate $R^\rightarrow: \mathcal{Y} \times \mathcal{Y}$ on this binary space is defined as:

$$R^\rightarrow(\mathbf{y}, \mathbf{y}') = \begin{cases} 1 & \text{Ham}(\mathbf{y}, \mathbf{y}') = 1 \\ -d \log_2 K & \mathbf{y} = \mathbf{y}' \\ 0 & \text{otherwise} \end{cases}. \quad (13)$$

This definition yields a symmetric CTMC that behaves as a time-homogeneous random walk on the hypercube. As established in Lemma 2 (proof in Appendix D), the process exhibits exponential convergence, mirroring the behavior of continuous OU processes.

Lemma 2 Suppose the transition rate function R^\rightarrow of the CTMC $\{\mathbf{y}_t^\rightarrow\}_{t=0}^T$ is defined by Eq. (13). Then, the underlying distribution q_t^\rightarrow of \mathbf{y}_t^\rightarrow satisfies $\text{KL}(q_t^\rightarrow \| q_\infty^\rightarrow) \leq e^{-t} \cdot d \log_2 K$.

3.3. Truncated Uniformization

Recall from Section 2 that the computational complexity of simulating a CTMC via uniformization is governed by the uniformization constant β (defined in Eq. (7)), which bounds the total transition rate out of any state. A tighter bound directly translates to improved simulation efficiency. Motivated by this, we analyze the time-dependent bound β_t for the time-inhomogeneous reverse CTMC

Algorithm 1 TRAINING DATA QUANTIZATION

- 1: **Input:** The training set $\mathcal{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$.
- 2: Initialize output set $\mathcal{Y} = \{\}$ and the parameters, e.g., L and l as shown in Lemma 1.
- 3: **for** $n = 1$ **to** N **do**
- 4: Quantize the training sample $\mathbf{x}^{(n)}$ to $\bar{\mathbf{y}}^{(n)}$ via

$$\bar{\mathbf{y}}^{(n)} = [\bar{\mathbf{y}}_0^{(n)}, \bar{\mathbf{y}}_1^{(n)}, \dots, \bar{\mathbf{y}}_{d-1}^{(n)}] \quad \text{where} \quad \bar{\mathbf{y}}_i^{(n)} = \lfloor (\mathbf{x}_i^{(n)} + L)/l \rfloor.$$

- 5: Append the set \mathcal{Y} with binary encoded $\mathbf{y}^{(n)} = \text{vBin}(\bar{\mathbf{y}}^{(n)})$ where $\mathbf{y}^{(n)} \in \{0, 1\}^{d \log_2 K}$.
 - 6: **end for**
 - 7: **return** \mathcal{Y} .
-

defined in Eq. (3). Specifically, when the forward transition rates are set according to Eq. (13), the resulting β_t satisfies the following lemma. The detailed proof is provided in Appendix E.1.

Lemma 3 Consider a CTMC whose transition rate function R^\rightarrow is defined as Eq. (13). Then, for any \mathbf{y} , the reverse transition rate function satisfies

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^\leftarrow(\mathbf{y}', \mathbf{y}) := R_t^\leftarrow(\mathbf{y}) \leq \beta_t := (2d \log_2 K) \cdot \max\{1, (T - t)^{-1}\}. \quad (14)$$

Therefore, it is important for us to divide the entire reverse process into W segments. With a proper segmentation, we can assign a tight upper bound β_{t_w} for $R_t^\leftarrow(\mathbf{y})$ when $t \in [t_{w-1}, t_w)$ and minimize the expectation of transition events, given by $\sum_w \beta_{t_w} \cdot (t_w - t_{w-1})$. In practice, since the exact form R_t^\leftarrow is intractable, we approximate it by minimizing Eq. (5):

$$R_t^\leftarrow(\mathbf{y}, \mathbf{y}') \approx \tilde{R}_t(\mathbf{y}, \mathbf{y}') = R^\rightarrow(\mathbf{y}', \mathbf{y}) \cdot \tilde{v}_{t, \mathbf{y}'}(\mathbf{y}).$$

Here, $\tilde{v}_{t, \mathbf{y}'}(\mathbf{y})$ can approximate the ideal density ratio $q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$ with high accuracy. However, this approximation may violate the desired global rate bound in Lemma 3. To address this, prior work (Chen and Ying, 2024; Liang et al., 2025) impose an estimated score boundedness assumption for discrete diffusion inference:

$$\sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_t(\mathbf{y}, \mathbf{y}') \leq Cd \log_2 K \cdot \max\{1, (T - t)^{-1}\} \quad (15)$$

We argue that this assumption can be safely removed by truncating the approximate transition rate function as follows:

$$\hat{R}_t(\mathbf{y}, \mathbf{y}') = \begin{cases} \tilde{R}_t(\mathbf{y}, \mathbf{y}') \cdot \beta_t / \tilde{R}_t(\mathbf{y}') & \tilde{R}_t(\mathbf{y}') > \beta_t \\ \tilde{R}_t(\mathbf{y}, \mathbf{y}') & \text{otherwise.} \end{cases}, \quad \forall \mathbf{y}' \neq \mathbf{y}, \quad (16)$$

and

$$\hat{R}_t(\mathbf{y}', \mathbf{y}') = - \sum_{\mathbf{y} \neq \mathbf{y}'} \hat{R}_t(\mathbf{y}, \mathbf{y}'). \quad (17)$$

It ensures that the total outgoing rate from any state does not exceed β_t , hence eliminating the need for explicit score bounds. Combining \hat{R}_t with the two-step uniformization mentioned in Section 2, we obtain a practical and efficient inference algorithm, summarized in Alg. 2. Here, e_i denotes the one-hot vector with a 1 at position i and 0 elsewhere, and mod is an element-wise operator.

Algorithm 2 INFERENCE PROCESS WITH TRUNCATED UNIFORMIZATION

- 1: **Input:** Total time T , a time partition $0 = t_0 < \dots < t_W = T - \delta$, parameters $\beta_{t_1}, \dots, \beta_{t_W}$ set as Eq. (14), a reverse transition rate function \hat{R}_t^{\leftarrow} obtained by the learnt score function $\tilde{v}_{t, \mathbf{y}'}(\cdot)$.
 - 2: Draw an initial sample $\hat{\mathbf{y}}_{t_0} \sim \text{Uniform}(\{0, 1\}^{d \log_2 K})$.
 - 3: **for** $w = 1$ **to** W **do**
 - 4: Draw $N \sim \text{Poisson}(\beta_{t_w}(t_w - t_{w-1}))$;
 - 5: Sample N points i.i.d. uniformly from $[t_{w-1}, t_w]$ and sort them as $\tau_1 < \tau_2 < \dots < \tau_N$;
 - 6: Set $\mathbf{z}_0 = \hat{\mathbf{y}}_{t_{w-1}}$;
 - 7: **for** $n = 1$ **to** N **do**
 - 8: Set

$$\mathbf{z}_n = \begin{cases} (\mathbf{z}_{n-1} + \mathbf{e}_i) \bmod 2, & w.p. \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n}(\mathbf{z}_{n-1} + \mathbf{e}_i, \mathbf{z}_{n-1}), \quad 0 \leq i \leq d \log_2 K - 1 \\ \mathbf{z}_{n-1}, & w.p. 1 - \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n}(\mathbf{z}_{n-1}). \end{cases}$$
 - 9: **end for**
 - 10: Set $\hat{\mathbf{y}}_{t_w} = \mathbf{z}_N$.
 - 11: **end for**
 - 12: Recover the cell index with $\bar{\mathbf{y}} = \text{vBin}^{-1}(\hat{\mathbf{y}}_{t_W})$ and uniformly draw a sample $\hat{\mathbf{x}}$ from $\text{Cell}(\bar{\mathbf{y}})$.
 - 13: **return** $\hat{\mathbf{x}}$.
-

Empirical Results. We empirically validate that QTD achieves acceleration over standard DDPM baselines based on reversing the OU process, using both synthetic and toy real-world datasets. Due to space constraints, detailed experimental results are deferred to Appendix G.

4. Theoretical Results

In this section, we first analyze the convergence of the discrete generated distribution $\hat{q}_{T-\delta}$ in Alg. 2, demonstrating that it converges to the discrete target distribution \bar{q}_* in KL divergence, contingent on accurate score estimation. A central technical innovation is the application of the infinite-time chain rule to bound the KL divergence between the marginal distributions of the ideal and approximated reverse processes. This technique serves as a robust alternative to the discrete Girsanov theorem, effectively circumventing the restrictive bounded score assumption. Subsequently, we extend these results to the continuous domain. Leveraging standard regularity assumptions from the diffusion model literature, we establish that the QTD-based method (Alg. 2) achieves ϵ -accuracy in TV distance to the continuous target p_* only using $\tilde{O}(d)$ discrete score evaluations. Finally, we provide a theoretical comparison between the proposed scheme and existing continuous inference algorithms, highlighting the superior scaling properties of our approach.

Truncated Uniformization Results. To analyze the convergence and the query complexity required to achieve a target TV distance in the discrete setting, we introduce the following assumption regarding the accuracy of the discrete score estimation:

[A1] Score Estimation Error. Let \mathcal{Y} be the discrete training set obtained by quantizing the continuous set \mathcal{X} via Alg. 1. We assume the discrete score function \tilde{v}_t , trained according to Eq. (5), satisfies a bounded estimation error condition, i.e., $L_{\text{SE}}(\hat{v}) \leq \epsilon_{\text{score}}^2$.

Lemma 4 Under Assumption [A1], if the reverse process is implemented via algorithm 2, then:

$$\text{KL} (q_{T-\delta}^{\leftarrow} \parallel \hat{q}_{T-\delta}) \leq \text{KL} (q_0^{\leftarrow} \parallel \hat{q}_0) + (T - \delta) \epsilon_{\text{score}}^2$$

This lemma establishes that the introduction of truncation mechanisms does not impede the convergence of the uniformization-type sampler toward the target distribution $q_{T-\delta}^{\leftarrow} \approx q_*$. A key challenge is that the standard discrete Girsanov Theorem is inapplicable to \hat{R}_t (Eq. (16)) due to the truncation operator. To address this, the proof of Lemma 4 follows the sketch outlined below:

1. First, we analyze the time derivative of the KL divergence using the chain rule, and then exchange the limit and expectation operations:

$$\frac{d\text{KL} (q_t^{\leftarrow} \parallel \hat{q}_t)}{dt} \leq \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \underbrace{\lim_{\Delta t \rightarrow 0} \left[\frac{\text{KL} (q_{t+\Delta t|t}^{\leftarrow}(\cdot|\mathbf{y}) \parallel \hat{q}_{t+\Delta t|t}(\cdot|\mathbf{y}))}{\Delta t} \right]}_{\text{Term 1}}.$$

2. The limit of the conditional KL divergence gap, i.e., Term 1, can be decomposed into two components, each relating to the discrepancy between the ideal and practical reverse transition rate matrices (R_t^{\leftarrow} and \hat{R}_t):

$$\lim_{\Delta t \rightarrow 0} \left[\sum_{\mathbf{y}' \neq \mathbf{y}} \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] = \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})},$$

and

$$\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right) \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] = \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}).$$

3. This formulation allows the truncation effects to be seamlessly incorporated into our analysis while coupling naturally with the learned transition rate matrix \tilde{R}_t employed during training:

$$\begin{aligned} & \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \\ &= \underbrace{\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y})}_{\text{related to } \epsilon_{\text{score}}} + \underbrace{\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{\tilde{R}_t(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - \tilde{R}_t(\mathbf{y})}_{\text{Term 2}}. \end{aligned}$$

4. Term 2 depends solely on the truncation operator. By exploiting the connection between \tilde{R}_t and \hat{R}_t , we show that this term is provably non-positive, thereby obviating the need for a bounded score assumption and completing the analysis.

Since Alg. 2 serves as a self-contained sampling algorithm for discrete diffusion models, we provide the following convergence guarantee for the truncated uniformization sampler. We omit its proof, as it is nearly identical to that of Theorem 6.

Reference	Algorithm	Assumptions	Complexity (for TV)
Chen and Ying (2024)	Uniformization	[A1], (15)	$\tilde{\mathcal{O}}(d)$
Zhang et al. (2024)	Euler-Method	[A1]	$\tilde{\mathcal{O}}(d^{4/3}\epsilon^{-4/3})$
Ren et al. (2025)	τ -leaping	[A1]	$\tilde{\mathcal{O}}(d\epsilon^{-1})$
This paper	Alg. 2	[A1]	$\tilde{\mathcal{O}}(d)$

Table 1: Comparison among discrete inference algorithms. Stopping time will be $T - \epsilon/d$ to guarantee the TV convergence with ϵ error tolerance.

Theorem 5 Suppose Assumption [A1] holds and we aim to draw samples from the discrete set $\{0, 1\}^{d \log_2 K}$ according to the distribution q_0 . If we train a discrete diffusion model satisfying

$$\epsilon_{score} \leq \frac{\epsilon}{\ln(d/\epsilon) + \ln \log_2 K} = \tilde{\mathcal{O}}(\epsilon),$$

and run Alg. 2, returning $\hat{\mathbf{y}}_{t_W}$, with

$$t_0 = 0, \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}), \quad t_W = T - \delta, \quad \text{and} \quad \beta_{t_w} = \frac{2d \log_2 K}{\min\{1, T - t_w\}},$$

where

$$T = \ln(d/\epsilon) + \ln \log_2 K \quad \text{and} \quad \delta \leq \frac{\epsilon}{d \log_2 K},$$

then the expected score estimation complexity of Alg. 2 is $\mathcal{O}(d \ln^2(d/\epsilon))$ to achieve $\text{TV}(q_0, \hat{q}_{t_W}) = o(\epsilon)$, where \hat{q}_{t_W} denotes the distribution of $\hat{\mathbf{y}}_{t_W}$.

In comparison to biased discrete inference, such as the Euler method (Zhang et al., 2024) and τ -leaping (Ren et al., 2025), which respectively require $\tilde{\mathcal{O}}(d^{4/3}\epsilon^{-4/3})$ and $\tilde{\mathcal{O}}(d\epsilon^{-1})$ complexity to ensure total variation convergence, truncated uniformization method only requires $\tilde{\mathcal{O}}(d)$ discrete score evaluations, significantly improving efficiency. Further distinguishing itself from standard uniformization methods, truncated uniformization removes the widely adopted assumption in Eq. (15), thereby significantly enhancing its practical applicability. Detailed comparison can be found in Table 1.

QTD-based Inference Results. To analyze the end-to-end convergence and gradient complexity required to achieve TV distance convergence, we make the following assumptions on p_* :

[A2] The second moment of p_* is bounded, i.e., $\mathbb{E}_{\mathbf{x} \sim p_*}[\|\mathbf{x}\|^2] \leq m_0$.

[A3] The energy function of p_* has bounded Hessian, i.e., $\|\nabla^2 \ln p_*\| \leq H$.

[A4] For any $\mathbf{u} \in \mathbb{R}^d$, there is a scalar sub-Gaussian tail, i.e.,

$$\mathbb{E}_{\mathbf{x} \sim p_*} \left[\exp \left(t \cdot \mathbf{x}^\top \mathbf{u} \right) \right] \leq \exp \left(\sigma^2 t^2 \|\mathbf{u}\|^2 / 2 \right).$$

Assumptions [A2] and [A3] encapsulate the minimal smoothness conditions in Chen et al. (2023a), and Assumption [A3] can be circumvented via early stopping (Benton et al., 2024a; Chen et al., 2023a; Li and Yan, 2024). Though our analysis also assumes light tails, Assumption [A4] (the σ -sub-Gaussian property) is invoked mainly to ensure clear convergence guarantees; exponential tails would suffice. We impose no isoperimetric constraints, and p_* need not be log-concave or unimodal. Meanwhile, Assumption [A1] aligns with standard practice for discrete score estimation (Zhang et al., 2024; Chen and Ying, 2024). Crucially, no smoothness or boundedness is required for the intermediate scores \tilde{v}_t , enforcing only minimal conditions on the score. Under these assumptions, we establish the following theorem, with the proof deferred to Appendix E.2.

