

ON THE COMPLEXITY THEORY OF MASKED DISCRETE DIFFUSION: FROM $\text{poly}(1/\epsilon)$ TO NEARLY ϵ -FREE

Anonymous authors

Paper under double-blind review

ABSTRACT

We study *masked discrete diffusion*—a flexible paradigm for text generation in which tokens are progressively corrupted by special mask symbols before being denoised. Although this approach has demonstrated strong empirical performance, its theoretical complexity in high-dimensional settings remains insufficiently understood. Existing analyses largely focus on *uniform* discrete diffusion, and more recent attempts addressing masked diffusion either (1) overlook widely used Euler samplers, (2) impose restrictive bounded-score assumptions, or (3) fail to showcase the advantages of masked discrete diffusion over its uniform counterpart. To address this gap, we show that Euler samplers can achieve ϵ -accuracy in total variation (TV) with $\tilde{O}(d^2\epsilon^{-3/2})$ discrete score evaluations, thereby providing the first rigorous analysis of typical Euler sampler in masked discrete diffusion. We then propose a *Mask-Aware Truncated Uniformization* (MATU) approach that both removes bounded-score assumptions and preserves unbiased discrete score approximation. By exploiting the property that each token can be unmasked at most once, MATU attains a nearly ϵ -free complexity of $O(d \ln d \cdot (1 - \epsilon^2))$. This result surpasses existing uniformization methods under uniform discrete diffusion, eliminating the $\ln(1/\epsilon)$ factor and substantially speeding up convergence. Our findings not only provide a rigorous theoretical foundation for masked discrete diffusion, showcasing its practical advantages over uniform diffusion for text generation, but also pave the way for future efforts to analyze diffusion-based language models developed under masking paradigm.

1 INTRODUCTION

Diffusion language models (Sohl-Dickstein et al., 2015; Hoogeboom et al.; Austin et al., 2021; Lou et al., 2024; Ou et al., 2024) have recently emerged as a powerful class of generative paradigms, frequently regarded as both complements and competitors to the auto-regressive based language models (Achiam et al., 2023; Touvron et al., 2023; Zhao et al., 2023). Whereas auto-regressive models learn the conditional distribution of the next token given a prefix, diffusion language models approximate the joint distribution of an entire token sequence through a noising–denoising process. This process transforms a potentially complex data distribution into a simpler prior distribution and then iteratively reconstructs it. In the forward (*noising*) direction, tokens are progressively replaced by special mask symbols, thereby mapping the data distribution to a one-hot stationary distribution. The reverse (*denoising*) direction then recovers the original text step by step by estimating discrete scores (i.e., density ratios) over the corrupted samples.

Although masked discrete diffusion has empirically outperformed uniform discrete diffusion (where the forward process admits a uniform stationary distribution) (Lou et al., 2024), analyzing and mitigating its computational overhead in high-dimensional settings remains challenging. As summarized in Table 3, most existing theoretical results focus on *uniform discrete diffusion*. In these analyses, Euler-type samplers approximate continuous-time scores by holding them constant over short intervals, leading to polynomial complexity in the total variation (TV) distance ϵ . Specifically, exponential-integrator methods (Zhang et al., 2024) require $\tilde{O}(\epsilon^{-2})$ steps, while τ -leaping methods (Campbell et al., 2022; Lou et al., 2024) and their higher-order variants (Ren et al., 2025) need $\tilde{O}(\epsilon^{-1})$ steps. Notably, uniformization-based techniques offer a promising approach, achieving $O(\ln(1/\epsilon))$ complexity by unbiasedly simulating the reverse Markov chain. In the context of

masked discrete diffusion, Liang et al. (2025a) rigorously examined ϵ -TV convergence, showing that τ -leaping can take $\tilde{O}(\epsilon^{-2})$ steps to converge and also improves upon the dimensional dependence found in uniform discrete diffusion. However, their stronger bounded-score assumptions make direct comparisons of algorithmic complexity with existing works (Chen & Ying, 2024; Huang et al., 2025) uncertain. Although uniformization can theoretically reach a complexity of $O(\ln(1/\epsilon))$ in their framework, it retains the same ϵ -dependence as uniform discrete diffusion and has yet to exhibit clear empirical benefits in masked diffusion. Finally, the analysis of the typical Euler sampler used in most empirical studies (Lou et al., 2024; Ou et al., 2024) is still not fully understood.

To address the theoretical challenges of masked discrete diffusion, we first analyze a typical Euler sampler that parallels the inference procedures used in many empirical studies (Lou et al., 2024; Ou et al., 2024). Our findings reveal that reaching ϵ -TV convergence in masked discrete diffusion with the typical Euler sampler requires $\tilde{O}(d^2\epsilon^{-3/2})$ discrete score evaluations. This result stands as the first rigorous analysis of the typical Euler method in masked discrete diffusion and demonstrates faster convergence than the τ -leaping approach (Liang et al., 2025a) under stringent accuracy demands. We then examine uniformization-based approaches for masked discrete diffusion, where uniformization converts a continuous-time Markov chain (CTMC) into a discrete-time Markov chain (DTMC) by sampling random Poisson jump times. This technique preserves the exact transition structure of the original CTMC and provides an unbiased simulation without time-step discretization error. To eliminate the bounded-score assumption used in previous uniformization analyses (Chen & Ying, 2024; Liang et al., 2025a), we propose a *Mask-Aware Truncated Uniformization* (MATU) method inspired by Huang et al. (2025). Under MATU, we rescale the outgoing transition rates of the reverse process according to the number of masked tokens in preceding states, naturally tightening enforcing boundedness in the discrete score estimator while preserving the unbiasedness of uniformization-based score approximation. We prove that MATU can reach the same ϵ -TV convergence at a nearly ϵ -free complexity, offering a significant speedup from $O(\ln(1/\epsilon))$ to $O(1 - \epsilon^2)$. The key insight is that uniformization in the masked setting explicitly identifies which tokens remain masked and require denoising, thereby avoiding the redundant denoising attempts that slow convergence in uniform discrete diffusion. Our main contributions are summarized as follows.