Theorem 6 *Suppose Assumptions [A1]–[A4] hold. If we apply Alg. 1 with*

$$L = \sigma \cdot \sqrt{2 \ln(2d/\epsilon)}, \quad l = \left[2H \cdot \left(\sigma \sqrt{2d \ln(2d/\epsilon)} + d + \sqrt{dm_0} \right) \right]^{-1} \cdot \epsilon, \quad \text{and} \quad K = 2L/l$$

to quantize p_ , and run Alg. 2 with the same hyperparameter settings as in Theorem 5, then the expected score estimation complexity of Alg. 2 is $\mathcal{O}(d \ln^2(d/\epsilon))$ to achieve $\text{TV}(p_*, \hat{p}) \leq 5\epsilon$, where \hat{p} denotes the distribution of the generated samples.*

Unlike conventional diffusion models that directly apply the noising–denoising procedure in Euclidean space, QTD achieves a state-of-the-art linear convergence rate with respect to the error tolerance ϵ , only requiring the additional mild sub-Gaussian assumption [A4]. Besides, we also provide Theorem 7, i.e., the early-stopping version of Theorem 6, in the following

Theorem 7 *Suppose Assumption [A1], [A2] and [A4] hold. If we apply Alg. 1 with*

$$L \geq 2\sqrt{2d}\sigma \cdot \sqrt{\ln(2d\sigma/\epsilon)}, \quad l = \epsilon/\sqrt{d} \quad \text{and} \quad K = 2L/l,$$

to quantize p_ and run Alg. 2 with the same hyperparameter settings as in Theorem 5. Then the expectation of score estimation complexity of Alg. 2 will be $\mathcal{O}(d \ln^2(d/\epsilon))$ to achieve $\text{TV}(\bar{p}_*, \hat{p}) \leq 2\epsilon$ where \hat{p} denotes the underlying distribution of generated samples, and \bar{p}_* approximates the data distribution in 2-Wasserstein distance, i.e., $W_2(p_*, \bar{p}_*) \leq 2\epsilon$.*

The detailed description and proof will be deferred to Appendix E.4. As summarized in Table 2, our results simultaneously achieve fast convergence and operate under mild assumptions. Prior continuous diffusion samplers either attain only polynomial complexity in ϵ —ranging from $\tilde{\mathcal{O}}(d/\epsilon^2)$ to $\tilde{\mathcal{O}}(d^{5/4}/\epsilon^{1/2})$ —or, when targeting logarithmic dependence on $1/\epsilon$, impose strong structural conditions such as smooth score requirements on the entire forward OU process (Chen et al., 2026) or the restriction that the target distribution takes the form of a bounded distribution convolved with a Gaussian (Gatmiry et al., 2026). In contrast, our method achieves a score estimation complexity of $\mathcal{O}(d \ln^2(d/\epsilon))$ —the best among all listed approaches—while requiring neither the Euclidean score estimation error assumption [A1] nor the **smooth score** condition on the forward marginals. Instead, our assumptions pertain only to the data distribution p_* itself and can be further relaxed via early stopping (Theorem 7), making our guarantees broadly applicable.

5. Conclusion and Limitation

In conclusion, we introduce a novel approach, QTD, which first quantizes the continuous data distribution into a discrete counterpart, and then applies a truncated uniformization procedure to

Table 2: Comparison with prior works simulating reverse particle SDEs, where **[A1]'** denotes the score estimation error trained in Euclidean space and **smooth score** denotes the smooth score assumption for the whole OU process (p_t) starting from p_* . Note that Assumptions **[A3]** is only about p_* and can be replaced by the early stopping trick. All complexities for TV convergence are achieved by assuming $\epsilon_{\text{score}} = \tilde{o}(\epsilon)$.

Results	Early-Stopping	Assumptions	Complexity (for TV)
Chen et al. (2023b)	No	[A1]' , [A2] , smooth score	$\tilde{\mathcal{O}}(d\epsilon^{-2})$
Chen et al. (2023a)	No	[A1]' , [A2] , [A3]	$\tilde{\mathcal{O}}(d^2\epsilon^{-2})$
Huang et al. (2024b)	No	[A1]' , [A2] , smooth score	$\tilde{\mathcal{O}}(d^{1/2}\epsilon^{-1})$
Gatmiry et al. (2026)	No	$p_* = p_0(\text{R-bounded}) * \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$	$\mathcal{O}((R/\sigma)^2 \ln^2(1/\epsilon))$
Chen et al. (2026)	No	[A1]' , [A2] , smooth score	$\mathcal{O}(d \ln^3(d/\epsilon))$
Our Theorem 6	No	[A1] , [A2] , [A3] , [A4]	$\mathcal{O}(d \ln^2(d/\epsilon))$
Chen et al. (2023a)	Yes	[A1]' , [A2]	$\tilde{\mathcal{O}}(d^2\epsilon^{-2})$
Benton et al. (2024a)	Yes	[A1]' , [A2]	$\tilde{\mathcal{O}}(d\epsilon^{-2})$
Li and Yan (2024)	Yes	[A1]' , [A2]	$\tilde{\mathcal{O}}(d\epsilon^{-1})$
Li and Cai (2024)	Yes	[A1]' , [A2]	$\tilde{\mathcal{O}}(d^{5/4}\epsilon^{-1/2})$
Our Theorem 7	Yes	[A1] , [A2] , [A4]	$\mathcal{O}(d \ln^2(d/\epsilon))$

achieve unbiased inference with improved score-evaluation complexity for continuous data generation. Beyond its state-of-the-art theoretical complexity, the truncated uniformization framework is of independent interest as an inference algorithm for discrete diffusion models, where it also attains top-tier theoretical complexity under minimal assumptions.

A key limitation of our approach is that achieving accelerated convergence without degrading generation quality requires the discrete score estimation error to be on par with the continuous score estimation error outlined by (Chen et al., 2023b; Benton et al., 2024a). While some works (Meng et al., 2022; Lou et al., 2024) have introduced discrete training objectives such as concrete score matching and denoising score entropy, no direct comparison between discrete and continuous score training has been conducted.

Acknowledgments

The research is supported by NSF Award CCF-2112665 (TILOS) and CDC-RFA-FT-23-0069 from the CDC’s Center for Forecasting and Outbreak Analytics. Difan Zou also acknowledges the support from Hong Kong GRF award 17307425 and ECS award 27309624.

References

- Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly d -linear convergence bounds for diffusion models via stochastic localization. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Joe Benton, Yuyang Shi, Valentin De Bortoli, George Deligiannidis, and Arnaud Doucet. From denoising diffusions to denoising markov models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 86(2):286–301, 2024b.
- Stéphane Boucheron, Gábor Lugosi, and Olivier Bousquet. Concentration inequalities. In *Summer school on machine learning*, pages 208–240. Springer, 2003.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- Fan Chen, Sinho Chewi, Constantinos Daskalakis, and Alexander Rakhlin. High-accuracy sampling for diffusion models and log-concave distributions. *arXiv preprint arXiv:2602.01338*, 2026.
- Hongrui Chen and Lexing Ying. Convergence analysis of discrete diffusion model: Exact implementation through uniformization. *arXiv preprint arXiv:2402.08095*, 2024.
- Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pages 4735–4763. PMLR, 2023a.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Khashayar Gatmiry, Sitan Chen, and Adil Salim. High-accuracy and dimension-free sampling with diffusions. *arXiv preprint arXiv:2601.10708*, 2026.
- Ye He, Kevin Rojas, and Molei Tao. Zeroth-order sampling methods for non-log-concave distributions: Alleviating metastability by denoising diffusion. *Advances in Neural Information Processing Systems*, 37:71122–71161, 2024.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Xunpeng Huang, Hanze Dong, Yifan Hao, Yi An Ma, and Tong Zhang. Reverse diffusion monte carlo. In *12th International Conference on Learning Representations, ICLR 2024*, 2024a.
- Xunpeng Huang, Difan Zou, Hanze Dong, Zhang Zhang, Yian Ma, and Tong Zhang. Reverse transition kernel: A flexible framework to accelerate diffusion inference. *Advances in Neural Information Processing Systems*, 37:95515–95578, 2024b.
- Gen Li and Changxiao Cai. Provable acceleration for diffusion models under minimal assumptions. *arXiv preprint arXiv:2410.23285*, 2024.

- Gen Li and Yuchen Jiao. Improved convergence rate for diffusion probabilistic models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Gen Li and Yuling Yan. $O(d/t)$ convergence theory for diffusion probabilistic models under minimal assumptions. *arXiv preprint arXiv:2409.18959*, 2024.
- Yuchen Liang, Renxiang Huang, Lifeng Lai, Ness Shroff, and Yingbin Liang. Absorb and converge: Provable convergence guarantee for absorbing discrete diffusion models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=VsZzTSyk5p>.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Proceedings of the 41st International Conference on Machine Learning*, pages 32819–32848, 2024.
- Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35: 34532–34545, 2022.
- Yinuo Ren, Haoxuan Chen, Yuchen Zhu, Wei Guo, Yongxin Chen, Grant M Rotskoff, Molei Tao, and Lexing Ying. Fast solvers for discrete diffusion models: Theory and applications of high-order algorithms. *arXiv preprint arXiv:2502.00234*, 2025.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr, 2015.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Nico M van Dijk. Approximate uniformization for continuous-time markov chains with an application to performability analysis. *Stochastic processes and their applications*, 40(2):339–357, 1992.
- Nico M van Dijk, Sem PJ van Brummelen, and Richard J Boucherie. Uniformization: Basics, extensions and applications. *Performance evaluation*, 118:8–32, 2018.
- Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- Zikun Zhang, Zixiang Chen, and Quanquan Gu. Convergence of score-based discrete diffusion models: A discrete-time analysis. *arXiv preprint arXiv:2410.02321*, 2024.

Contents

1	Introduction	2
2	Preliminaries	3
3	Quantized Transition Diffusion	5
3.1	Histogram Approximation	6
3.2	The Discrete Forward Process and Binary Encoding	7
3.3	Truncated Uniformization	8
4	Theoretical Results	10
5	Conclusion and Limitation	13
A	Notation Summary	18
B	Forward and Reverse Processes of Discrete Diffusion Models	19
B.1	The Forward Process of Discrete Diffusion Models	19
B.2	Proof of Eq. (3)	22
B.3	Proof of Eq. (4)	24
C	Proof of Lemma 1	25
D	Proof of Lemma 2	28
E	Supplementary Proofs for the Discrete Reverse Process	31
E.1	Proof of Lemma 3	31
E.2	Proof of Theorem 6	32
E.3	The Proof of Lemma 4	36
E.4	Proof of Theorem 7	40
F	Technical Lemmas	43
G	Empirical Results	44
G.1	Experiments on Synthetic Data	44
G.2	Experiments on MNIST	46

Appendix A. Notation Summary

We summarize all notations used in the main paper and appendix in Table 3.

Table 3: Summary of key notations used in the paper.

Symbol	Description
Cube (L)	Bounded cube $[-L, L]^d$ covering high-probability mass of p_*
Cell (i_0, \dots, i_{d-1})	Quantization cell (hypercubes) defined by coordinate bins, Eq. (10)
\mathcal{Y}	Binary discrete space $\{0, 1\}^{d \log_2 K}$
$\bar{\mathcal{Y}}$	Grid index space $\{0, \dots, K-1\}^d$
$\text{vBin}(\cdot)$	Mapping from grid index $\bar{\mathcal{Y}}$ to binary code \mathcal{Y}
$p_* \propto \exp(-f_*)$	Target continuous distribution in \mathbb{R}^d
\tilde{p}_*	Truncated and renormalized version of p_* over Cube (L), Eq. (9)
\bar{p}_*	Histogram approximation to \tilde{p}_* over Cube (L), Eq. (11)
\bar{q}_*	Discrete distribution on $\bar{\mathcal{Y}} = \{0, \dots, K-1\}^d$ induced by \bar{p}_* , Eq. (12)
q_*	Discrete distribution on $\mathcal{Y} = \{0, 1\}^{d \log_2 K}$, $q_* = \bar{q}_* \circ \text{vBin}^{-1}$
\mathbf{y}_t^\rightarrow	Forward-time CTMC on \mathcal{Y}
q_t^\rightarrow	Marginal distribution of forward process at time t , i.e., $\mathbf{y}_t^\rightarrow \sim q_t^\rightarrow$
$q_{t',t}^\rightarrow$	Joint distribution of $(\mathbf{y}_{t'}^\rightarrow, \mathbf{y}_t^\rightarrow)$
q_∞^\rightarrow	Stationary distribution of the forward CTMC (uniform distribution)
$q_{t' t}^\rightarrow(\mathbf{y}' \mathbf{y})$	Conditional transition probability in forward process, Eq. (1)
\mathbf{y}_t^\leftarrow	Reverse-time CTMC defined by $q_t^\leftarrow := q_{T-t}^\rightarrow, \mathbf{y}_t^\leftarrow \sim q_t^\leftarrow$
q_t^\leftarrow	Marginal distribution of reverse process at time t , $q_t^\leftarrow = q_{T-t}^\rightarrow$
$q_{t',t}^\leftarrow$	Joint distribution of $(\mathbf{y}_{t'}^\leftarrow, \mathbf{y}_t^\leftarrow)$
$q_{t' t}^\leftarrow(\mathbf{y}' \mathbf{y})$	Conditional transition probability of the ideal reverse process
$\hat{q}_{t+\Delta t t}(\mathbf{y}' \mathbf{y})$	Practical reverse conditional probability, Eq. (8)
$R^\rightarrow(\mathbf{y}, \mathbf{y}')$	Forward transition rate from state \mathbf{y}' to \mathbf{y} , Eq. (2), and Eq. (13). This follows the ordering of the conditional distribution $p(\mathbf{y} \mathbf{y}')$, which is the <i>transpose</i> of the convention used in some other works.
$R_t^\leftarrow(\mathbf{y}, \mathbf{y}')$	Reverse transition rate at time t from state \mathbf{y}' to \mathbf{y} , $R_t^\leftarrow(\mathbf{y}, \mathbf{y}') := R^\rightarrow(\mathbf{y}', \mathbf{y}) \cdot \frac{q_t^\leftarrow(\mathbf{y})}{q_t^\leftarrow(\mathbf{y}')}$, Eq. (3)
$\tilde{R}_t(\mathbf{y}, \mathbf{y}')$	Estimated reverse transition rate using the learned density ratio, $\tilde{R}_t(\mathbf{y}, \mathbf{y}') = R^\rightarrow(\mathbf{y}', \mathbf{y}) \cdot \tilde{v}_{t, \mathbf{y}'}(\mathbf{y})$, Eq. (5)
$\hat{R}_t(\cdot, \cdot)$	Truncated version of $\tilde{R}_t(\cdot, \cdot)$ with threshold β_t , Eq. (16)
$R_t^\leftarrow(\mathbf{y}), \tilde{R}_t(\mathbf{y}), \hat{R}_t(\mathbf{y})$	Total reverse transition rate out of state \mathbf{y} for each rate type, defined as $R(\mathbf{y}) := \sum_{\mathbf{y}' \neq \mathbf{y}} R(\mathbf{y}', \mathbf{y})$ with $R \in \{R_t^\leftarrow, \tilde{R}_t, \hat{R}_t\}$
β_t	Upper bound on $R_t^\leftarrow(\mathbf{y})$, $\beta_t = 2d \log_2 K \max\{1, (T-t)^{-1}\}$, Eq. (14)
$v_{t, \mathbf{y}'}(\mathbf{y})$	Density ratio $q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$
$\tilde{v}_{t, \mathbf{y}'}(\mathbf{y})$	Learned approximation to $v_{t, \mathbf{y}'}(\mathbf{y}) = q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$
$L_{\text{SE}}(\hat{v})$	Score entropy loss used to train \tilde{v} , Eq. (5)
\mathbf{e}_i	One-hot vector with a 1 at position i and 0 elsewhere
l	Width of each quantization cell
$K = 2L/l$	Number of quantization bins per dimension

Appendix B. Forward and Reverse Processes of Discrete Diffusion Models

In order to simplify the notation in this section, we introduce some new notations as supplementary to Section 2. Since we consider the discrete diffusion on \mathcal{Y} , we defined the inner product on this discrete space for two functions as

$$\langle f, g \rangle_{\mathcal{Y}} := \sum_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{y}) \cdot g(\mathbf{y}).$$

Besides, the delta on \mathcal{Y} is defined as

$$\delta_{\mathbf{y}}(\mathbf{y}') = \begin{cases} 1 & \mathbf{y}' = \mathbf{y} \\ 0 & \text{otherwise} \end{cases}.$$

B.1. The Forward Process of Discrete Diffusion Models

In this section, we refine the introduction about the forward process of discrete diffusion in Section 2 with the same notations. In general, the time-homogeneous CTMC can be described by a Markov semigroup $\mathcal{Q}_t^{\rightarrow}$ defined as:

$$\mathcal{Q}_t^{\rightarrow}[f](\mathbf{y}) = \mathbb{E}[f(\mathbf{y}_t) | \mathbf{y}_0 = \mathbf{y}] = \left\langle f, q_{t|0}^{\rightarrow}(\cdot | \mathbf{y}) \right\rangle_{\mathcal{Y}} \quad (18)$$

where the function $f: \mathcal{Y} \rightarrow \mathbb{R}$. Due to the definition, the infinitesimal operator $\mathcal{L}^{\rightarrow}$ of the time homogeneous $\mathcal{Q}_t^{\rightarrow}$ is denoted as

$$\mathcal{L}^{\rightarrow}[f](\mathbf{y}) = \lim_{t \rightarrow 0} \left[\frac{\mathcal{Q}_t^{\rightarrow}[f] - f}{t} \right](\mathbf{y}) = \left\langle f, \partial_t q_{t|0}^{\rightarrow}(\cdot | \mathbf{y}) \Big|_{t=0} \right\rangle_{\mathcal{Y}} := \langle f, R^{\rightarrow}(\cdot, \mathbf{y}) \rangle_{\mathcal{Y}} \quad (19)$$

where

$$R^{\rightarrow}(\mathbf{y}', \mathbf{y}) := \partial_t q_{t|0}^{\rightarrow}(\mathbf{y}' | \mathbf{y}) \Big|_{t=0} = \lim_{t \rightarrow 0} \left[\frac{q_{t|0}^{\rightarrow}(\mathbf{y}' | \mathbf{y}) - \delta_{\mathbf{y}}(\mathbf{y}')}{t} \right]. \quad (20)$$