- We present the first rigorous theoretical analysis of typical Euler samplers for masked discrete diffusion. Achieving ϵ -TV convergence requires $\tilde{O}(d^2\epsilon^{-3/2})$ discrete score evaluations, surpassing τ -leaping (Liang et al., 2025a) in high-accuracy settings.
- We propose a new method called *Mask-Aware Truncated Uniformization* (MATU). Unlike simply applying uniformization to masked discrete diffusion (Liang et al., 2025a), our approach leverages a truncation on the outgoing rate, thereby removing the need for a score-bounded assumption. Moreover, our truncation is adaptive to the number of masked tokens, in contrast to Huang et al. (2025) which relies on a uniform constant, thus making full use of masked discrete diffusion properties.
- By leveraging the property that tokens cannot be unmasked multiple times, MATU significantly accelerates convergence on the discrete space $\{1, 2, \dots, K\}^d$. Specifically, to reach ϵ -TV convergence, MATU uses an expected number of discrete score calls on the order of

$$O(d \cdot (1 - \epsilon^2/d) + d \ln d).$$

Compared to uniformization-based sampler in uniform discrete diffusion (Huang et al., 2025; Liang et al., 2025a), this result improves upon the $O(\ln(1/\epsilon))$ rate and surpasses the linear convergence limitation. Moreover, the dependence on both vocabulary size K and dimension d aligns with state-of-the-art performance (Zhang et al., 2024).

2 PRELIMINARIES

In this section, we establish the notation and setup for both forward and reverse Markov processes in general discrete diffusion models. We discuss marginal and conditional distributions, the transition rate function, neural-network-parameterized discrete scores (density ratios), and a standard training objective. We also present the commonly adopted assumption on score estimation error, which underlies many theoretical and empirical works (Zhang et al., 2024; Lou et al., 2024; Chen & Ying, 2024; Huang et al., 2025; Liang et al., 2025a). A comprehensive summary of the notation can be found in Table 2 of Appendix A.

The forward process notations. In this paper, we consider discrete distributions over $\mathcal{Y} = \{1, 2, \dots, K\}^d$. For any functions $f, g : \mathcal{Y} \rightarrow \mathbb{R}$, we define their inner product as

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

Given a target distribution q_* , we define a forward Markov process $\{\mathbf{y}_t^{\rightarrow}\}_{t=0}^T$ with $q_0^{\rightarrow} = q_*$, which converges to a stationary distribution q_{∞}^{\rightarrow} as $T \rightarrow \infty$. We denote by q_t^{\rightarrow} its marginal at time t , and use $q_{t',t}^{\rightarrow}(\mathbf{y}', \mathbf{y})$ and $q_{t'|t}^{\rightarrow}(\mathbf{y}'|\mathbf{y})$ to represent the joint and conditional distributions over times t' and t , respectively:

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

Both masked and uniform discrete diffusion models treat this forward process as a time-homogeneous CTMC with transition rate function $R^{\rightarrow} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ which denotes the instantaneous transition rate from \mathbf{y}' to \mathbf{y} . Formally,

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

where $\delta_{\mathbf{y}'}(\mathbf{y}) = 1$ if $\mathbf{y} = \mathbf{y}'$ and 0 otherwise. We further define $R^{\rightarrow}(\mathbf{y}') := \sum_{\mathbf{y} \neq \mathbf{y}'} R^{\rightarrow}(\mathbf{y}, \mathbf{y}')$ as the outgoing rate, which denotes the instantaneous transition rate from \mathbf{y}' to all other feasible states. Under this condition, the discrete forward process follows

$$\frac{dq_t^{\rightarrow}}{dt}(\mathbf{y}|\mathbf{y}_0) = \left\langle R^{\rightarrow}(\mathbf{y}, \cdot), q_{t|s}^{\rightarrow}(\cdot|\mathbf{y}_0) \right\rangle_{\mathcal{Y}}, \quad \frac{dq_t^{\rightarrow}}{dt}(\mathbf{y}) = \langle R^{\rightarrow}(\mathbf{y}, \cdot), q_t^{\rightarrow}(\cdot) \rangle_{\mathcal{Y}}. \quad (2)$$

More details and derivation can be found in Appendix B.