According to the time-homogeneous property, we have

$$q_{t+\Delta t|t}^{\rightarrow}(\mathbf{y}' | \mathbf{y}) = \delta_{\mathbf{y}}(\mathbf{y}') + \Delta t \cdot R^{\rightarrow}(\mathbf{y}', \mathbf{y}) + o(\Delta t)$$

for any t . Here, the transition rate function R^{\rightarrow} must satisfy

$$R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \geq 0 \text{ when } \mathbf{y}' \neq \mathbf{y} \quad \text{and} \quad R^{\rightarrow}(\mathbf{y}', \mathbf{y}') = - \sum_{\mathbf{y} \neq \mathbf{y}'} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \leq 0 \quad (21)$$

due to the definition Eq. (20). Under this setting, we can provide the dynamic of $q_{t|0}$ for any t . Specifically, we have

$$\begin{aligned} \partial_t \mathcal{Q}_t^{\rightarrow}[f](\mathbf{y}) &= \mathcal{Q}_t^{\rightarrow}[\mathcal{L}f](\mathbf{y}) = \left\langle \mathcal{L}f, q_{t|0}^{\rightarrow}(\cdot | \mathbf{y}) \right\rangle_{\mathcal{Y}} = \sum_{\mathbf{y}' \in \mathcal{Y}} \mathcal{L}^{\rightarrow}[f](\mathbf{y}') \cdot q_{t|0}^{\rightarrow}(\mathbf{y}' | \mathbf{y}) \\ &= \sum_{\mathbf{y}' \in \mathcal{Y}} \left[\sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} f(\tilde{\mathbf{y}}) \cdot R^{\rightarrow}(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{t|0}(\mathbf{y}' | \mathbf{y}) \right] = \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} \left[f(\tilde{\mathbf{y}}) \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} R^{\rightarrow}(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{t|0}(\mathbf{y}' | \mathbf{y}) \right], \end{aligned}$$

where the first inequality follows from the semigroup property. Combining with the fact

$$\partial_t \mathcal{Q}_t^\rightarrow[f](\mathbf{y}) = \left\langle f, \partial_t q_{t|0}^\rightarrow(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}}$$

derived from Eq. (18), we have

$$\partial_t q_{t|0}^\rightarrow(\tilde{\mathbf{y}}|\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{t|0}^\rightarrow(\mathbf{y}'|\mathbf{y}) = \left\langle R(\tilde{\mathbf{y}}, \cdot), q_{t|0}^\rightarrow(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}}.$$

According to the time-homogeneous property, the above equation can be easily extended to

$$\partial_t q_{t|s}^\rightarrow(\tilde{\mathbf{y}}|\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{t|s}^\rightarrow(\mathbf{y}'|\mathbf{y}) = \left\langle R(\tilde{\mathbf{y}}, \cdot), q_{t|s}^\rightarrow(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}}. \quad (22)$$

Combining with Bayes' Theorem, the transition of the marginal distribution is

$$\frac{dq_t^\rightarrow}{dt}(\mathbf{y}) = \langle R(\mathbf{y}, \cdot), q_t^\rightarrow \rangle_{\mathcal{Y}}. \quad (23)$$

Matrix Presentation. Suppose the support set \mathcal{Y} of q_t^\rightarrow be written as $\mathcal{Y} = \{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{|\mathcal{X}|}\}$, we may consider the marginal distribution q_s^\rightarrow to be a vector, i.e.,

$$\mathbf{q}_t^\rightarrow = [q_t(\mathbf{y}_0), q_t(\mathbf{y}_1), \dots, q_t(\mathbf{y}_{|\mathcal{Y}|-1})],$$

conditional transition probability function $q_{t|s}^\rightarrow$ to be a matrix, i.e.,

$$\mathbf{Q}_{t|s}^\rightarrow = \begin{bmatrix} q_{t|s}^\rightarrow(\mathbf{y}_0|\mathbf{y}_0) & q_{t|s}^\rightarrow(\mathbf{y}_0|\mathbf{y}_1) & \dots & q_{t|s}^\rightarrow(\mathbf{y}_0|\mathbf{y}_{|\mathcal{Y}|-1}) \\ q_{t|s}^\rightarrow(\mathbf{y}_1|\mathbf{y}_0) & q_{t|s}^\rightarrow(\mathbf{y}_1|\mathbf{y}_1) & \dots & q_{t|s}^\rightarrow(\mathbf{y}_1|\mathbf{y}_{|\mathcal{Y}|-1}) \\ \dots & \dots & \dots & \dots \\ q_{t|s}^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}|\mathbf{y}_0) & q_{t|s}^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}|\mathbf{y}_1) & \dots & q_{t|s}^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}|\mathbf{y}_{|\mathcal{Y}|-1}) \end{bmatrix}.$$

Similarly, the function R can also be presented as

$$\mathbf{R}^\rightarrow = \begin{bmatrix} R^\rightarrow(\mathbf{y}_0, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_0, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_0, \mathbf{y}_{|\mathcal{Y}|-1}) \\ R^\rightarrow(\mathbf{y}_1, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_1, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_1, \mathbf{y}_{|\mathcal{Y}|-1}) \\ \dots & \dots & \dots & \dots \\ R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_{|\mathcal{Y}|-1}) \end{bmatrix}. \quad (24)$$

Under this condition, Eq. (23) can be written as

$$d\mathbf{q}_t^\rightarrow/dt = \mathbf{R}^\rightarrow \cdot \mathbf{q}_t^\rightarrow \quad (25)$$

matching the usual presentation shown in [Chen and Ying \(2024\)](#); [Zhang et al. \(2024\)](#). Besides, Eq. (21) shown in Section 2 can also be presented as $\mathbf{1} \cdot \mathbf{R} = \mathbf{0}$.

The following lemma gives the closed-form expression for the probability transition kernel of the forward process, which also suggests an efficient implementation.

Lemma 8 (Forward transition kernel) Consider the forward CTMC, i.e., $\{\mathbf{y}_t\}_{t=0}^T$ with the infinitesimal operator \mathbb{R}^\rightarrow given in Eq. (13). Then for any two timestamps $s \leq t$, the forward transition probability satisfies

$$q_{t|s}^\rightarrow(\mathbf{y}|\mathbf{y}') = 2^{-d \log_2 K} \cdot \prod_{i=0}^{d \log_2 K - 1} \left[1 + (-1)^{|\mathbf{y}_i - \mathbf{y}'_i|} \cdot e^{-2(t-s)} \right].$$

Remark 9 The transition probability in Lemma 8 factorizes across coordinates. This means that the forward transition can be implemented as $d \log_2 K$ independent bit-wise updates. Specifically, for each coordinate i , flip \mathbf{y}'_i with probability $\frac{1 - e^{-2(t-s)}}{2}$ to obtain \mathbf{y}_i .

Proof Combining Eq. (24) and Eq. (25), the dynamic of marginal distribution q_t^\rightarrow can be written as a matrix-vector product, i.e.,

$$dq_t^\rightarrow/dt = \mathbf{R}^\rightarrow \cdot q_t^\rightarrow$$

where

$$\mathbf{R}^\rightarrow = \begin{bmatrix} R^\rightarrow(\mathbf{y}_0, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_0, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_0, \mathbf{y}_{|\mathcal{Y}|-1}) \\ R^\rightarrow(\mathbf{y}_1, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_1, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_1, \mathbf{y}_{|\mathcal{Y}|-1}) \\ \dots & \dots & \dots & \dots \\ R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_0) & R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_1) & \dots & R^\rightarrow(\mathbf{y}_{|\mathcal{Y}|-1}, \mathbf{y}_{|\mathcal{Y}|-1}) \end{bmatrix}.$$

Here, \mathbf{R}^\rightarrow can be decomposed into the sum

$$\mathbf{R}^\rightarrow = \sum_{i=0}^{d \log_2 K - 1} \mathbf{R}_i^\rightarrow,$$

we first note that the state space is $\{0, 1\}^{d \log_2 K}$, where each coordinate can flip independently. Hence, each coordinate contributes its own ‘‘flip’’ component to the overall generator \mathbf{R}^\rightarrow . Concretely, let us label the coordinates $0, \dots, d \log_2 K - 1$, and consider the generator corresponding to a single coordinate i . Such a generator only acts nontrivially on the i th coordinate, which can flip from 0 to 1 or 1 to 0, while all other coordinates remain unchanged.

Each ‘‘flip’’ for coordinate i can be represented by a 2×2 generator matrix (reflecting the two possible states, 0 or 1). Moreover, since the flipping of different coordinates occurs independently, we adopt the tensor-product (or Kronecker-product) structure, placing the 2×2 flip matrix in the i th position and 2×2 identity matrices in all other positions. Hence, each \mathbf{R}_i^\rightarrow is of the form

$$\mathbf{R}_i^\rightarrow = \mathbf{I} \otimes \dots \otimes \mathbf{A} \otimes \dots \otimes \mathbf{I},$$

where

$$\mathbf{A} := \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}$$

is a generator of the flip in the i th coordinate, and \mathbf{I} is the 2×2 identity in all coordinates. By the Kolmogorov forward equation, we have

$$\mathbf{Q}_{t|s}^\rightarrow = \exp((t-s)\mathbf{R}^\rightarrow) = \exp((t-s)\mathbf{A})^{\otimes d} = \begin{bmatrix} \frac{1 + e^{-2(t-s)}}{2} & \frac{1 - e^{-2(t-s)}}{2} \\ \frac{1 - e^{-2(t-s)}}{2} & \frac{1 + e^{-2(t-s)}}{2} \end{bmatrix}^{\otimes d},$$

which implies

$$q_{t|s}^{\rightarrow}(\mathbf{y}|\mathbf{y}') = 2^{-d \log_2 K} \cdot \prod_{i=0}^{d \log_2 K - 1} \left[1 + (-1)^{|\mathbf{y}_i - \mathbf{y}'_i|} \cdot e^{-2(t-s)} \right] \quad \text{and} \quad \mathbf{y}, \mathbf{y}' \in \mathcal{Y}.$$

Hence, the proof is completed. \blacksquare

Figure 3 visualizes the evolution of transition probabilities under different forward processes. The tridiagonal CTMC (second row) can be viewed as a discrete analogue of the normalized Gaussian transition (first row), where the domain $[0, 1]$ is quantized into 8 bins. The tridiagonal structure results in slow mixing, as transitions are restricted to immediate neighbors. At small time steps (e.g., $t = 0.01$, first column), the transition kernel satisfies $\mathbf{Q}_{t+\Delta t|t}^{\rightarrow} \approx \mathbf{I} + \Delta t \cdot \mathbf{R}^{\rightarrow}$, so the sparsity of the transition kernel closely reflects that of the rate matrix \mathbf{R}^{\rightarrow} . For efficient simulation of the reverse process, defined by $R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') := R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_t^{\leftarrow}(\mathbf{y})}{q_t^{\leftarrow}(\mathbf{y}')}$ as Eq. (3), it is essential that \mathbf{R}^{\rightarrow} remains sparse. While the dense forward process (third row) mixes rapidly, it incurs high computational cost per step when simulating the reverse process. In contrast, the hypercube structure (fourth row) achieves a favorable balance: it enables efficient long-range transitions for fast mixing while preserving an $\mathcal{O}(\log |\mathcal{Y}|)$ sparse structure for efficient implementation.

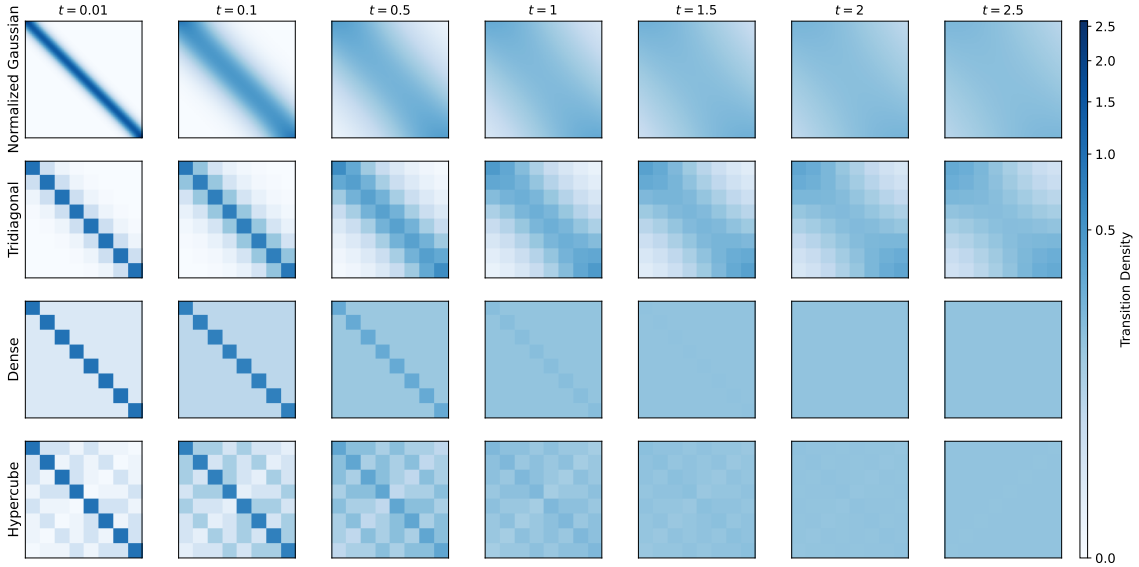


Figure 3: Heatmaps of the probability transition at different time steps t for four diffusion processes: a continuous normalized Gaussian kernel on $[0, 1]$ (top row), and discrete CTMCs over $|\mathcal{Y}| = 8$ states based on tridiagonal, dense, and hypercube transition rate matrices (bottom three rows).

B.2. Proof of Eq. (3)

Proof For any $t \in [0, T]$, the marginal, joint, and conditional distribution w.r.t. $\{\mathbf{y}_t^{\leftarrow}\}$ are denoted as

$$\mathbf{y}_t^{\leftarrow} \sim q_t^{\leftarrow}, \quad (\mathbf{y}_t^{\leftarrow}, \mathbf{y}_{t'}^{\leftarrow}) \sim q_{t,t'}^{\leftarrow}, \quad \text{and} \quad q_{t'|t}^{\leftarrow} = q_{t',t}^{\leftarrow} / q_t^{\leftarrow},$$

which have $q_t^{\leftarrow} = q_{T-t}^{\rightarrow}$. Considering the reverse CTMC, we have

$$\begin{aligned} \frac{dq_t^{\leftarrow}}{dt}(\mathbf{y}) &= \lim_{\Delta t \rightarrow 0} \frac{q_{t+\Delta t}^{\leftarrow}(\mathbf{y}) - q_t^{\leftarrow}(\mathbf{y})}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}'} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') \cdot q_t^{\leftarrow}(\mathbf{y}') - q_t^{\leftarrow}(\mathbf{y})}{\Delta t} \\ &= \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}'} \left(q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'(\mathbf{y})} \right) \cdot q_t^{\leftarrow}(\mathbf{y}')}{\Delta t}. \end{aligned}$$

Therefore, due to the following definition

$$R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') := \lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'(\mathbf{y})}}{\Delta t} \right],$$

we have

$$\frac{dq_t^{\leftarrow}}{dt}(\mathbf{y}) = \langle R_t^{\leftarrow}(\mathbf{y}, \cdot), q_t^{\leftarrow}(\cdot) \rangle_{\mathcal{Y}}.$$

Then, we start to check the property of $R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}')$ when $\mathbf{y} \neq \mathbf{y}'$ in the following, which has

$$R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') = \lim_{\Delta t \rightarrow 0} \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}')}{\Delta t} = \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}').$$

Specifically, with the Bayes Theorem, it has

$$\begin{aligned} \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}') &= -1 \cdot \lim_{s \rightarrow t} \partial_{T-t} q_{T-t|T-s}^{\rightarrow}(\mathbf{y}|\mathbf{y}') \\ &= -1 \cdot \lim_{s \rightarrow t} \partial_{T-t} \left[\frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) \cdot q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} \right] \\ &= \underbrace{-\lim_{s \rightarrow t} \partial_{T-t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')}}_{\text{Term 1}} - \underbrace{\lim_{s \rightarrow t} \frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} \cdot \partial_{T-t} q_{T-t}^{\rightarrow}(\mathbf{y})}_{\text{Term 2}}. \end{aligned} \tag{26}$$

For Term 1 of Eq. (26), we have

$$\text{Term 1} = -R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')},$$

which follows from the Kolmogorov backward theorem (Lemma 19) and Eq. (19):

$$\partial_{T-t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) = -\mathcal{L}^{\rightarrow}[q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\cdot)](\mathbf{y}) = -\left\langle q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\cdot), R^{\rightarrow}(\cdot, \mathbf{y}) \right\rangle_{\mathcal{Y}} = R^{\rightarrow}(\mathbf{y}', \mathbf{y}).$$

Here, the last equation follows from the fact

$$\lim_{s \rightarrow t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\tilde{\mathbf{y}}) = 1 \quad \text{only when} \quad \mathbf{y}' = \tilde{\mathbf{y}}.$$

Additionally, with the fact $\mathbf{y}' \neq \mathbf{y}$, it has

$$\lim_{s \rightarrow t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}') = 0.$$

That means Term 2 of Eq. (26) will definitely be 0. Hence, the proof is completed. \blacksquare