The reverse process notations. To sample from $q_* = q_0^{\rightarrow}$, discrete diffusion models define a reverse process $\{\mathbf{y}_t^{\leftarrow}\}_{t=0}^T$ such that $\mathbf{y}_t^{\leftarrow} \sim q_t^{\leftarrow} = q_{T-t}^{\rightarrow}$ and $(\mathbf{y}_{t'}^{\leftarrow}, \mathbf{y}_t^{\leftarrow}) \sim q_{t',t}^{\leftarrow}$. By Lemma 1 (proof in Appendix B.2), this time-inhomogeneous Markov chain satisfies:

Lemma 1 (Adapted from Eqs. (3) and (4) of Huang et al. (2025)). *The probability mass function q_t^{\leftarrow} in the reverse process follows*

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

and the reverse transition function R_t^{\leftarrow} arises as the infinitesimal operator of the reverse process:

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

while the outgoing rate is $R_t^{\leftarrow}(\mathbf{y}') = \sum_{\mathbf{y} \neq \mathbf{y}'} R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}')$.

Under this formulation, the reverse transition rate R_t^{\leftarrow} depends on the forward transition rate R^{\rightarrow} as well as the *discrete score*, defined as the density ratio $q_t^{\leftarrow}(\mathbf{y}) / q_t^{\leftarrow}(\mathbf{y}')$. Since this ratio is generally intractable, it is approximated in practice by a neural network \tilde{v} :

$$\tilde{v}_{t,\mathbf{y}'}(\cdot) \approx v_{t,\mathbf{y}'}(\cdot) = q_t^{\leftarrow}(\cdot) / q_t^{\leftarrow}(\mathbf{y}'), \quad (5)$$

yielding an approximate reverse transition rate \tilde{R}_t^{\leftarrow} via Eq. (3). To train \tilde{v} , one typically uses the *score entropy* loss (Lou et al., 2024; Benton et al., 2024),

$$L_{\text{SE}}(\tilde{v}) = \frac{1}{T} \int_0^T \mathbb{E}_{\mathbf{y}_t \sim q_t^{\rightarrow}} \left[\sum_{\mathbf{y} \neq \mathbf{y}_t} R^{\rightarrow}(\mathbf{y}_t, \mathbf{y}) D_{\phi} (v_{T-t,\mathbf{y}_t}(\mathbf{y}) \| \tilde{v}_{T-t,\mathbf{y}_t}(\mathbf{y})) \right] dt, \quad (6)$$

where $D_{\phi}(\cdot \| \cdot)$ is the Bregman divergence associated with $\phi(c) = c \ln c$. As in continuous diffusion (Chen et al., 2023), practitioners often replace L_{SE} by *implicit* or *denoising score entropy* (Lou et al., 2024; Benton et al., 2024) for more tractable optimization but invariant minimum.

General Assumptions. To analyze both convergence properties and the computational effort required for achieving TV distance convergence in practical settings, we assume the score entropy loss will be upper-bounded. Formally:

[A1] Score approximation error. The discrete score \tilde{v}_t obtained from Eq. (6) is well-trained, and its estimation error is small enough so that $L_{SE}(\tilde{v}) \leq \epsilon_{\text{score}}^2$.

This assumption is standard in theoretical inference research (Chen & Ying, 2024; Zhang et al., 2024; Lou et al., 2024), where it is commonly presumed that the score can be trained arbitrarily well such that $\epsilon_{\text{score}} \leq \epsilon$ for any desired $\epsilon > 0$.

3 THE FORWARD PROCESS OF MASKED DISCRETE DIFFUSION

In this section, we instantiate the masked discrete diffusion from the framework outlined in Section 2. We then construct a family of auxiliary distributions that approach the ideal forward marginal distribution exponentially quickly as time progresses. This construction leverages the forward transition kernel of masked discrete diffusion for any $0 < s < t < T$, and can be used as an alternative to the reverse initialization proposed by Liang et al. (2025a).

Additional settings. Following Ou et al. (2024), we adopt a diffusion-based language modeling framework. Our vocabulary is $\{1, 2, \dots, K\}$, where K denotes the mask token. We aim to generate a length- d sequence (sentence) $\mathbf{y} \in \mathcal{Y} = \{1, 2, \dots, K\}^d$. The number of mask tokens in specific sentence \mathbf{y} and the Hamming distance between two sentences (\mathbf{y} and \mathbf{y}') are denoted as

$$\text{numK}(\mathbf{y}) := \sum_{i=1}^d \delta_K(\mathbf{y}_i) \quad \text{and} \quad \text{Ham}(\mathbf{y}, \mathbf{y}') = d - \sum_{i=1}^d \delta_{\mathbf{y}_i}(\mathbf{y}'_i)$$

respectively. Generally, we suppose the mask token is never observed in target distribution:

[A2] No mask in the target distribution. The target distribution $q_0^\rightarrow = q_*: \mathcal{Y} \rightarrow \mathbb{R}$ assigns positive probability only to those sequences without any mask tokens, i.e. $q_*(\mathbf{y}) > 0$ if and only if $\text{numK}(\mathbf{y}) = 0$.