B.3. Proof of Eq. (4)

Proof [Adapted from Proposition 1 of [Campbell et al. \(2022\)](#)] The RHS of Eq. (4) satisfies

$$\lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'(\mathbf{y})}}{\Delta t} \right] = \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}').$$

Besides, we have

$$\begin{aligned} \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}') &= \lim_{s \rightarrow t} \partial_t \left[q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} \right] \\ &= \lim_{s \rightarrow t} \left[\partial_t (q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y})) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} + q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) \cdot \frac{\partial_t q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} \right]. \end{aligned}$$

When $\mathbf{y} \neq \mathbf{y}'$, we have

$$\lim_{s \rightarrow t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) = 0,$$

which implies

$$\lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}') = \lim_{s \rightarrow t} \partial_t (q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y})) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-s}^{\rightarrow}(\mathbf{y}')} = R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')}.$$

The last equation follows from the Kolmogorov backward theorem, i.e., Lemma 19 and Eq. (19)

$$\partial_{T-t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) = -\mathcal{L}^{\rightarrow}[q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\cdot)](\mathbf{y}) = -\left\langle q_{T-s|T-t}^{\rightarrow}(\mathbf{y}'|\cdot), R^{\rightarrow}(\cdot, \mathbf{y}) \right\rangle_{\mathbf{y}} = R^{\rightarrow}(\mathbf{y}', \mathbf{y}).$$

Combining with Eq. (3), we have

$$\lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'(\mathbf{y})}}{\Delta t} \right] = \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}') = R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')} = R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') \quad (27)$$

when $\mathbf{y}' \neq \mathbf{y}$. Besides, we have

$$\begin{aligned} \sum_{\mathbf{y} \in \mathcal{Y}} R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') &= \sum_{\mathbf{y} \in \mathcal{Y}} R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')} \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} \lim_{\Delta t \rightarrow 0} \left[\frac{q_{T-t+\Delta t|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) - \delta_{\mathbf{y}'(\mathbf{y})}}{\Delta t} \right] \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')} = \sum_{\mathbf{y} \in \mathcal{Y}} \lim_{\Delta t \rightarrow 0} \left[\frac{q_{T-t+\Delta t|T-t}^{\rightarrow}(\mathbf{y}'|\mathbf{y}') - \delta_{\mathbf{y}'(\mathbf{y}')}}{\Delta t} \right] = 0, \end{aligned}$$

which means

$$R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}') = - \sum_{\mathbf{y} \neq \mathbf{y}'} R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') = \lim_{\Delta t \rightarrow 0} - \left[\frac{1 - \sum_{\mathbf{y} \neq \mathbf{y}'} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}')}{\Delta t} \right],$$

where the last inequality follows from Eq. (27). Hence, the proof is completed. \blacksquare

Appendix C. Proof of Lemma 1

Lemma 10 Suppose the data distribution p_* is σ sub-Gaussian, by choosing $L \geq \sigma \cdot \sqrt{2 \ln(2d/\epsilon)}$, the TV distance between p_* and \tilde{p}_* defined in Eq. (9) will be smaller than ϵ , i.e., $\text{TV}(p_*, \tilde{p}_*) \leq \epsilon$.

Proof When p_* satisfies σ sub-Gaussian properties, i.e.,

$$\mathbb{E}_{\mathbf{x} \sim p_*} [\exp(l \langle \mathbf{x}, \mathbf{u} \rangle)] \leq \exp\left(\frac{\sigma^2 l^2 \cdot \|\mathbf{u}\|^2}{2}\right).$$

By choosing $\mathbf{u} = \mathbf{e}_i$, we can easily found that each dimension of \mathbf{x} will be σ sub-Gaussian, i.e.,

$$\mathbb{E}_{\mathbf{x}_i \sim p_{*,i}} [\exp(l \mathbf{x}_i \mathbf{u}_i)] \leq \exp\left(\frac{\sigma^2 l^2 \cdot \|\mathbf{u}_i\|^2}{2}\right).$$

According to the sub-Gaussian properties for each coordinate, we have

$$\mathbb{P}_i [|\mathbf{x}_i| \geq l] \leq 2 \exp\left(-\frac{l^2}{2\sigma^2}\right).$$

With the union bound, we have

$$\mathbb{P}_* \left[\max_{1 \leq i \leq d} |\mathbf{x}_i| > L \right] \leq \sum_{i=1}^d \mathbb{P}_* [|\mathbf{x}_i| > L] \leq 2d \cdot \exp\left(-\frac{L^2}{2\sigma^2}\right).$$

Under this condition, by supposing

$$2d \cdot \exp\left(-\frac{L^2}{2\sigma^2}\right) \leq \epsilon \quad \Leftrightarrow \quad L \geq \sigma \cdot \sqrt{2 \ln \frac{2d}{\epsilon}}, \quad (28)$$

we have $\mathbb{P}_* [\max_{1 \leq i \leq d} |\mathbf{x}_i| \geq L] \leq \epsilon$. Under this condition, the total variation distance between \tilde{p}_* and p_* can be upper bounded by

$$\begin{aligned} \text{TV}(p_*, \tilde{p}_*) &= \frac{1}{2} \int_{\mathbb{R}^d} |p_*(\mathbf{x}) - \tilde{p}_*(\mathbf{x})| d\mathbf{x} \\ &= \frac{1}{2} \int_{\mathbf{x} \in \text{Cube}(L)} (\tilde{p}_*(\mathbf{x}) - p_*(\mathbf{x})) d\mathbf{x} + \frac{1}{2} \int_{\mathbf{x} \notin \text{Cube}(L)} p_*(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{2} \left[1 - \int_{\mathbf{x} \in \text{Cube}(L)} p_*(\mathbf{x}) d\mathbf{x} \right] + \frac{1}{2} \int_{\mathbf{x} \notin \text{Cube}(L)} p_*(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbf{x} \notin \text{Cube}(L)} p_*(\mathbf{x}) d\mathbf{x} \leq \epsilon \end{aligned} \quad (29)$$

where the last inequality follows from Eq. (28). Hence, the proof is completed. \blacksquare

Lemma 11 Suppose the distribution \tilde{p}_* defined in Eq. (9) satisfies H -smoothness, by choosing

$$l \leq (2HL + \|\nabla f_*(\mathbf{0})\|)^{-1} \cdot d^{-1/2} \epsilon,$$

the TV distance satisfies $\text{TV}(\tilde{p}_*, \bar{p}_*) \leq 2\epsilon$ where \bar{p}_* is defined in Eq. (11).

Proof By Lagrange's mean value theorem, for each cell $\text{Cell}(i_0, i_1, \dots, i_{d-1})$, there exists a point $\bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})$ such that

$$\tilde{p}_*(\bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}}) = \frac{\int_{\mathbf{u} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} \tilde{p}_*(\mathbf{u}) d\mathbf{u}}{l^d}.$$

Therefore, the piecewise constant density \bar{p}_* satisfies $\bar{p}_*(\mathbf{x}) = \tilde{p}_*(\bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}})$, for any $\mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})$.

We now aim to bound the difference $|\tilde{p}_*(\mathbf{u}) - \tilde{p}_*(\mathbf{x})|$ for any $\mathbf{u}, \mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})$, using H -smoothness. Later, we will choose $\mathbf{u} = \bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}}$ to bound the total variation distance between \tilde{p}_* and \bar{p}_* .

According to the construction of \tilde{p}_* , i.e., Eq. (9), we have

$$\frac{\tilde{p}_*(\mathbf{u})}{\tilde{p}_*(\mathbf{x})} = \frac{p_*(\mathbf{u})}{p_*(\mathbf{x})} = \exp(f_*(\mathbf{x}) - f_*(\mathbf{u})). \quad (30)$$

With H -smoothness, i.e., $\|\nabla^2 f_*\| \leq H$, we have

$$\begin{aligned} f_*(\mathbf{x}) - f_*(\mathbf{u}) &\leq \nabla f_*(\mathbf{u}) \cdot (\mathbf{x} - \mathbf{u}) + \frac{H}{2} \cdot \|\mathbf{u} - \mathbf{x}\|^2 \\ &\leq \|\nabla f_*(\mathbf{u})\| \cdot \|\mathbf{x} - \mathbf{u}\| + \frac{H}{2} \cdot \|\mathbf{x} - \mathbf{u}\|^2. \end{aligned} \quad (31)$$

Since $\mathbf{u}, \mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})$, and each cell is an axis-aligned hypercube of side length l , we have

$$\|\mathbf{x} - \mathbf{u}\|^2 = \sum_{i=1}^d \|\mathbf{x}_i - \mathbf{u}_i\|^2 \leq dl^2.$$

Let $G_0 := \|\nabla f_*(\mathbf{0})\|$. Then we have

$$\|\nabla f_*(\mathbf{u})\| \leq \|\nabla f_*(\mathbf{u}) - \nabla f_*(\mathbf{0})\| + G_0 \leq H \cdot 2L + G_0,$$

where the last inequality follows from $\mathbf{u} \in \text{Cube}(L)$. Therefore, by requiring

$$l \leq \frac{\epsilon}{\sqrt{d} \cdot (2HL + G_0)},$$

and $\epsilon \leq 8HL^2$ without loss of generality, we will have $l \leq \sqrt{2\epsilon/(dH)}$, which means

$$\|\nabla f_*(\mathbf{u})\| \cdot \|\mathbf{x} - \mathbf{u}\| + \frac{H}{2} \cdot \|\mathbf{x} - \mathbf{u}\|^2 \leq (2HL + G_0) \cdot \sqrt{dl} + \frac{H}{2} \cdot dl^2 \leq 2\epsilon. \quad (32)$$

Plugging Eq. (31) and Eq. (32) into Eq. (30), we have

$$\frac{\tilde{p}_*(\mathbf{u})}{\tilde{p}_*(\mathbf{x})} \leq \exp(2\epsilon) \leq (1 + 4\epsilon). \quad (33)$$

With a similar technique, we have

$$-(f_*(\mathbf{x}) - f_*(\mathbf{u})) = f_*(\mathbf{u}) - f_*(\mathbf{x}) \leq \|\nabla f_*(\mathbf{x})\| \cdot \|\mathbf{x} - \mathbf{u}\| + \frac{H}{2} \cdot \|\mathbf{x} - \mathbf{u}\|^2.$$

Under the same setting, it implies

$$\frac{\tilde{p}_*(\mathbf{x})}{\tilde{p}_*(\mathbf{u})} \leq \exp(2\epsilon) \quad \Leftrightarrow \quad \frac{\tilde{p}_*(\mathbf{u})}{\tilde{p}_*(\mathbf{x})} \geq \exp(-2\epsilon) \geq 1 - 2\epsilon. \quad (34)$$

Combining Eq. (33) with Eq. (34), we have

$$1 - 4\epsilon \leq \frac{\tilde{p}_*(\mathbf{u})}{\tilde{p}_*(\mathbf{x})} \leq 1 + 4\epsilon. \quad (35)$$

Hence we are able to control the TV distance between \tilde{p}_* and \bar{p}_* , i.e.,

$$\begin{aligned} \text{TV}(\bar{p}_*, \tilde{p}_*) &= \frac{1}{2} \int_{\mathbf{x} \in \text{Cube}(L)} |\bar{p}_*(\mathbf{x}) - \tilde{p}_*(\mathbf{x})| \, d\mathbf{x} \\ &= \frac{1}{2} \sum_{i_0, i_1, \dots, i_{d-1}} \int_{\mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} |\bar{p}_*(\mathbf{x}) - \tilde{p}_*(\mathbf{x})| \, d\mathbf{x} \\ &= \frac{1}{2} \sum_{i_0, i_1, \dots, i_{d-1}} \int_{\mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} |\tilde{p}_*(\bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}}) - \tilde{p}_*(\mathbf{x})| \, d\mathbf{x} \\ &= \frac{1}{2} \sum_{i_0, i_1, \dots, i_{d-1}} \int_{\mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} \tilde{p}_*(\mathbf{x}) \left| \frac{\tilde{p}_*(\bar{\mathbf{x}}_{i_0, i_1, \dots, i_{d-1}})}{\tilde{p}_*(\mathbf{x})} - 1 \right| \, d\mathbf{x} \\ &\leq \frac{1}{2} \sum_{i_0, i_1, \dots, i_{d-1}} \int_{\mathbf{x} \in \text{Cell}(i_0, i_1, \dots, i_{d-1})} \tilde{p}_*(\mathbf{x}) 4\epsilon \, d\mathbf{x} \\ &= 2\epsilon, \end{aligned} \quad (36)$$

where the last inequality follows from Eq. (35). Hence, the proof is completed. \blacksquare

Lemma 12 Suppose the data distribution p_* satisfy Assumption [A2]–[A3], we have

$$\|\nabla f_*(\mathbf{0})\|^2 \leq 2Hd + 2H^2m_0$$

Proof We start with the following inequality

$$\begin{aligned} \|\nabla f_*(\mathbf{0})\|^2 &= \int_{\mathbf{x} \in \mathbb{R}^d} p_*(\mathbf{x}) \|\nabla f_*(\mathbf{0})\|^2 \, d\mathbf{x} \\ &\leq 2 \int_{\mathbf{x} \in \mathbb{R}^d} p_*(\mathbf{x}) \|\nabla f_*(\mathbf{x})\|^2 \, d\mathbf{x} + 2 \int_{\mathbf{x} \in \mathbb{R}^d} p_*(\mathbf{x}) \|\nabla f_*(\mathbf{0}) - \nabla f_*(\mathbf{x})\|^2 \, d\mathbf{x} \\ &\leq 2Hd + 2H^2 \int_{\mathbf{x} \in \mathbb{R}^d} p_*(\mathbf{x}) \|\mathbf{x}\|^2 \, d\mathbf{x} = 2Hd + 2H^2m_0 \end{aligned}$$

where the second inequality follows from Lemma 20 and Assumption [A3] and the last inequality follows from Assumption [A2]. Hence, the proof is completed. \blacksquare

Proof [Proof of Lemma 1] The TV distance between the original data distribution p_* and the histogram approximation \bar{p}_* can be written as

$$\text{TV}(p_*, \bar{p}_*) \leq \text{TV}(p_*, \tilde{p}_*) + \text{TV}(\tilde{p}_*, \bar{p}_*).$$

Following from Lemma 10, we will have $\text{TV}(p_*, \tilde{p}_*)$ by choosing

$$L \geq \sigma \cdot \sqrt{2 \ln(2d/\epsilon)}. \quad (37)$$

Moreover, with the quantization shown in Eq. (11), it has $\text{TV}(\tilde{p}_*, \bar{p}_*) \leq 2\epsilon$ by choosing

$$l \leq (2HL + \|\nabla f_*(\mathbf{0})\|)^{-1} \cdot d^{-1/2}\epsilon, \quad (38)$$

which follows from Lemma 11. Combining Eq. (37) with Eq. (38), if we set

$$L = \sigma \cdot \sqrt{2 \ln(2d/\epsilon)} \quad \text{and} \quad l := \frac{\epsilon}{2H(L\sqrt{d} + d + \sqrt{dm_0})}$$

and l satisfies

$$\begin{aligned} l &\leq \frac{\epsilon}{(2HL + 2\sqrt{Hd} + 2H\sqrt{m_0})\sqrt{d}} \leq \frac{\epsilon}{(2HL + \sqrt{2Hd} + 2H^2m_0)\sqrt{d}} \\ &\leq (2HL + \|\nabla f_*(\mathbf{0})\|)^{-1} \cdot d^{-1/2}\epsilon \end{aligned}$$

where the last inequality follows from Lemma 12. That means

$$l = \Omega \left(\left[2H \cdot \left(\sigma \sqrt{2d \ln(2d/\epsilon)} + d + \sqrt{dm_0} \right) \right]^{-1} \cdot \epsilon \right),$$

it will have $\text{TV}(p_*, \bar{p}_*) \leq 3\epsilon$. Hence, the proof is completed. \blacksquare

Appendix D. Proof of Lemma 2

To make our analysis clear, we define the variables, the random variables, and the marginal density derived by a specific ordered set $S \subseteq \{0, 1, \dots, d \log_2 K - 1\}$. Specifically, we have

$$\mathbf{y}_S = \sum_{i=0}^{|S|-1} \mathbf{e}_i \cdot y_{S_i} \quad \text{and} \quad \mathbf{y}_{t,S} = \sum_{i=0}^{|S|-1} \mathbf{e}_i \cdot y_{t,S_i}$$

where there are

$$\mathbf{y} = [y_0, y_1, \dots, y_{d \log_2 K - 1}] \quad \text{and} \quad \mathbf{y}_t = [y_{t,0}, y_{t,1}, \dots, y_{t,d \log_2 K - 1}].$$

Suppose $\mathbf{y}_t \sim q_t$. The underlying distribution of $\mathbf{y}_{t,S}$ is denoted as

$$q_{t,S}(\mathbf{y}_S) = \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} q_t(\tilde{\mathbf{y}}) \cdot \mathbf{1}_{\mathbf{y}_S}(\tilde{\mathbf{y}}_S).$$

Lemma 13 (Modified log-Sobolev inequality for the forward process) *Suppose the transition rate function R^\rightarrow of the CTMC $\{\mathbf{y}_t^\rightarrow\}_{t=0}^T$ be defined as Eq. (13). CTMC satisfies modified log-Sobolev inequality with a constant 2, that is to say, for any $f \in \mathbb{L}_2(q_\infty^\rightarrow)$, it has*

$$\text{Ent}_{q_\infty^\rightarrow}[f] \leq \mathcal{E}(f, \ln f)$$

where Ent and \mathcal{E} denote the entropy and the Dirichlet functional.