Masked discrete diffusion instantiation and approximation. We begin by specifying the absorbing forward transition rate function for masked discrete diffusion:

$$R^\rightarrow(\mathbf{y}, \mathbf{y}') = \begin{cases} 1 & \text{if } \text{Ham}(\mathbf{y}, \mathbf{y}') = 1 \text{ and } \mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')} = K \\ -\sum_{i=1}^d [1 - \delta_K(\mathbf{y}_i)] & \text{if } \mathbf{y} = \mathbf{y}' \\ 0 & \text{otherwise} \end{cases}. \quad (7)$$

Here, $\text{DiffIdx}(\mathbf{y}, \mathbf{y}')$ denotes the single coordinate where \mathbf{y} and \mathbf{y}' differ. Under this transition rule, each non-masked coordinate tends to become masked at an exponential rate. Concretely, for any $0 < s < t < T$, the forward transition kernel satisfies

$$q_{t|s}^\rightarrow(\mathbf{y}|\mathbf{y}') = \prod_{i=1}^d \left[\delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i) + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_0(\mathbf{y}_i - \mathbf{y}'_i) \cdot e^{-(t-s)} + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_K(\mathbf{y}_i) \cdot (1 - e^{-(t-s)}) \right], \quad (8)$$

as shown in Lemma 8. To approximate the forward marginal distribution q_t^\rightarrow at time t , we exploit this exponential decay by modeling each non-mask coordinate under a uniform distribution and masking coordinates at a constant rate. Specifically, we define

$$\tilde{q}_t(\mathbf{y}) \propto \prod_{i=1}^d \exp(-t \cdot [1 - \delta_K(\mathbf{y}_i)]) = \exp(-t \cdot [d - \text{numK}(\mathbf{y})]). \quad (9)$$

so that \tilde{q}_t factorizes over coordinates and is straightforward to sample from. Moreover, as established in Lemma 2, the KL divergence between q_t^\rightarrow and \tilde{q}_t decreases exponentially with t .

Lemma 2 (Exponentially decreasing KL divergence between q_t^\rightarrow and \tilde{q}_t). *Suppose the CTMC $\{\mathbf{y}_t^\rightarrow\}_{t=0}^T$ has transition rates R^\rightarrow from Eq. (7), with $\mathbf{y}_t^\rightarrow \sim q_t^\rightarrow$. Let \tilde{q}_t be the approximation of q_t^\rightarrow defined by Eq. (9). Then,*

$$\text{KL}(q_t^\rightarrow \parallel \tilde{q}_t) \leq (1 + e^{-t})^d - 1.$$

Consequently, to ensure $\text{KL}(q_t^\rightarrow \parallel \tilde{q}_t) \leq \epsilon$, it suffices to choose $t \geq \ln(4d/\epsilon)$.

From Lemma 2, the running time T required for \tilde{q}_T to approximate q_T^\rightarrow falls on the order of $\mathcal{O}(\ln(d/\epsilon))$. It precisely matches the forward mixing time for uniform discrete diffusion (Chen & Ying, 2024; Zhang et al., 2024; Huang et al., 2025) and continuous diffusion (Chen et al., 2023) converging to their stationary distributions. Although the final results exhibit a similar convergence rate, the underlying analytical techniques differ substantially because the one-hot stationary distribution of masked discrete diffusion does not satisfy the modified log-Sobolev condition. Further technical details are deferred to Appendix B.3.

4 EULER SAMPLER IN MASKED DISCRETE DIFFUSION

This section first introduces the Euler sampler in masked discrete diffusion, widely used for its parallel coordinate updates when reverse transition can be factorized coordinate-wise. We then extend it to handle more general reverse marginals with unknown correlations, and show how to control accumulative errors by introducing the exponential integrator as the auxiliary process. Finally, we provide convergence and complexity guarantees for achieving ϵ -TV convergence.

Typical Euler samplers and their extensions. Euler-type samplers have become increasingly popular in empirical studies (Lou et al., 2024; Ou et al., 2024) because their parallel-friendly updates often run faster than traditional auto-regressive models. Let $\{\hat{\mathbf{y}}_t\}_{t=0}^T$ denote the practical reverse process, whose marginal, joint, and conditional distributions satisfy:

$$\hat{\mathbf{y}} \sim \hat{q}_t, \quad (\hat{\mathbf{y}}_{t'}, \hat{\mathbf{y}}_t) \sim \hat{q}_{t',t}, \quad \text{and} \quad \hat{q}_{t'|t}(\mathbf{y}'|\mathbf{y}) = \hat{q}_{t',t}(\mathbf{y}', \mathbf{y})/\hat{q}_t(\mathbf{y}) \quad \text{where } t' \geq t.$$

A key assumption is that the reverse transition for each coordinate is conditionally independent:

$$\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) \propto \prod_{i=1}^d \hat{q}_{t+\Delta t|t}^{(i)}(\mathbf{y}[\{i\} \rightarrow \{\mathbf{y}'_i\}]|\mathbf{y}), \quad (10)$$

where the token revision function

$$\mathbf{y}[S: \rightarrow Y' \subseteq \mathcal{Y}^{|S|}] = \sum_{i=1}^d e_i \cdot \mathbf{1}[i \notin S] \cdot \mathbf{y}_i + \sum_{j=1}^{|S|} e_{s_j} \cdot Y'_j$$

indicates that the coordinates of \mathbf{y} indexed by the set S are replaced by the corresponding values in Y' . Then, each non-masked token can be updated independently in the reverse-time direction. Specifically, by discretizing Eq. (4) from Lemma 1, the update for the i th coordinate takes the form:

$$\hat{q}_{t+h|t}^{(i)}(\mathbf{y}[\{i\} \rightarrow \{\mathbf{y}'_i\}]|\mathbf{y}) = \delta_{\mathbf{y}_i}(\mathbf{y}'_i) + h \cdot R^\rightarrow(\mathbf{y}, \mathbf{y}[\{i\} \rightarrow \{\mathbf{y}'_i\}]) \cdot \tilde{v}_{t,\mathbf{y}}(\mathbf{y}[\{i\} \rightarrow \{\mathbf{y}'_i\}]).$$

Since $\text{Ham}(\mathbf{y}, \mathbf{y}[\{i\} \rightarrow \mathbf{y}'_i]) = 1$, the definition of R^\rightarrow in Eq. (7) ensures that $R^\rightarrow(\mathbf{y}, \mathbf{y}[\{i\} \rightarrow \mathbf{y}'_i]) \neq 0$. Hence, $\hat{q}_{t+h|t}^{(i)}(\mathbf{y}[\{i\} \rightarrow k]|\mathbf{y})$ for any non-mask token $k \neq K$, enabling all coordinates to be updated in parallel.

However, if the assumption in Eq. (10) does not hold, parallel updates become invalid. A practical alternative is to discretize Eq. (4) jointly, leading to the sequential update:

$$\hat{q}_{t+h|t}(\mathbf{y}'|\mathbf{y}) \propto \delta_{\mathbf{y}}(\mathbf{y}') + h \cdot \tilde{R}_t(\mathbf{y}', \mathbf{y}) = \delta_{\mathbf{y}}(\mathbf{y}') + h \cdot R^\rightarrow(\mathbf{y}, \mathbf{y}') \cdot \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') \quad (11)$$

where $\hat{q}_{t+\Delta t|t}(\mathbf{y}'|\mathbf{y}) \neq 0$ only if $R^\rightarrow(\mathbf{y}, \mathbf{y}') \neq 0$, which implies $\text{Ham}(\mathbf{y}, \mathbf{y}') = 1$ (see Eq. (7)). Consequently, at most one masked token could be denoised per update. In the subsequent analysis, we consider the Euler sampler using Eq. (11) in this more general setting.

Theoretical results. For the Euler sampler, the construction of the training loss, e.g., *denoising score entropy*, will be related to the step size h and share the same minimum with

$$L_{\text{DisSE}}(\tilde{v}) := \frac{1}{T - \delta} \sum_{k=0}^{n-1} \int_{kh}^{(k+1)h} \mathbb{E}_{\mathbf{y}_t \sim q_t^{\leftarrow}} \left[\sum_{\mathbf{y} \neq \mathbf{y}_t} R^{\rightarrow}(\mathbf{y}_t, \mathbf{y}) D_{\phi}(v_{kh, \mathbf{y}_t}(\mathbf{y}) || \tilde{v}_{kh, \mathbf{y}_t}(\mathbf{y})) \right] dt.$$

Correspondingly, to suppose the neural score estimator well approximates the discrete score only requires the following score estimation assumption, milder than Assumption [A1], i.e.,

[A1]- Score approximation error. The discrete score \tilde{v}_t obtained from Eq. (6) is well-trained, and its estimation error is small enough so that $L_{\text{DisSE}}(\tilde{v}) \leq \epsilon_{\text{score}}^2$.

Then, we summarize the convergence and complexity of Euler sampler (with proof in Section C.1).

Theorem 1. Suppose Assumption [A1]-, [A2] and Assumption 2 of Liang et al. (2025a) hold, implement Euler sampler with Eq. (11), if we require

$$T = \ln(4d/\epsilon^2), \quad h \lesssim \min \left\{ \frac{\epsilon}{K^2 d^2 \log(d/\epsilon)}, \frac{\epsilon^{\frac{3}{2}}}{d \sqrt{\log(d/\epsilon)}} \right\}, \quad \text{and} \quad \epsilon_{\text{score}} \leq \tilde{o}(\epsilon^2/d),$$

the Euler sampler will achieve $\text{TV}(p_*, \hat{p}) \leq 2\epsilon$ by requiring iterations to at an $\tilde{O}(d^2 \epsilon^{-3/2})$ level.

Compared to the τ -leaping method analyzed in Liang et al. (2025a), Euler-based approaches can be more effective in high-accuracy settings (e.g., $\epsilon \leq d^{-2}$). However, establishing a clear advantage over uniform discrete diffusion remains challenging. Due to time-discretization errors in discrete score estimation, Euler-based inference incurs polynomial complexity in both the dimensionality d and the error tolerance ϵ , which is still be worse than that in uniformization-based samplers.

5 TRUNCATED UNIFORMIZATION IN MASKED DISCRETE DIFFUSION

This section extends the truncated uniformization sampler of Huang et al. (2025) to masked discrete diffusion. We first revisit the core principle of unbiased reverse process simulation via uniformization. Next, we show that the expected complexity of uniformization-based inference depends critically on the outgoing rates of the reverse transition, and that masked discrete diffusion naturally offers smaller outgoing rates than its uniform counterpart, leading to faster convergence. We then introduce *Mask-Aware Truncated Uniformization* (MATU), which rescales the outgoing rates to eliminate the bounded-score assumption while preserving unbiased reverse process simulation. Finally, we provide theoretical results on MATU’s convergence and computational complexity, and compare these findings with existing approaches in the literature.