Proof We start from the setting of the transition rate matrix of the forward process shown in Eq. (13). Combining with the Eq. (19), the infinitesimal generator for the forward process can be obtained, i.e.,

$$\mathcal{L}^{\rightarrow}[f](\mathbf{y}) = \langle f, R^{\rightarrow}(\cdot, \mathbf{y}) \rangle_{\mathcal{Y}}. \quad (39)$$

To verify the modified log-Sobolev inequality, we first require to calculate the Dirichlet functional $\mathcal{E}(f, \ln f)$. Here \mathcal{E} denotes the Dirichlet functional

$$\mathcal{E}(f, g) := \int \Gamma(f, g) dq_{\infty}^{\rightarrow},$$

where q_{∞} denotes the invariant measure of this forward process and Γ denotes the carré du champ operator, i.e.,

$$\Gamma(f, g) := \frac{1}{2} (\mathcal{L}[f \cdot g] - f \cdot \mathcal{L}[g] - g \cdot \mathcal{L}[f]).$$

Specifically, presenting the transition rate matrix to be a matrix version Eq. (25), we have

$$dq_t^{\rightarrow}/dt = \mathbf{R}^{\rightarrow} \cdot \mathbf{q}_t^{\rightarrow}$$

Combining the fact $\mathbf{1} \cdot \mathbf{R} = \mathbf{0}$ and \mathbf{R} is symmetric, the RHS of the above equation satisfies

$$\mathbf{R}^{\rightarrow} \cdot 2^{-d \log_2 K} \cdot \mathbf{1} = \mathbf{0},$$

which implies the uniform distribution coincides with the invariant measure of q_{∞}^{\rightarrow} . Then, for the Dirichlet functional, it has

$$\begin{aligned} \mathcal{E}(f, \ln f) &= \frac{1}{2} \int \mathcal{L}[f \cdot \ln f](\mathbf{y}) - f(\mathbf{y}) \cdot \mathcal{L}[\ln f](\mathbf{y}) - \ln f(\mathbf{y}) \cdot \mathcal{L}[f](\mathbf{y}) dq_{\infty}(\mathbf{y}) \\ &= \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{Y}} q_{\infty}(\mathbf{y}) \cdot \left[\sum_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{y}') \ln f(\mathbf{y}') \cdot R(\mathbf{y}', \mathbf{y}) - f(\mathbf{y}) \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} \ln f(\mathbf{y}') \cdot R(\mathbf{y}', \mathbf{y}) \right. \\ &\quad \left. - \ln f(\mathbf{y}) \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{y}') \cdot R(\mathbf{y}', \mathbf{y}) + f(\mathbf{y}) \cdot \ln f(\mathbf{y}) \cdot \underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}} R(\mathbf{y}', \mathbf{y})}_{=0} \right] \\ &= \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{\mathbf{y}' \in \mathcal{Y}} q_{\infty}(\mathbf{y}) (f(\mathbf{y}) - f(\mathbf{y}')) \cdot R(\mathbf{y}', \mathbf{y}) \cdot (\ln f(\mathbf{y}) - \ln f(\mathbf{y}')). \end{aligned}$$

Plugging the definition of \mathbf{R} into the above equation, we have

$$\begin{aligned} \mathcal{E}(f, \ln f) &= \frac{1}{2} \cdot \sum_{\mathbf{y} \in \mathcal{Y}} q_{\infty}(\mathbf{y}) \cdot \sum_{i=0}^{d \log_2 K - 1} \sum_{\tilde{y}_i \in \{0,1\}} (f(\mathbf{y}) - f(\mathbf{y} + (\tilde{y}_i - y_i) \cdot \mathbf{e}_i)) \\ &\quad \cdot (\ln f(\mathbf{y}) - \ln f(\mathbf{y} + (\tilde{y}_i - y_i) \cdot \mathbf{e}_i)). \end{aligned} \quad (40)$$

Then, we consider $\text{Ent}_{q_\infty^\rightarrow}[f]$, which satisfies

$$\begin{aligned} \text{Ent}_{q_\infty^\rightarrow}[f] &= \mathbb{E}_{\mathbf{y} \sim q_\infty^\rightarrow} [f(\mathbf{y}) \ln f(\mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim q_\infty^\rightarrow} [f(\mathbf{y})] \ln (\mathbb{E}_{\mathbf{y} \sim q_\infty^\rightarrow} [f(\mathbf{y})]) \\ &\leq \sum_{i=0}^{d \log_2 K - 1} \mathbb{E}_{\mathbf{y}_{[0:i-1, i+1:d \log_2 K-1]}} \left[\underbrace{\mathbb{E}_{\mathbf{y}_i} [f(\mathbf{y}) \ln f(\mathbf{y})] - \mathbb{E}_{\mathbf{y}_i} [f(\mathbf{y})] \ln (\mathbb{E}_{\mathbf{y}_i} [f(\mathbf{y})])}_{\text{Term 1}} \right]. \end{aligned} \quad (41)$$

due to the sub-additivity of the entropy, i.e., Lemma 17. Term 1 of Eq. (41) satisfies

$$\begin{aligned} \text{Term 1} &= \sum_{\mathbf{y}_i \in \{0,1\}} q_{\infty,i}^\rightarrow(\mathbf{y}_i) \cdot f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) \ln f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) \\ &\quad - \sum_{\mathbf{y}_i \in \{0,1\}} q_{\infty,i}^\rightarrow(\mathbf{y}_i) f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) \\ &\quad \cdot \ln \left(\sum_{\tilde{\mathbf{y}}_i \in \{0,1\}} q_{\infty,i}^\rightarrow(\tilde{\mathbf{y}}_i) f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1}) \right) \\ &\leq \sum_{\mathbf{y}_i \in \{0,1\}} q_{\infty,i}^\rightarrow(\mathbf{y}_i) \cdot f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) \\ &\quad \cdot \sum_{\tilde{\mathbf{y}}_i \in \{0,1\}} \left[\frac{\ln f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1})}{2} - \frac{\ln f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1})}{2} \right] \\ &\leq \frac{1}{2} \sum_{\mathbf{y}_i, \tilde{\mathbf{y}}_i \in \{0,1\}} q_{\infty,i}^\rightarrow(\mathbf{y}_i) \cdot (f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) - f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1})) \\ &\quad \cdot (\ln f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) - \ln f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1})), \end{aligned}$$

where the first inequality follows from the concavity of the logarithm function, and the last inequality follows from

$$\begin{aligned} &\sum_{\mathbf{y}, \tilde{\mathbf{y}}} f(\mathbf{y}) \cdot (\ln f(\mathbf{y}) - \ln f(\tilde{\mathbf{y}})) = \sum_{\tilde{\mathbf{y}}} f(\tilde{\mathbf{y}}) \cdot (\ln f(\tilde{\mathbf{y}}) - \ln f(\mathbf{y})) \\ &= \frac{1}{2} \sum_{\mathbf{y}, \tilde{\mathbf{y}}} (f(\mathbf{y}) - f(\tilde{\mathbf{y}})) \cdot (\ln f(\mathbf{y}) - \ln f(\tilde{\mathbf{y}})) \end{aligned}$$

and $q_\infty^\rightarrow(\cdot)$ is a constant function. Then, plugging this inequality into Eq. (41), we have

$$\begin{aligned} \text{Ent}_{q_\infty^\rightarrow}[f] &\leq \frac{1}{2} \cdot \sum_{i=0}^{d \log_2 K - 1} \left[\sum_{\mathbf{y}_{[0:i-1, i+1:d \log_2 K-1]}} q_{\infty, [0:i-1, i+1:d \log_2 K-1]}^\rightarrow(\mathbf{y}_{[0:i-1, i+1:d \log_2 K-1]}) \right. \\ &\quad \sum_{\mathbf{y}_i} q_{\infty,i}^\rightarrow(\mathbf{y}_i) \sum_{\tilde{\mathbf{y}}_i} (f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) - f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1})) \\ &\quad \cdot (\ln f(\mathbf{y}_{0:i-1}, \mathbf{y}_i, \mathbf{y}_{i+1:d \log_2 K-1}) - \ln f(\mathbf{y}_{0:i-1}, \tilde{\mathbf{y}}_i, \mathbf{y}_{i+1:d \log_2 K-1})) \left. \right] \\ &= \frac{1}{2} \cdot \sum_{\mathbf{y}} q_\infty^\rightarrow(\mathbf{y}) \cdot \sum_{i=0}^{d \log_2 K - 1} \sum_{\tilde{\mathbf{y}}_i} (f(\mathbf{y}) - f(\mathbf{y} + (\tilde{\mathbf{y}}_i - \mathbf{y}_i) \cdot \mathbf{e}_i)) \\ &\quad \cdot (\ln f(\mathbf{y}) - \ln f(\mathbf{y} + (\tilde{\mathbf{y}}_i - \mathbf{y}_i) \cdot \mathbf{e}_i)) \end{aligned} \quad (42)$$

Comparing Eq. (42) and Eq. (40), it satisfies

$$\text{Ent}_{q_\infty} [f] \leq \frac{C_{\text{LSI}}}{2} \cdot \mathcal{E}(f, \ln f)$$

by choosing $C_{\text{LSI}} = 2$. ■

Proof [Proof of Lemma 2] We investigate the dynamic of KL divergence between q_t^\rightarrow and q_∞^\rightarrow in the forward process. Specifically, we have

$$\begin{aligned} \frac{d\text{KL}(q_t^\rightarrow \| q_\infty^\rightarrow)}{dt} &= \sum_{\mathbf{y} \in \mathcal{Y}} \frac{dq_t^\rightarrow(\mathbf{y})}{dt} \cdot \ln \frac{q_t^\rightarrow(\mathbf{y})}{q_\infty^\rightarrow(\mathbf{y})} = \sum_{\mathbf{y} \in \mathcal{Y}} \ln \frac{q_t^\rightarrow(\mathbf{y})}{q_\infty^\rightarrow(\mathbf{y})} \left(\sum_{\mathbf{y}_0 \in \mathcal{Y}} R(\mathbf{y}, \mathbf{y}_0) \cdot q_t^\rightarrow(\mathbf{y}_0) \right) \\ &= \sum_{\mathbf{y}_0} q_\infty^\rightarrow(\mathbf{y}_0) \cdot \frac{q_t^\rightarrow(\mathbf{y}_0)}{q_\infty^\rightarrow(\mathbf{y}_0)} \cdot \sum_{\mathbf{y}} \ln \frac{q_t^\rightarrow(\mathbf{y})}{q_\infty^\rightarrow(\mathbf{y})} \cdot R^\rightarrow(\mathbf{y}, \mathbf{y}_0) \\ &= \sum_{\mathbf{y}_0} q_\infty^\rightarrow(\mathbf{y}_0) \cdot \frac{q_t^\rightarrow(\mathbf{y}_0)}{q_\infty^\rightarrow(\mathbf{y}_0)} \cdot \mathcal{L}[\ln \frac{q_t^\rightarrow}{q_\infty^\rightarrow}](\mathbf{y}') = -\mathcal{E} \left(\frac{q_t^\rightarrow}{q_\infty^\rightarrow}, \ln \frac{q_t^\rightarrow}{q_\infty^\rightarrow} \right) \end{aligned}$$

Due to Lemma 13, we have

$$\frac{d\text{KL}(q_t^\rightarrow \| q_\infty^\rightarrow)}{dt} = -\mathcal{E} \left(\frac{q_t^\rightarrow}{q_\infty^\rightarrow}, \ln \frac{q_t^\rightarrow}{q_\infty^\rightarrow} \right) \leq \text{Ent}_{q_\infty} \left[\frac{q_t^\rightarrow}{q_\infty^\rightarrow} \right] = -\text{KL}(q_t^\rightarrow \| q_\infty^\rightarrow).$$

According to the Gronwall's theorem, we have

$$\text{KL}(q_t^\rightarrow \| q_\infty^\rightarrow) \leq e^{-t} \cdot \text{KL}(q_0^\rightarrow \| q_\infty^\rightarrow).$$

Combining with the following initialization error bound,

$$\text{KL}(q_0^\rightarrow \| q_\infty^\rightarrow) = \sum_{\mathbf{y} \in \mathcal{Y}} q_0^\rightarrow(\mathbf{y}) \ln \frac{q_0^\rightarrow(\mathbf{y})}{2^{-d \log_2 K}} \leq d \log_2 K.$$

Hence, the proof is completed. ■

Appendix E. Supplementary Proofs for the Discrete Reverse Process

E.1. Proof of Lemma 3

Proof [Proof of Lemma 3 (adapted from Proposition 5 of Chen and Ying (2024))] Suppose the transition rate function R^\rightarrow of the CTMC $\{\mathbf{y}_t^\rightarrow\}_{t=0}^T$ be defined as Eq. (13), the marginal distribution at time t can be written as

$$q_t^\rightarrow(\mathbf{y}) = \sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^\rightarrow(\mathbf{y}_0) \cdot q_{t|0}^\rightarrow(\mathbf{y}|\mathbf{y}_0).$$

Define the plus operator as follows

$$\mathbf{y} \oplus \mathbf{e}_i = [y_0, y_1, \dots, y_{i-1}, (y_i + 1) \bmod 2, y_{i+1}, \dots, y_{d \log_2 K - 1}],$$

then we have

$$\frac{q_t^{\rightarrow}(\mathbf{y} \oplus \mathbf{e}_i)}{q_t^{\rightarrow}(\mathbf{y})} = \frac{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{t|0}^{\rightarrow}(\mathbf{y} + \mathbf{e}_i | \mathbf{y}_0)}{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0)} = \frac{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0) \cdot \frac{q_{t|0}^{\rightarrow}(\mathbf{y} + \mathbf{e}_i | \mathbf{y}_0)}{q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0)}}{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0)}.$$

According to Bayes Theorem, we have

$$q_{0|t}^{\rightarrow}(\mathbf{y}_0 | \mathbf{y}) \cdot q_t^{\rightarrow}(\mathbf{y}) = q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0) \cdot q_0^{\rightarrow}(\mathbf{y}_0) \quad \Leftrightarrow \quad q_{0|t}^{\rightarrow}(\mathbf{y}_0 | \mathbf{y}) \propto q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0) \cdot q_0^{\rightarrow}(\mathbf{y}_0),$$

which implies

$$\frac{q_t^{\rightarrow}(\mathbf{y} \oplus \mathbf{e}_i)}{q_t^{\rightarrow}(\mathbf{y})} = \mathbb{E}_{\mathbf{y}_0 \sim q_{0|t}^{\rightarrow}(\cdot | \mathbf{y})} \left[\frac{q_{t|0}^{\rightarrow}(\mathbf{y} + \mathbf{e}_i | \mathbf{y}_0)}{q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0)} \right].$$

With Lemma 8, we have

$$\frac{q_{t|0}^{\rightarrow}(\mathbf{y} + \mathbf{e}_i | \mathbf{y}_0)}{q_{t|0}^{\rightarrow}(\mathbf{y} | \mathbf{y}_0)} = \frac{1 + (-1)^{|(y_i + 1 - y_{0,i}) \bmod 2|} \cdot e^{-2t}}{1 + (-1)^{|(y_i - y_{0,i}) \bmod 2|} \cdot e^{-2t}} \leq \frac{1 + e^{-2t}}{1 - e^{-2t}},$$

which means

$$\frac{q_t^{\rightarrow}(\mathbf{y} \oplus \mathbf{e}_i)}{q_t^{\rightarrow}(\mathbf{y})} \leq \frac{1 + e^{-2t}}{1 - e^{-2t}} \leq 1 + t^{-1}.$$

Therefore, if we consider the transition rate matrix of the reverse process, i.e.,

$$R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) := R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot \frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})}$$

provided by Eq (3), it has

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) = \sum_{i=0}^{d \log_2 K - 1} \frac{q_t^{\leftarrow}(\mathbf{y} \oplus \mathbf{e}_i)}{q_t^{\leftarrow}(\mathbf{y})} = \sum_{i=0}^{d \log_2 K - 1} \frac{q_{T-t}^{\rightarrow}(\mathbf{y} \oplus \mathbf{e}_i)}{q_{T-t}^{\rightarrow}(\mathbf{y})} \leq (d \log_2 K) \cdot (1 + (T-t)^{-1}).$$

Hence, the proof is completed. ■

E.2. Proof of Theorem 6

The ultimate target of Alg. 2 is to generate sample $\hat{\mathbf{x}}$ and require its underlying distribution \hat{p} to be close to the continuous data distribution p_* . However, Alg. 2 can be divided into two parts:

1. **Truncated Uniformization:** Generate a discrete sample following $\hat{q}_{T-\delta} = \hat{q}_{t_W}$ which approximates q_* , which is from Step. 2 to Step. 10.
2. Mapping the generated discrete data to the corresponding cell in Euclidean space and uniformly drawing a sample from the cell, which is from Step. 12

All the following notations correspond to those mentioned in Alg. 2.

Lemma 14 *Suppose we have a timestamp sequence satisfying*

$$t_0 = 0 \quad \text{and} \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}),$$

then we know the sequence $\{t_w\}_{w=0}^W$ is strict increasing and $t_W < T$ for any W .