Uniformization and the expected number of discrete score calls. Consider a time-dependent reverse transition rate R_t^{\leftarrow} defined over the interval $[a, b]$. The evolution of the ideal reverse process for any \mathbf{y}, \mathbf{y}' can be described by

$$q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y}) = \begin{cases} \Delta t \cdot R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}), & \mathbf{y}' \neq \mathbf{y}, \\ 1 - \Delta t \cdot R_t^{\leftarrow}(\mathbf{y}), & \mathbf{y}' = \mathbf{y}, \end{cases} \quad \text{as } \Delta t \rightarrow 0, \quad (12)$$

following Eq. (4). If the total outgoing rate—denoting the instantaneous transition rate from \mathbf{y} to all other feasible states—is uniformly bounded by some β , i.e.,

$$R_t^{\leftarrow}(\mathbf{y}) = \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \leq \beta_t \leq \max_{t \in [a, b]} \beta_t = \beta, \quad (13)$$

then with probability $1 - \Delta t \cdot \beta$, the particle remains in the same state in each infinitesimal time step, thus requiring no additional score computation.

Based on this observation, the standard *uniformization* method (van Dijk, 1992; van Dijk et al., 2018; Chen & Ying, 2024) simulates the reverse dynamics over $[a, b]$ by iterating the following two-step procedure in the limit $\Delta t \rightarrow 0$:

1. Sample whether a transition occurs with probability $\Delta t \cdot \beta$.
2. If a transition occurs, move $\mathbf{y}_t^{\leftarrow}$ from \mathbf{y} to \mathbf{y}' with probability

$$M_t(\mathbf{y}' | \mathbf{y}) = \begin{cases} \beta^{-1} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}), & \mathbf{y}' \neq \mathbf{y}, \\ 1 - \beta^{-1} R_t^{\leftarrow}(\mathbf{y}), & \text{otherwise.} \end{cases} \quad (14)$$

Under this update scheme, the reverse transitions of uniformization will be equivalent to Eq. (12) exactly and introduce no time-discretization error (see Appendix D.2 for details). Moreover, since the number of transitions (and hence the number of discrete score computations) over $[a, b]$ follows a Poisson distribution with mean $\beta \cdot (b - a)$, any tighter bound on $R_t^{\leftarrow}(\mathbf{y})$ reduces β and thereby lowers the expected inference complexity.

The comparison of computational complexity and outgoing rate. By the previous discussion of uniformization, the expected number of discrete score calls over the time interval $[0, T]$ can be approximated by

$$\sum_{w=1}^W \max_{t \in [t_{w-1}, t_w]} \beta_t \cdot (t_w - t_{w-1}) \stackrel{W \rightarrow \infty}{\approx} \int_{t=0}^T \beta_t dt, \quad (15)$$

where $[t_0, t_1, \dots, t_W]$ is a partition of $[0, T]$. In uniform discrete diffusion, Chen & Ying (2024); Huang et al. (2025) show that the ideal reverse process satisfies

$$\beta_t := 2K \cdot d \cdot \max\{1, (T - t)^{-1}\} \leq \beta := 2K \cdot d \cdot \max\{1, (T - b)^{-1}\} \quad \forall t \in [a, b], \quad (16)$$

providing a uniform upper bound on the total outgoing rate $R_t^{\leftarrow}(\mathbf{y})$.

For *masked* discrete diffusion, Lemma 3 (with proof in Appendix D.1) shows that the outgoing rate can be bounded instead by

Lemma 3 (Bound of the outgoing rate). *Consider a CTMC whose transition rate function R^{\rightarrow} is defined as Eq. (7). Then, for any \mathbf{y} , the reverse transition rate function satisfies*

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) = R_t^{\leftarrow}(\mathbf{y}) \leq \beta_t(\mathbf{y}) := \frac{\text{numK}(\mathbf{y}) \cdot K}{e^{(T-t)} - 1}. \quad (17)$$

Compared to (16), this bound explicitly depends on $\text{numK}(\mathbf{y})$, the number of mask tokens in \mathbf{y} . Since $\text{numK}(\mathbf{y}) \leq d$, it is strictly smaller than the uniform bound in (16). Furthermore, $\text{numK}(\mathbf{y})$ decreases monotonically as the reverse process proceeds, which progressively enlarges the gap in outgoing rate between masked and uniform discrete diffusion. Because a lower outgoing rate implies fewer expected discrete score evaluations for each time t , masked discrete diffusion can be significantly more computationally efficient.

From an empirical perspective, a central observation is: *during inference, masked discrete diffusion only updates (denoises) masked tokens, whereas uniform discrete diffusion attempts to re-denoise tokens that have already been denoised.* Hence, in masked discrete diffusion, particles are more likely to remain unchanged at each step, leading to a smaller outgoing rate (and thus smaller β_t) over $[0, T]$. Consequently, fewer discrete score evaluations are required, underscoring the computational advantages of masked compared to uniform discrete diffusion.