Proof According to the timestamp setting, i.e.,

$$t_0 = 0 \quad \text{and} \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}),$$

solve for t_{k+1} , we have

$$t_{w+1} = \frac{0.5T + t_w}{1.5} = \frac{T + 2t_w}{3}.$$

Then, we consider the difference:

$$t_{w+1} - t_w = \frac{T + 2t_w}{3} - t_w = \frac{T + 2t_w - 3t_w}{3} = \frac{T - t_w}{3}.$$

If $T - t_w > 0$, then we have

$$t_{w+1} - t_w = \frac{T - t_w}{3} > 0,$$

which shows $t_{w+1} > t_w$. Thus, as long as $t_w < T$, the sequence is strictly increasing.

Moreover, due to the fact $t_0 = 0 < T$, we can prove that $t_w < T$ for all w . Specifically, assume $t_w < T$; then

$$t_{w+1} = \frac{T + 2t_w}{3} < \frac{T + 2T}{3} = T.$$

Therefore, $t_{w+1} < T$ as well, completing the induction. Hence t_w remains below T for all w , and the sequence $\{t_w\}$ is strictly increasing. \blacksquare

Lemma 15 *Suppose the reverse process is divided into W segments with endpoints $\{t_w\}_{w=0}^W$ satisfying*

$$t_0 = 0, \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}) \quad \text{and} \quad t_W = T - \delta,$$

if we set

$$\beta_{t_w} := 2d \log_2 K / \min\{1, T - t_w\}$$

then we have

$$\sum_{k=1}^W \beta_{t_k} \cdot (t_k - t_{k-1}) \leq 2d \log_2 K \cdot (T + \ln(1/\delta))$$

Proof [Adapted from Theorem 6 of [Chen and Ying \(2024\)](#)] Suppose there exist time steps t_0, t_1, \dots, t_W such that $T - t_w = s_w$ for each $w = 0, \dots, W$. According to Lemma 14, we know $\{t_w\}_{w=0}^W$ is a increasing sequence, if we set

$$s_w := T - t_w,$$

then it can be expected that $s_0 > s_1 > \dots > s_W \geq \delta > 0$. According to the choice of β_w , it has

$$\beta_w = \frac{Cd \log_2 K}{\min(1, s_w)}, \quad \text{and} \quad s_{w-1} - s_w > 0.$$

For w such that $\delta \leq s_w < 1$, notice that $\min(1, s_w) = s_w$, we have $\beta_w = Cd \log_2 K / s_w$ and

$$\sum_{w:\delta \leq s_w < 1} \beta_w \cdot (t_w - t_{w-1}) = \sum_{w:\delta \leq s_w < 1} \beta_w (s_{w-1} - s_w) = \sum_{w:\delta \leq s_w < 1} \frac{Cd \log_2 K}{s_w} (s_{w-1} - s_w).$$

Because $1/s$ is a decreasing function for $s > 0$, we have

$$\frac{1}{s_w} \leq \frac{1}{s} \quad \text{for all } s \in [s_w, s_{w-1}],$$

which implies

$$\frac{Cd}{s_w} (s_{w-1} - s_w) \leq Cd \log_2 K \int_{s_w}^{s_{w-1}} \frac{1}{s} ds.$$

Hence,

$$\sum_{w:\delta \leq s_w < 1} \frac{Cd \log_2 K}{s_w} (s_{w-1} - s_w) \leq Cd \log_2 K \sum_{w:\delta \leq s_w < 1} \int_{s_w}^{s_{w-1}} \frac{1}{s} ds = Cd \log_2 K \int_{\delta}^1 \frac{1}{s} ds.$$

Evaluating the integral on the right gives

$$Cd \log_2 K \int_{\delta}^1 \frac{1}{s} ds = Cd \log_2 K [\ln(s)]_{\delta}^1 = Cd \log_2 K \ln(1/\delta).$$

Therefore, we have established the exact upper bound

$$\sum_{k:\delta \leq s_k < 1} \lambda_k (s_{k-1} - s_k) \leq Cd \log_2 K \ln(1/\delta).$$

For $s_w \geq 1$, notice that $\min(1, s_w) = 1$, we have $\beta_w = Cd \log_2 K$.

$$\begin{aligned} \sum_{w:1 \leq s_w \leq T} \beta_w \cdot (t_w - t_{w-1}) &= \sum_{w:1 \leq s_w \leq T} \beta_w (s_{w-1} - s_w) \\ &= \sum_{w:1 \leq s_w \leq T} Cd \log_2 K \cdot (s_{w-1} - s_w) \leq Cd \log_2 K \cdot (T - 1). \end{aligned}$$

Combining the two parts, we have

$$\begin{aligned} \sum_{w=1}^W \beta_w \cdot (t_w - t_{w-1}) &= \sum_{w=1}^W \beta_w \cdot (s_{w-1} - s_w) \\ &= \sum_{w:\delta \leq s_w < 1} \beta_w (s_{w-1} - s_w) + \sum_{w:1 \leq s_w \leq T} \beta_w (s_{w-1} - s_w) \leq Cd \log_2 K \cdot (T + \ln(1/\delta)). \end{aligned}$$

Hence, the proof is completed. ■

Lemma 16 *Following the notations shown in Section 2, we have*

$$\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] = \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}).$$

Proof Since we have required $\Delta t \rightarrow 0$, that is to say

$$\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) \rightarrow \hat{q}_{t|t}(\mathbf{y}'|\mathbf{y}) = 0 \quad \text{and} \quad q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \rightarrow q_{t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) = 0 \quad \forall \mathbf{y}' \neq \mathbf{y},$$

which automatically makes

$$\left| \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right| \leq \frac{1}{2} < 1.$$

Under this condition, we have

$$\begin{aligned} \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} &= \ln \left[1 + \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] \\ &= \sum_{i=1}^{\infty} \frac{(-1)^{i+1}}{i} \cdot \left[\frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right]^i, \end{aligned}$$

which implies (with the dominated convergence theorem)

$$\begin{aligned} &\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] \\ &= \sum_{i=1}^{\infty} \frac{(-1)^{i+1}}{i} \cdot \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{\Delta t} \\ &\quad \cdot \lim_{\Delta t \rightarrow 0} \frac{\left(\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right) \right)^{i-1}}{\left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) \right)^i}. \end{aligned}$$

Only when $i = 1$, we have

$$\lim_{\Delta t \rightarrow 0} \frac{\left(\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right) \right)^{i-1}}{\left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) \right)^i} = 1,$$

otherwise it will be equivalent to 0. Therefore, we have

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] &= \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{\Delta t} \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{R}_t(\mathbf{y}', \mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right) = \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}). \end{aligned}$$

Hence, the proof is completed. ■

E.3. The Proof of Lemma 4

Proof We start from the dynamic of KL divergence with the time growth in the reverse process, i.e.,

$$\begin{aligned} \frac{d\text{KL}(q_t^\leftarrow \parallel \hat{q}_t)}{dt} &= \lim_{\Delta t \rightarrow 0} \left[\frac{\text{KL}(q_{t+\Delta t}^\leftarrow \parallel \hat{q}_{t+\Delta t}) - \text{KL}(q_t^\leftarrow \parallel \hat{q}_t)}{\Delta t} \right] \\ &\leq \lim_{\Delta t \rightarrow 0} \left[\frac{\mathbb{E}_{\mathbf{y} \sim q_t^\leftarrow} \left[\text{KL}(q_{t+\Delta t|t}^\leftarrow(\cdot|\mathbf{y}) \parallel \hat{q}_{t+\Delta t|t}(\cdot|\mathbf{y})) \right]}{\Delta t} \right] \end{aligned}$$

where the inequality follows from the chain rule of KL divergence. Under this condition, we have

$$\frac{d\text{KL}(q_t^\leftarrow \parallel \hat{q}_t)}{dt} \leq \sum_{\mathbf{y} \in \mathcal{Y}} q_t^\leftarrow(\mathbf{y}) \cdot \underbrace{\lim_{\Delta t \rightarrow 0} \left[\frac{\text{KL}(q_{t+\Delta t|t}^\leftarrow(\cdot|\mathbf{y}) \parallel \hat{q}_{t+\Delta t|t}(\cdot|\mathbf{y}))}{\Delta t} \right]}_{\text{Term 1}}. \quad (43)$$

For each $\mathbf{y} \in \mathcal{Y}$, we focus on Term 1 of Eq. (43), and have

$$\begin{aligned} \text{Term 1} &= \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y}) \cdot \ln \frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] \\ &= \lim_{\Delta t \rightarrow 0} \underbrace{\left[\sum_{\mathbf{y}' \neq \mathbf{y}} \frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right]}_{\text{Term 1.1}} + \\ &\quad \underbrace{\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y}) \right) \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right]}_{\text{Term 1.2}}. \end{aligned} \quad (44)$$

For Term 1.1, we have

$$\begin{aligned} \text{Term 1.1} &= \sum_{\mathbf{y}' \neq \mathbf{y}} \lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\Delta t} \right] \cdot \lim_{\Delta t \rightarrow 0} \left[\ln \frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^\leftarrow(\mathbf{y}', \mathbf{y}) \cdot \ln \left[\lim_{\Delta t \rightarrow 0} \left(\frac{q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{\Delta t} \cdot \frac{\Delta t}{\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right) \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^\leftarrow(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^\leftarrow(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})}, \end{aligned} \quad (45)$$

where the second equation follows from the composition rule of the limit calculation. For Term 1.2, we have

$$\begin{aligned} \text{Term 1.2} &= \lim_{\Delta t \rightarrow 0} \left[1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y}) \right] \cdot \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^\leftarrow(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y})} \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{R}_t(\mathbf{y}', \mathbf{y}) - R_t^\leftarrow(\mathbf{y}', \mathbf{y}) \right) = \hat{R}_t(\mathbf{y}) - R_t^\leftarrow(\mathbf{y}) \end{aligned} \quad (46)$$

where the first inequality follows from Lemma 16. Plugging Eq. (45), Eq. (46) and Eq. (44), into Eq. (43) we have

$$\frac{d\text{KL}(q_t^{\leftarrow} \parallel \hat{q}_t)}{dt} \leq \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \right). \quad (47)$$

For any $\mathbf{y} \in \mathcal{Y}$, we have

$$\begin{aligned} & \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \tilde{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \\ &+ \underbrace{\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{\tilde{R}_t(\mathbf{y}', \mathbf{y})}{\hat{R}_t(\mathbf{y}', \mathbf{y})} + \hat{R}_t(\mathbf{y}) - \tilde{R}_t(\mathbf{y})}_{\text{Term 2}}. \end{aligned} \quad (48)$$

When $\tilde{R}_t(\mathbf{y}) \leq \beta_t$, we have

$$\hat{R}_t(\mathbf{y}', \mathbf{y}) = \tilde{R}_t(\mathbf{y}', \mathbf{y}) \quad \text{and} \quad \hat{R}_t(\mathbf{y}) = \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{R}_t(\mathbf{y}', \mathbf{y}) = \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_t(\mathbf{y}', \mathbf{y}) = \tilde{R}_t(\mathbf{y})$$

which implies Term 2 = 0 in Eq. (48). Otherwise, we have

$$\frac{\hat{R}_t(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} = \frac{\beta_t}{\tilde{R}_t(\mathbf{y})} \quad \text{and} \quad \frac{\hat{R}_t(\mathbf{y})}{\tilde{R}_t(\mathbf{y})} = \frac{\beta_t}{\tilde{R}_t(\mathbf{y})},$$

which implies

$$\begin{aligned} \text{Term 2} &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{\tilde{R}_t(\mathbf{y})}{\beta_t} + \beta_t - \tilde{R}_t(\mathbf{y}) \\ &= R_t^{\leftarrow}(\mathbf{y}) \cdot \ln \left[1 + \frac{\tilde{R}_t(\mathbf{y}) - \beta_t}{\beta_t} \right] + \beta_t - \tilde{R}_t(\mathbf{y}) \leq \beta_t \cdot \left[\frac{\tilde{R}_t(\mathbf{y}) - \beta_t}{\beta_t} \right] + \beta_t - \tilde{R}_t(\mathbf{y}) = 0. \end{aligned}$$

Combining with Eq. (48) and Eq. (47), we have

$$\begin{aligned} \frac{d\text{KL}(q_t^{\leftarrow} \parallel \hat{q}_t)}{dt} &\leq \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \tilde{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \right) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_t(\mathbf{y}', \mathbf{y}) - \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot \left[-\frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} + \hat{v}_{t, \mathbf{y}}(\mathbf{y}') + \frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} \ln \frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y}) \hat{v}_{t, \mathbf{y}}(\mathbf{y}')} \right] \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') D_{\phi} \left(\frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} \parallel \hat{v}_{t, \mathbf{y}}(\mathbf{y}') \right), \end{aligned} \quad (49)$$

where D_ϕ is the Bregman divergence with $\phi(c) = c \ln c$ (as Eq. (5)), and the last equation follows from the definition of Bregman divergence:

$$D_\phi(u||v) = \phi(u) - \phi(v) - \langle \nabla \phi(v), u - v \rangle = u \ln \frac{u}{v} - u + v.$$

Then, by Eq. (5) and Assumption [A1], we have

$$\int_0^{T-\delta} \text{dKL}(q_t^{\leftarrow} || \hat{q}_t) \leq (T - \delta) \epsilon_{\text{score}}^2.$$

Hence, the proof is completed. ■

Bounding $\text{TV}(q_*, q_\delta^{\rightarrow})$ We adopt the proof strategy of Theorem 6 in Chen and Ying (2024). Consider the forward process $(X_t)_{t \geq 0}$. By the coupling characterization of the total variation distance, we have

$$\text{TV}(q_*, q_\delta^{\rightarrow}) := \inf_{\gamma \in \Gamma(q_*, q_\delta^{\rightarrow})} \mathbb{P}_{(u,v) \sim \gamma}[u \neq v] \leq \mathbb{P}(X_0 \neq X_\delta),$$

where $\Gamma(q_*, q_\delta^{\rightarrow})$ is the set of all couplings of $(q_*, q_\delta^{\rightarrow})$, and the inequality holds because (X_0, X_δ) gives a coupling of $(q_*, q_\delta^{\rightarrow})$.

By the transition kernel given in (Chen and Ying, 2024, Proposition 3), we have

$$\mathbb{P}(X_0 = X_\delta) = \frac{1}{2^{d \log_2 K}} \prod_{i=1}^{d \log_2 K} (1 + (-1)^0 e^{-2\delta})^{d \log_2 K} = \left(\frac{1 + e^{-2\delta}}{2} \right)^{d \log_2 K} \geq e^{-\delta d \log_2 K},$$

where the inequality holds due to the convexity of the exponential function. Thus,

$$\text{TV}(q_*, q_\delta^{\rightarrow}) \leq 1 - e^{-\delta d \log_2 K} \quad (50)$$

Proof [Proof of Theorem 6] We start from the quantization algorithm, i.e., Alg. 1. Since the data distribution p_* is supposed to satisfy Assumption [A2]–[A4], by introducing Lemma 1, the histogram-like approximation \bar{p}_* will be close to p_* , i.e.,

$$\text{TV}(\bar{p}_*, p_*) \leq 3\epsilon$$

by choosing

$$L = \sigma \cdot \sqrt{2 \ln(2d/\epsilon)} \quad \text{and} \quad l = \left[2H \cdot \left(\sigma \sqrt{2d \ln(2d/\epsilon)} + d + \sqrt{dm_0} \right) \right]^{-1} \cdot \epsilon.$$

Under this condition, we have

$$\begin{aligned} K &= \frac{2L}{l} = 4H \cdot \left[2\sigma^2 d^{1/2} \cdot \ln \frac{2d}{\epsilon} + \sigma d \cdot \sqrt{2 \ln \frac{2d}{\epsilon}} + d^{1/2} m_0^{1/2} \cdot \sqrt{2 \ln \frac{2d}{\epsilon}} \right] \cdot \epsilon^{-1} \\ &\leq 24H\sigma^2 dm_0 \epsilon^{-1} \cdot \ln(2d/\epsilon) \end{aligned}$$

where the last inequality follows from $\sigma \geq 1$ and $m_0 \geq 1$ without loss of generality. Then, after the training, the implementation of Alg. 2 requires $\bar{N} \sim \text{Poisson}(\bar{\beta})$ steps.