Mask-aware truncation and algorithm proposal. In practice, we approximate the reverse transition rate $R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})$ by a learned neural score $\tilde{v}_{t,\mathbf{y}}(\mathbf{y}')$, yielding

$$\tilde{R}_t(\mathbf{y}', \mathbf{y}) = R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \tilde{v}_{t,\mathbf{y}}(\mathbf{y}'),$$

as dictated by Lemma 1 and Eq. (5). Because \tilde{v} is a learned estimator, the outgoing rate $\tilde{R}_t(\mathbf{y})$ may have no explicit upper bounds, complicating control over the expected number of discrete score evaluations. To mitigate unbounded transition rates, prior work typically imposes a bounded-score assumption on $\tilde{R}_t(\mathbf{y})$, restricting it to remain below a fixed constant (Liang et al., 2025a) or to grow as a function of the inference time (Chen & Ying, 2024). However, such assumptions can severely impact inference efficiency because the chosen upper bound β directly governs Step 2 of uniformization, as described in Eq. (14). When β is unknown, it can be treated as a hyperparameter.

Algorithm 1 MASK-AWARE TRUNCATED UNIFORMIZATION (MATU)

```

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. (17), 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  $\hat{\mathbf{y}}_{t_0} = [K, K, \dots, K]$ .
3: for  $w = 1$  to  $W$  do
4:   Choose  $\beta_{t_w} = K \cdot \text{numK}(\hat{\mathbf{y}}_{t_{w-1}}) / (e^{T-t_w} - 1)$ 
5:   Draw  $N \sim \text{Poisson}(\beta_{t_w}(t_w - t_{w-1}))$ ;
6:   Sample  $N$  points i.i.d. uniformly from  $[t_{w-1}, t_w]$  and sort them as  $\tau_1 < \tau_2 < \dots < \tau_N$ ;
7:   Set  $\mathbf{z}_0 = \hat{\mathbf{y}}_{t_{w-1}}$ ;
8:   for  $n = 1$  to  $N$  do
9:     Find the index set  $\mathcal{M}$  of [MASK] token appeared in random vector  $\mathbf{z}_{n-1}$ 
10:    For any  $i \in \mathcal{M}$  and  $k \in \{1, 2, \dots, K-1\}$ , update  $\mathbf{z}_{n-1}$  with
        
$$\mathbf{z}_n = \begin{cases} \mathbf{z}_{n-1}[\mathbf{z}_i: K \rightarrow k] & w.p. \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n, \mathbf{z}_0}(\mathbf{z}_{n-1}[\mathbf{z}_i: K \rightarrow k], \mathbf{z}_{n-1}), \\ \mathbf{z}_{n-1}, & w.p. 1 - \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n, \mathbf{z}_0}(\mathbf{z}_{n-1}). \end{cases}$$

11:   end for
12:   Set  $\hat{\mathbf{y}}_{t_w} = \mathbf{z}_N$ .
13: end for
14: return  $\hat{\mathbf{y}}_{t_W}$ .

```

Setting β too small may yield an infeasible probability $1 - \beta^{-1} \hat{R}_t(\mathbf{y}) < 0$, forcing the algorithm to fail; setting it too large preserves feasibility but inflates complexity in direct proportion to β . Thus, tightening this bounding scheme is crucial for balancing both correctness and computational efficiency in uniformization-based inference.

Motivated by Huang et al. (2025), we propose a *mask-aware truncation* scheme to rescale the practical outgoing rate $\tilde{R}_t(\mathbf{y}', \mathbf{y})$. This ensures that the non time-discretization property is preserved without additional cost, even when $\tilde{R}_t(\mathbf{y})$ becomes large. Specifically, consider simulating the reverse process over the (w -th) time segment $[t_{w-1}, t_w]$, assuming the state at time t_{w-1} is $\hat{\mathbf{y}}_{t_{w-1}} = \mathbf{y}_{t_{w-1}}$. Following from the monotonicity of $(e^{T-t} - 1)^{-1}$ and $\text{numK}(\hat{\mathbf{y}}_t)$ in Lemma 3, the mask-aware truncation is chosen as $\beta_{t_w}(\mathbf{y}_{t_{w-1}})$, then we set

$$\hat{R}_{t, \mathbf{y}_{t_{w-1}}}(\mathbf{y}, \mathbf{y}') = \begin{cases} \tilde{R}_t(\mathbf{y}, \mathbf{y}') \beta_{t_w}(\mathbf{y}_{t_{w-1}}) / \tilde{R}_t(\mathbf{y}'), & \text{if } \tilde{R}_t(\mathbf{y}') > \beta_{t_w}(\mathbf{y}_{t_{w-1}}), \\ \tilde{R}_t(\mathbf{y}, \mathbf{y}'), & \text{otherwise,} \end{cases} \quad \forall \mathbf{y}' \neq \mathbf{y}, \quad (18)$$

and

$$\hat{R}_{t, \mathbf{y}_{t_{w-1}}}(\mathbf{y}', \mathbf{y}') = - \sum_{\mathbf{y} \neq \mathbf{y}'} \hat{R}_{t, \mathbf{y}_{t_{w-1}}}(\mathbf{y}, \mathbf{y}'). \quad (19)$$

With these truncations, the corrected outgoing rate will be definitely upper bounded by $\beta_{t_w}(\mathbf{y}_{t_{w-1}})$. Then, we obtain a practical and efficient inference algorithm, summarized in Alg. 1.