Proof of bound of the expectation of \bar{N} , i.e., $\bar{\beta}$. According to Lemma 15, if we set

$$t_0 = 0, \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}) \quad \text{and} \quad t_W = T - \delta,$$

for the time partitions,

$$\beta_{t_w} := 2d \log_2 K / \min\{1, T - t_w\}$$

for the intermediate Poisson, then it has

$$\begin{aligned} \bar{\beta} &= \sum_{k=1}^W \beta_{t_w} \cdot (t_w - t_{w-1}) \leq 2d \log_2 K \cdot (T + \ln(1/\delta)) \\ &\leq 2d \cdot [\log_2(24H\sigma^2) + \log_2(dm_0/\epsilon) + \log_2[\ln(2d/\epsilon)]] \cdot (T + \ln(1/\delta)). \end{aligned}$$

Proof of the TV distance bound. Since our truncated uniformization, i.e., Alg. 2, exactly simulates the reversed process, from Lemma 4, the KL divergence gap between $q_{T-\delta}^{\leftarrow} = q_{\delta}^{\rightarrow}$ and $\hat{q}_{T-\delta}$ is bounded by the KL divergence as follows:

$$\begin{aligned} \text{KL}(q_{T-\delta}^{\leftarrow} \parallel \hat{q}_{T-\delta}) &\leq \text{KL}(q_0^{\leftarrow} \parallel \hat{q}_0) + (T - \delta) \epsilon_{\text{score}}^2 \\ &\leq e^{-T} \cdot d \log_2 K + (T - \delta) \epsilon_{\text{score}}^2 = \epsilon^2 + (\ln(d/\epsilon) + \ln \log_2 K)^2 \cdot \epsilon_{\text{score}}^2 \leq 2\epsilon^2 \end{aligned} \quad (51)$$

where the second inequality follows from Lemma 2, the third inequality establishes when T is chosen as

$$T = \ln(d/\epsilon) + \ln \log_2 K,$$

and the last inequality is established when we have

$$\epsilon_{\text{score}} = \frac{\epsilon}{\ln(d/\epsilon) + \ln \log_2 K} = \tilde{O}(\epsilon).$$

Under this condition, due to Pinsker's inequality, Eq. (51) can be relaxed to

$$\text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) \leq \sqrt{\frac{\text{KL}(q_{T-\delta}^{\leftarrow} \parallel \hat{q}_{T-\delta})}{2}} \leq \epsilon.$$

Then we have

$$\begin{aligned} \text{TV}(q_*, \hat{q}_{T-\delta}) &\leq \text{TV}(q_*, q_{T-\delta}^{\leftarrow}) + \text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) \\ &= \text{TV}(q_*, q_{\delta}^{\rightarrow}) + \text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) = 1 - e^{-\delta d \log_2 K} + \epsilon \leq 2\epsilon \end{aligned}$$

where the second equation follows from Eq. (50) and the last inequality is established by requiring

$$\delta \leq \frac{\epsilon}{d \cdot \log_2 K} \quad \Leftrightarrow \quad \delta d \log_2 K \leq \epsilon.$$

Under this condition, we have

$$\delta d \log_2 K \leq \epsilon \leq \ln \frac{1}{1 - \epsilon} \quad \Rightarrow \quad 1 - e^{-\delta d \log_2 K} \leq \epsilon.$$

Suppose the underlying distributions of $\bar{\mathbf{y}}, \hat{\mathbf{x}}$ are \bar{q}, \hat{p} respectively, due to the connection between \hat{p}, \bar{p}_* and \bar{q}, \bar{q}_* shown in Eq. (12), we have

$$\text{TV}(\bar{p}_*, \hat{p}) = \int |\bar{p}_*(\mathbf{x}) - \hat{p}(\mathbf{x})| d\mathbf{x} = \sum_{\bar{\mathbf{y}} \in \bar{\mathcal{Y}}} |\bar{q}(\bar{\mathbf{y}}) - \bar{q}_*(\bar{\mathbf{y}})| = \sum_{\mathbf{y} \in \mathcal{Y}} |\hat{q}_{T-\delta}(\mathbf{y}) - \hat{q}_*(\mathbf{y})| \leq 2\epsilon.$$

Combining this result with Lemma 1, we have

$$\text{TV}(p_*, \hat{p}) \leq \text{TV}(p_*, \bar{p}_*) + \text{TV}(\bar{p}_*, \hat{p}) \leq 3\epsilon + 2\epsilon \leq 5\epsilon.$$

Hence, the proof is completed. ■

E.4. Proof of Theorem 7.

Proof We start from proving $W_2(p_*, \bar{p}_*) \leq 2\epsilon$ by analyzing the quantization Alg. 1, and then prove $\text{TV}(\bar{p}_*, \hat{p}) \leq 2\epsilon$. We emphasize that the whole analysis does not require the Lipschitzness assumption.

Proof of $W_2(p_*, \bar{p}_*) \leq 2\epsilon$: We prove $W_2(p_*, \tilde{p}_*) \leq \epsilon$ and $W_2(\bar{p}_*, \tilde{p}_*) \leq \epsilon$ both by considering the coupling, and then use triangle inequality to conclude. First, under the tail truncation step, i.e., Eq. (9) in our paper to construct \tilde{p}_* , i.e.,

$$\tilde{p}_*(x) := \frac{p_*(x)}{\int_{x \in \text{Cube}(L)} p_*(x) dx} \quad \forall x \in \text{Cube}(L).$$

To control the Wasserstein 2 distance between \tilde{p}_* and p_* , we construct the explicit coupling (\mathbf{x}, \mathbf{y}) in the following procedure: Let $\mathbf{x} \sim p_*$, and independently sample an auxiliary random variable $\mathbf{z} \sim \tilde{p}_*$, where $\mathbf{z} \perp \mathbf{x}$. Let $\mathbf{y} = \mathbf{x} \cdot \mathbf{1}\{\mathbf{x} \in \text{Cube}(L)\} + \mathbf{z} \cdot \mathbf{1}\{\mathbf{x} \notin \text{Cube}(L)\}$. Then, for any measurable $A \subseteq \text{Cube}(L)$, we have

$$\begin{aligned} \mathbb{P}(\mathbf{y} \in A) &= \mathbb{P}(\mathbf{x} \in A, \mathbf{x} \in \text{Cube}(L)) + \mathbb{P}(\mathbf{x} \notin \text{Cube}(L), \mathbf{z} \in A) \\ &= \mathbb{P}(\mathbf{x} \in A) + \mathbb{P}(\mathbf{x} \notin \text{Cube}(L)) \cdot \mathbb{P}(\mathbf{z} \in A) \\ &= p_*(A) + (1 - p_*(\text{Cube}(L))) \cdot \tilde{p}_*(A) \\ &= p_*(\text{Cube}(L)) \cdot \tilde{p}_*(A) + (1 - p_*(\text{Cube}(L))) \cdot \tilde{p}_*(A) = \tilde{p}_*(A). \end{aligned}$$

Therefore, the marginal distribution of \mathbf{y} is \tilde{p}_* .

Hence $(\mathbf{x}, \mathbf{y}) \sim \gamma$ gives a coupling of (p_*, \tilde{p}_*) . Then, we have

$$W_2^2(\tilde{p}_*, p_*) \leq \mathbb{E}_\gamma [\|\mathbf{x} - \mathbf{y}\|^2] \leq \underbrace{\mathbb{E}_\gamma [\|\mathbf{x} - \mathbf{y}\|^2 \mathbf{1}\{\mathbf{x} \in \text{Cube}(L)\}]}_{=0} + \mathbb{E}_\gamma [\|\mathbf{x} - \mathbf{y}\|^2 \mathbf{1}\{\mathbf{x} \notin \text{Cube}(L)\}]. \quad (52)$$

For the last term, since it has $\|\mathbf{y}\|^2 \leq dL^2 \leq \|\mathbf{x}\|^2$, we have

$$\mathbb{E}_\gamma [\|\mathbf{x} - \mathbf{y}\|^2 \mathbf{1}\{\mathbf{x} \notin \text{Cube}(L)\}] \leq 2\mathbb{E} [(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2) \mathbf{1}\{\mathbf{x} \notin \text{Cube}(L)\}] \leq 4\mathbb{E} [\|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\}].$$

Therefore, to upper bound $W_2^2(\tilde{p}_*, p_*)$, we bound $\mathbb{E} [\|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\}]$. We use the general identity (i.e., Layer cake representation) for non-negative random variables $W > 0$:

$$\mathbb{E}[W] = \int_0^\infty \mathbb{P}(W \geq t) dt.$$

By taking $W = \|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\}$, we have

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\}] &= \int_0^\infty \mathbb{P}(\|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\} \geq t) dt = \int_0^\infty \mathbb{P}(\|\mathbf{x}\|^2 \geq t, \|\mathbf{x}\| \geq L) dt \\ &= \int_0^{L^2} \mathbb{P}(\|\mathbf{x}\|^2 \geq L^2) dt + \int_{L^2}^\infty \mathbb{P}(\|\mathbf{x}\|^2 \geq t) dt. \end{aligned} \quad (53)$$

By choosing $u = e_i$ in Assumption [\[A4\]](#), we have that each dimension of x is σ sub-Gaussian, i.e.,

$$\mathbb{E}_{\mathbf{x}_i \sim p_{*,i}} [\exp(r\mathbf{x}_i)] \leq \exp\left(\frac{\sigma^2 r^2}{2}\right).$$

Thus we have $\mathbb{P}_i[|\mathbf{x}_i| \geq r] \leq 2 \exp\left(-\frac{r^2}{2\sigma^2}\right)$. With the union bound, we have for any $t > 0$,

$$\mathbb{P}_{\mathbf{x} \sim p_*}(\|\mathbf{x}\|^2 \geq t) \leq \mathbb{P}_{\mathbf{x} \sim p_*} \left[\max_{1 \leq i \leq d} |\mathbf{x}_i| > \sqrt{t/d} \right] \leq \sum_{i=1}^d \mathbb{P}_{\mathbf{x} \sim p_*} \left[|\mathbf{x}_i| > \sqrt{t/d} \right] \leq 2d \cdot \exp\left(-\frac{t}{2d\sigma^2}\right).$$

Therefore, the upper bound of Eq. (53) can be written as

$$\mathbb{E} [\|\mathbf{x}\|^2 \mathbf{1}\{\|\mathbf{x}\| \geq L\}] = \int_0^{L^2} \mathbb{P}(\|\mathbf{x}\|^2 \geq L^2) dt + \int_{L^2}^\infty \mathbb{P}(\|\mathbf{x}\|^2 \geq t) dt \leq 2d(L^2 + 2d\sigma^2) \exp\left(-\frac{L^2}{2d\sigma^2}\right).$$

Plugging this result into Eq. (52), we have

$$W_2^2(\tilde{p}_*, p_*) \leq 8d(L^2 + 2d\sigma^2) \exp\left(-\frac{L^2}{2d\sigma^2}\right).$$

Therefore, by choosing

$$L \geq \sqrt{2d}\sigma \cdot \sqrt{4 \ln \frac{2d\sigma}{\epsilon}},$$

we have

$$8d(L + 2d\sigma^2) \exp\left(-\frac{L^2}{2d\sigma^2}\right) \leq \epsilon^2.$$

Thus, we have $W_2(\tilde{p}_*, p_*) \leq \epsilon$.

Next, we consider the construction of \bar{p}_* in Eq. (11): \bar{p}_* is a piecewise constant function that averages the \tilde{p}_* over each quantization cell. Therefore, the total mass of each cell remains the same, i.e., $\bar{p}_*(C_k) = \tilde{p}_*(C_k)$ for any cell C_k . Therefore, we construct the coupling of (\bar{p}_*, \tilde{p}_*) by considering each cell C_k : within the cell C_k , choose the independent coupling γ_k that transports the mass distribution from $\tilde{p}_*(\cdot|C_k)$ to the uniform measure $\bar{p}_*(\cdot|C_k)$ only within C_k . Specifically, construct the coupling (X, Y) as follows:

1. Sample the cell index according to the PMF $\mathbb{P}(\text{index} = k) = \tilde{p}_*(C_k)$.
2. Conditioned on $K = k$, sample X from the conditional law $\bar{p}_*(\cdot|C_k)$; sample Y , independently of X , from the conditional law $\tilde{p}_*(\cdot|C_k)$.

Since $\bar{p}_*(C_k) = \tilde{p}_*(C_k)$, we have

$$\begin{aligned} p_X(x) &= \bar{p}_*(x | C_k) \tilde{p}_*(C_k) = \bar{p}_*(x | C_k) \bar{p}_*(C_k) = \bar{p}_*(x) \\ \text{and } p_Y(y) &= \tilde{p}_*(y | C_k) \tilde{p}_*(C_k) = \tilde{p}_*(y). \end{aligned}$$

Hence, (\mathbf{x}, \mathbf{y}) gives a coupling of (\bar{p}_*, \tilde{p}_*) . Conditioned on each cell C_k , we have

$$\|\mathbf{x} - \mathbf{y}\| \leq \sqrt{d} \cdot l, \quad \mathbf{x} \sim \tilde{p}_*(\cdot | C_k) \quad \text{and} \quad \mathbf{y} \sim \bar{p}_*(\cdot | C_k),$$

where l is the diameter of cell. Since all probability mass stays in its original cell, we get

$$\mathbb{E}[\|X - Y\|^2] \leq dl^2 \quad \Rightarrow \quad W_2^2(\tilde{p}_*, \bar{p}_*) \leq dl^2,$$

which implies $W_2(\tilde{p}_*, \bar{p}_*) \leq \epsilon$ by requiring $l = \epsilon/\sqrt{d}$. Such a choice is independent of the smoothness of p_* . Due to the triangle inequality, we have $W_2(p_*, \bar{p}_*) \leq 2\epsilon$.

Proof of the TV distance bound. Since our truncated uniformization, i.e., Alg. 2, exactly simulates the reversed process, from Lemma 4, the KL divergence gap between $q_{T-\delta}^{\leftarrow} = q_{\delta}^{\rightarrow}$ and $\hat{q}_{T-\delta}$ is bounded by the KL divergence as follows:

$$\begin{aligned} \text{KL}(q_{T-\delta}^{\leftarrow} \| \hat{q}_{T-\delta}) &\leq \text{KL}(q_0^{\leftarrow} \| \hat{q}_0) + (T - \delta) \epsilon_{\text{score}}^2 \\ &\leq e^{-T} \cdot d \log_2 K + (T - \delta) \epsilon_{\text{score}}^2 = \epsilon^2 + (\ln(d/\epsilon) + \ln \log_2 K)^2 \cdot \epsilon_{\text{score}}^2 \leq 2\epsilon^2, \end{aligned} \quad (54)$$

where the second inequality follows from Lemma 2, the third inequality establishes when T is chosen as

$$T = \ln(d/\epsilon) + \ln \log_2 K,$$

and the last inequality is established when we have

$$\epsilon_{\text{score}} = \frac{\epsilon}{\ln(d/\epsilon) + \ln \log_2 K} = \tilde{O}(\epsilon).$$

Under this condition, due to Pinsker's inequality, Eq. (54) can be relaxed to

$$\text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) \leq \sqrt{\frac{\text{KL}(q_{T-\delta}^{\leftarrow} \| \hat{q}_{T-\delta})}{2}} \leq \epsilon.$$

Then we have

$$\begin{aligned} \text{TV}(q_*, \hat{q}_{T-\delta}) &\leq \text{TV}(q_*, q_{T-\delta}^{\leftarrow}) + \text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) \\ &= \text{TV}(q_*, q_{\delta}^{\rightarrow}) + \text{TV}(q_{T-\delta}^{\leftarrow}, \hat{q}_{T-\delta}) = 1 - e^{-\delta d \log_2 K} + \epsilon \leq 2\epsilon \end{aligned}$$

where the second equation follows from Eq. (50) and the last inequality is established by requiring

$$\delta \leq \frac{\epsilon}{d \cdot \log_2 K} \quad \Leftrightarrow \quad \delta d \log_2 K \leq \epsilon.$$

Under this condition, we have

$$\delta d \log_2 K \leq \epsilon \leq \ln \frac{1}{1 - \epsilon} \quad \Rightarrow \quad 1 - e^{-\delta d \log_2 K} \leq \epsilon.$$

Suppose the underlying distributions of $\bar{\mathbf{y}}, \hat{\mathbf{x}}$ are \bar{q}, \hat{p} respectively, due to the connection between \hat{p}, \bar{p}_* and \bar{q}, \bar{q}_* shown in Eq. (12), we have

$$\text{TV}(\bar{p}_*, \hat{p}) = \int |\bar{p}_*(\mathbf{x}) - \hat{p}(\mathbf{x})| d\mathbf{x} = \sum_{\bar{\mathbf{y}} \in \bar{\mathcal{Y}}} |\bar{q}(\bar{\mathbf{y}}) - \bar{q}_*(\bar{\mathbf{y}})| = \sum_{\mathbf{y} \in \mathcal{Y}} |\hat{q}_{T-\delta}(\mathbf{y}) - \hat{q}_*(\mathbf{y})| \leq 2\epsilon.$$

Proof of bound of the expectation of \bar{N} , i.e., $\bar{\beta}$. After the training, the implementation of Alg. 2 requires $\bar{N} \sim \text{Poisson}(\bar{\beta})$ steps. According to Lemma 15, if we set

$$t_0 = 0, \quad t_{w+1} - t_w = 0.5 \cdot (T - t_{w+1}) \quad \text{and} \quad t_W = T - \delta,$$

for the time partitions,

$$\beta_{t_w} := 2d \log_2 K / \min\{1, T - t_w\}$$

for the intermediate Poisson, then it has

$$\begin{aligned} \bar{\beta} &= \sum_{k=1}^W \beta_{t_k} \cdot (t_k - t_{k-1}) \leq 2d \log_2 K \cdot (T + \ln(1/\delta)) \\ &= 2d \cdot \left[\log_2(4\sqrt{2d}\sigma \cdot \sqrt{\ln(2d\sigma/\epsilon)} \cdot \sqrt{d}/\epsilon) \right] \cdot (T + \ln(1/\delta)) \end{aligned}$$

where the last equation is due to $L = 2\sqrt{2d}\sigma \cdot \sqrt{\ln(2d\sigma/\epsilon)}$, $l = \epsilon/\sqrt{d}$, and $K = 2L/l$. Therefore, the proof is completed. \blacksquare

Appendix F. Technical Lemmas

Lemma 17 (Theorem 4.10 of Boucheron et al. (2003)) *Let $\Phi(x) = x \ln x$ for $x > 0$ and $\Phi(0) = 0$. Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be independent random variables taking values in a countable set \mathcal{X} and let $f: \mathcal{X} \rightarrow [0, \infty)$. We have*

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n} [\Phi(f(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n))] - \Phi(\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n} [f(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)]) \\ &\leq \sum_{i=1}^n \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n} [\mathbb{E}_{\mathbf{x}_i} [\Phi(f(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n))] - \Phi(\mathbb{E}_{\mathbf{x}_i} [f(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)])]. \end{aligned}$$

Lemma 18 (Chain rule of TV) *Consider four random variables, $\mathbf{x}, \mathbf{z}, \tilde{\mathbf{x}}, \tilde{\mathbf{z}}$, whose underlying distributions are denoted as p_x, p_z, q_x, q_z . Suppose $p_{x,z}$ and $q_{x,z}$ denotes the densities of joint distributions of (\mathbf{x}, \mathbf{z}) and $(\tilde{\mathbf{x}}, \tilde{\mathbf{z}})$, which we write in terms of the conditionals and marginals as*

$$\begin{aligned} p_{x,z}(\mathbf{x}, \mathbf{z}) &= p_{x|z}(\mathbf{x}|\mathbf{z}) \cdot p_z(\mathbf{z}) = p_{z|x}(\mathbf{z}|\mathbf{x}) \cdot p_x(\mathbf{x}) \\ q_{x,z}(\mathbf{x}, \mathbf{z}) &= q_{x|z}(\mathbf{x}|\mathbf{z}) \cdot q_z(\mathbf{z}) = q_{z|x}(\mathbf{z}|\mathbf{x}) \cdot q_x(\mathbf{x}). \end{aligned}$$

then we have

$$\begin{aligned} \text{TV}(p_{x,z}, q_{x,z}) &\leq \min \left\{ \text{TV}(p_z, q_z) + \mathbb{E}_{\mathbf{z} \sim p_z} [\text{TV}(p_{x|z}(\cdot|\mathbf{z}), q_{x|z}(\cdot|\mathbf{z}))], \right. \\ &\quad \left. \text{TV}(p_x, q_x) + \mathbb{E}_{\mathbf{x} \sim p_x} [\text{TV}(p_{z|x}(\cdot|\mathbf{x}), q_{z|x}(\cdot|\mathbf{x}))] \right\}. \end{aligned}$$

Besides, we have

$$\text{TV}(p_x, q_x) \leq \text{TV}(p_{x,z}, q_{x,z}).$$

Lemma 19 (Backward Kolmogorov equation) *Suppose the infinitesimal operator of a Markov semigroup is \mathcal{L} . If we denote the transition density from $\mathbf{y}_s = \mathbf{y}$ to $\mathbf{y}_t = \mathbf{y}'$ as $p_{t|s}(\mathbf{y}'|\mathbf{y})$, then it solves the backward Kolmogorov equation*

$$-\frac{\partial p_{t|s}(\mathbf{y}'|\mathbf{y})}{\partial s} = \mathcal{L} [p_{t|s}(\mathbf{y}'|\cdot)](\mathbf{y}), \quad p_{s|s}(\mathbf{y}'|\mathbf{y}) = \delta(\mathbf{y}' - \mathbf{y}).$$

Lemma 20 (Lemma 11 in Vempala and Wibisono (2019)) *Suppose the density function satisfies $p \propto \exp(-f)$ where f is H -smooth, i.e., [A3]. Then, it has*

$$\mathbb{E}_{\mathbf{x} \sim p} \left[\|\nabla f(\mathbf{x})\|^2 \right] \leq Hd.$$

Appendix G. Empirical Results

To evaluate the efficiency and accuracy of QTD relative to DDPM, we conduct a synthetic experiment on a 1D Gaussian Mixture Model (GMM) in Section G.1. On the other hand, although Table 2 indicates that the complexity of most diffusion inference methods scales with d , in practice the number of inference iterations is typically much smaller than the particle dimension. To demonstrate the practical effectiveness of QTD, we introduce a multi-flip trick that yields an inexact variant of QTD, allowing the number of iterations to be smaller than d . The implementation details and corresponding empirical results are presented in Section G.2.

G.1. Experiments on Synthetic Data

Data Setup. We define the target distribution as a 2-component GMM, discretized into $K = 128$ categories over the range $[-4, 4]$. For each of the 10 independent experimental seeds, we randomly sample the target parameters: component means $\mu_1 \sim \mathcal{U}[-1.5, -1.0]$ and $\mu_2 \sim \mathcal{U}[1.0, 1.5]$, standard deviations $\sigma_1, \sigma_2 \sim \mathcal{U}[0.4, 1.0]$, and mixing weights $w_1 \sim \mathcal{U}[0.3, 0.7]$ (with $w_2 = 1 - w_1$).

Experimental Setup and Architecture. We employ three-layer MLP backbones with comparable model capacity and ReLU activations to parameterize the discrete and continuous scores, respectively. QTD is trained by minimizing the discrete score entropy (Eq. 5), while DDPM is trained using the standard denoising diffusion objective. For the forward processes, we set the mixing time $T = 2.0$ for QTD and a 1000-step linear schedule for DDPM. We train both models for 20000 iterations using the Adam optimizer with a learning rate of 1×10^{-3} and a batch size of 256.

Inference and Evaluation. During inference, we draw 5000 independent samples from each model. To achieve a similar generation performance, DDPM is executed with 1000 denoising steps (constant step size $\Delta t = 10^{-3}$), while QTD utilizes an early stopping threshold of $\delta = 10^{-4}$. The performance is quantified using the Total Variation (TV) distance between the empirical distribution \hat{p} and the true discretized GMM p :

$$\text{TV}(\hat{p}, p) = \frac{1}{2} \sum_{k=1}^K |\hat{p}_k - p_k|.$$

The TV distance is calculated across the 5000 samples for each of the 10 target distribution seeds.

Results. As shown in Figure 4, QTD demonstrates significant computational advantages. While both methods eventually converge to similar TV distances, QTD achieves rapid stabilization within the first 100–150 Number of Function Evaluations (NFE)*. In contrast, DDPM requires the full 1000 denoising steps to reach a comparable level of accuracy. This behavior is further illustrated in

*. Since QTD sampling terminates adaptively, the NFE varies across trajectories. The notation mean max NFE in Figure 4 denotes the mean (over 10 seeds) of the maximum NFE (over 5000 trajectories per seed) as the computational budget for each run.

the histogram evolution (Figure 5), where the QTD empirical distribution matches the gray target GMM much earlier in the generative process than the DDPM baseline.

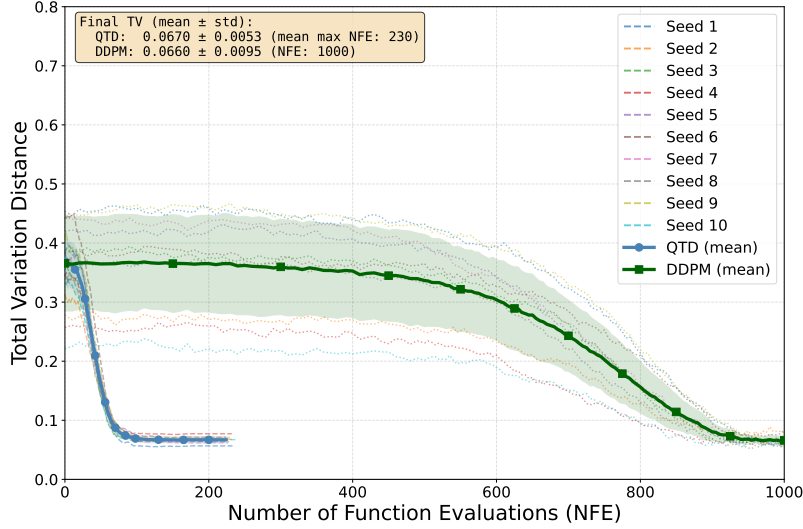


Figure 4: Convergence comparison between QTD and DDPM on a 1D Gaussian Mixture Model. We plot TV distance from the target distribution against NFE. Dashed lines show individual runs across 10 random seeds; solid lines show the mean with shaded standard deviation. QTD (blue) converges rapidly within the first 100 NFE and stabilizes around 150 NFE, whereas DDPM (green) requires the full 1000 denoising steps. Both methods achieve similar final TV distances, but QTD demonstrates significantly improved computational efficiency.

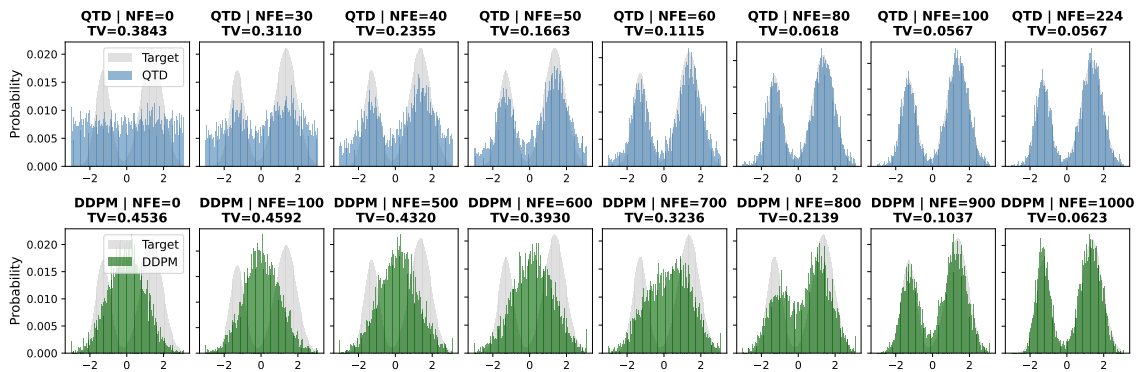


Figure 5: Histogram evolution comparison between QTD (first row) and DDPM (second row). Histograms show the empirical distribution (colored) versus the target Gaussian mixture model (gray) at various numbers of function evaluation (NFE) checkpoints. QTD achieves rapid convergence within around 100 NFE, while DDPM requires the full 1000 denoising steps to match the target distribution.

Algorithm 3 INFERENCE PROCESS WITH MULTIFLIP TRUNCATED UNIFORMIZATION

Require: Total time T , time partition $0 = t_0 < \dots < t_W = T - \delta$, parameters $\beta_{t_1}, \dots, \beta_{t_W}$ set as Eq. (14), reverse transition rate \hat{R}_t^{\leftarrow} from score $\tilde{v}_{t, \mathbf{y}'}(\cdot)$, Poisson factor λ , flip candidates M .

- 1: Draw $\hat{\mathbf{y}}_{t_0} \sim \text{Uniform}(\{0, 1\}^{d \log_2 K})$.
 - 2: **for** $w = 1$ **to** W **do**
 - 3: Draw $N \sim \text{Poisson}(\lambda \cdot \beta_{t_w} (t_w - t_{w-1}))$;
 - 4: Sample N points i.i.d. from $[t_{w-1}, t_w]$, sort as $\tau_1 < \tau_2 < \dots < \tau_N$;
 - 5: Set $\mathbf{z}_0 = \hat{\mathbf{y}}_{t_{w-1}}$;
 - 6: **for** $n = 1$ **to** N **do**
 - 7: Evaluate $\hat{R}_{\tau_n}(\mathbf{z}_{n-1} + \mathbf{e}_i, \mathbf{z}_{n-1})$ for all i using reference state \mathbf{z}_{n-1} ;
 - 8: Let $\hat{R}_{\tau_n}(\mathbf{z}_{n-1}) = \sum_i \hat{R}_{\tau_n}(\mathbf{z}_{n-1} + \mathbf{e}_i, \mathbf{z}_{n-1})$ and $p_n = \min\left(1, \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n}(\mathbf{z}_{n-1})\right)$;
 - 9: Sample M bit positions i.i.d. $i_1, \dots, i_M \sim \text{Cat}\left(\frac{\hat{R}_{\tau_n}(\mathbf{z}_{n-1} + \mathbf{e}_i, \mathbf{z}_{n-1})}{\hat{R}_{\tau_n}(\mathbf{z}_{n-1})}\right)$;
 - 10: Draw $A_1, \dots, A_M \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p_n)$;
 - 11: Collect accepted flip set $\mathcal{F}_n = \{i_m : A_m = 1, m = 1, \dots, M\}$;
 - 12: Set $\mathbf{z}_n = \mathbf{z}_{n-1} \oplus \mathbf{f}_n$, where $\mathbf{f}_n[j] = 1[j \in \mathcal{F}_n]$;
 - 13: **end for**
 - 14: Set $\hat{\mathbf{y}}_{t_w} = \mathbf{z}_N$;
 - 15: **end for**
 - 16: Recover $\bar{\mathbf{y}} = \text{vBin}^{-1}(\hat{\mathbf{y}}_{t_W})$ and draw $\hat{\mathbf{x}} \sim \text{Uniform}(\text{Cell}(\bar{\mathbf{y}}))$.
 - 17: **return** $\hat{\mathbf{x}}$.
-

G.2. Experiments on MNIST

Data Setup. We use the standard MNIST training split (60000 images) for training and the test split (10000 images) for FID evaluation. Each 28×28 grayscale image is first downsampled to 16×16 via bilinear interpolation and then keep $K = 256$ discrete gray levels using uniformly spaced thresholds over $[0, 1]$. For training data construction, QTD uses the quantized pixel values directly as discrete tokens, whereas DDPM re-normalizes them to $[-1, 1]$ as continuous inputs. During FID evaluation, all generated images are passed through the same quantization pipeline prior to feature extraction to ensure a fair comparison.

Experimental Setup and Architecture. Both models share the same UNet2DModel as the backbone with channel widths $[128, 256, 512]$, two residual blocks per level, a time-embedding dimension of 256, and a single self-attention layer at the penultimate encoder resolution (4×4). QTD parameterizes per-bit log-scores and is trained by minimizing the discrete score entropy loss, while DDPM is trained with the standard MSE ϵ -prediction objective. For the forward process, QTD uses a discrete uniform noising process with mixing time $T = 1.0$ and start time $t_0 = 0.0$; DDPM uses a 7200-step linear Gaussian noise schedule. Both models are trained for 10 epochs using the AdamW optimizer with a learning rate of 2×10^{-4} , a cosine warm-up of 500 steps, a batch size of 256, and gradient clipping at 1.0.

Inference and Evaluation. During inference, we draw 10,000 independent samples from each model. DDPM is executed with 1200, 2400, and 3600 denoising steps. QTD is evaluated under QTD-MultiFlip (Alg. 3), which uses multi-flip Truncated Uniformization with M simultaneous bit

flips and a Poisson acceleration factor of λ ; use an early-stopping threshold of δ . Image quality is quantified using the Fréchet Inception Distance (FID) between the empirical generated distribution \hat{p} and the real MNIST test distribution p :

$$\text{FID}(\hat{p}, p) = \|\mu_{\hat{p}} - \mu_p\|^2 + \text{Tr}\left(\Sigma_{\hat{p}} + \Sigma_p - 2(\Sigma_{\hat{p}}\Sigma_p)^{1/2}\right), \tag{55}$$

where μ and Σ are the mean and covariance of InceptionV3 pool₃ activations (dimension 2048) computed over 10000 images.

Results. Table 4 reports FID scores computed from 10000 generated samples at convergence. QTD-MultiFlip (Alg. 3) consistently achieves the lowest FID across all inference budgets, demonstrating that the multi-flip proposal accelerates convergence to the target distribution. Even with a modest number of inference steps, QTD matches or surpasses DDPM in generation quality while requiring significantly fewer network function evaluations, consistent with the computational advantages observed in the 1D synthetic setting.

Table 4: FID comparison on MNIST between DDPM and QTD (lower is better). All scores are computed from 10000 generated samples.

Inference Steps	FID (↓)	
	DDPM (Ho et al., 2020)	QTD (Alg. 3)
1200	4.358	4.147
2400	4.397	4.134
3600	4.330	4.081

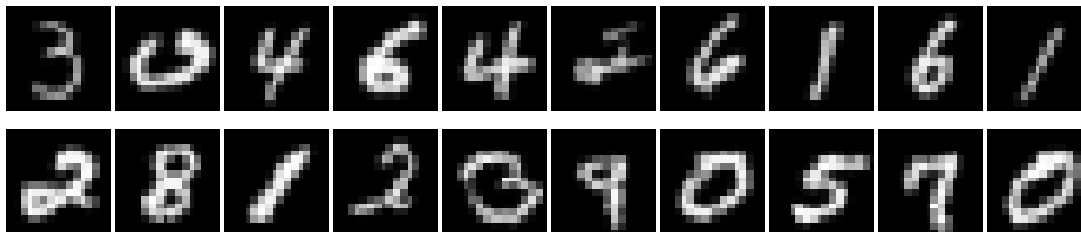


Figure 6: Random samples generated on MNIST. Top row: DDPM. Bottom row: QTD.