Theoretical results. We summarize the convergence and complexity of Algorithm 1 for approximating q_* in Theorem 2 (proved in Appendices D.2 and D.3).

Theorem 2 (Combination of Theorem 3 and Theorem 4). *Suppose Assumption [A1] and [A2] hold, for Alg. 1, if we require*

$$T = \ln(4d/\epsilon^2), \quad \delta \leq d^{-1}\epsilon, \quad \epsilon_{\text{score}} \leq T^{-1/2}\epsilon, \quad \epsilon < 1,$$

and the partition of the reverse process satisfies

$$\eta = \epsilon/2d, \quad W = (T - \delta)/\eta, \quad t_0 = 0, \quad t_W = T - \delta, \quad t_w - t_{w-1} = \eta \quad \forall w \in \{1, 2, \dots, W\}$$

the expectation of iteration/score estimation complexity of Alg. 1 will be upper bounded by

$$2K(d - \epsilon^2/4) + 12Kd \ln d \quad (20)$$

to achieve $\text{TV}(p_, \hat{p}) \leq 2\epsilon$ where \hat{p} denotes the underlying distribution of generated samples.*

Table 1: Comparison with prior works simulating reverse particle SDEs, where **[A3]** denotes the bounded-score assumption used in Chen & Ying (2024) and **[A3]+** denotes the bounded-score assumption used in Liang et al. (2025a) which is a little bit stronger than **[A3]** due to the time-invariant requirement. All complexities are on TV convergence (or TV convergence deduced from KL convergence via Pinsker’s inequality, e.g., Ren et al. (2024)), which are achieved by assuming $\epsilon_{\text{score}} = \tilde{o}(\epsilon)$ and setting early-stopping parameters $\delta = \epsilon/d$. Besides, the complexity presented by $\tilde{O}(\cdot)$ means the \ln dependencies are omitted.

Results	Forward Type	Inference Sampler	Assumptions	Complexity
Zhang et al. (2024)	Uniformed	Exponential Integrator	[A1], [A3]	$\tilde{O}(d^{5/3}\epsilon^{-2})$
Ren et al. (2024)	Uniformed	τ -leaping	[A1],[A3]	$\tilde{O}(d^2\epsilon^{-2})$
Chen & Ying (2024)	Uniformed	Uniformization	[A1],[A3]	$O(d\ln(d/\epsilon))$
Huang et al. (2025)	Uniformed	Truncated Uniformization	[A1]	$O(d\ln(d/\epsilon))$
Theorem 1	Masked	Typical Euler	[A1],[A2],[A3]+	$\tilde{O}(d^2\epsilon^{-3/2})$
Liang et al. (2025a)	Masked	τ -leaping	[A1],[A2],[A3]+	$O(d\epsilon^{-2})$
Liang et al. (2025a)	Masked	Uniformization	[A1],[A2],[A3]	$O(d\ln(d/\epsilon))$
Theorem 2	Masked	MATU	[A1],[A2]	$O(d\ln d)$

From the above theorem, Eq. (20) might appear to enable exact inference by setting $\epsilon = 0$. However, this would require infinite mixing time T , perfect score estimates ($\epsilon_{\text{score}} = 0$), and infinitely many intervals W , which is infeasible. Meanwhile, although each interval has length $\eta = \epsilon/(2d)$ —leading to $\text{poly}(d/\epsilon)$ intervals in the reverse process—the total discrete score calls remain nearly independent of ϵ , since many intervals involve no state transitions (see Eq. (15)). Thus, small intervals are used primarily to match the accurate outgoing rate upper bound, without inflating complexity.

Then, We provide a complexity comparison in Table 3. MATU achieves a SOTA for both the ϵ -free complexity and the assumption without bounded-score estimator. Compared with existing uniformization-based method, Alg 1 achieves an $O(\ln(1/\epsilon))$ speedup, primarily because each token is denoised at most once in masked diffusion, whereas uniform diffusion renoises tokens multiple times. Formally, masked diffusion leverages the monotonic decrease of masked tokens, which cancels the growing outgoing rate:

$$\begin{aligned} \mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\mathbf{y}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] &\approx \sum_{w=1}^W \mathbb{E}[\text{numK}(\mathbf{y}_{T-t_{w-1}}^{\rightarrow})] \cdot K \cdot \frac{e^{-(T-t_w)}}{1 - e^{-(T-t_w)}} \cdot \eta \\ &= \sum_{w=1}^W d \cdot \underbrace{(1 - e^{-(T-t_{w-1})})}_{\text{decreasing factor}} \cdot K \cdot \underbrace{(1 - e^{-(T-t_w)})^{-1}}_{\text{increasing factor}} \cdot e^{-(T-t_w)} \cdot \eta \leq CKd \cdot \sum_{w=1}^W e^{-(T-t_w)} \cdot \eta, \end{aligned}$$

where the factor $e^{-(T-t_w)}$ keeps complexity low. In uniform diffusion, the same factor remains but grows with $1/(T - t_w)$, leading to a higher order overall:

$$\mathbb{E} \left[\sum_{w=1}^W \beta_{t_w} \cdot (t_w - t_{w-1}) \right] \lesssim CKd \cdot \sum_{w=1}^W \max\{1, (T - t_w)^{-1}\} \cdot \eta.$$

Since the integral $\int(1/t) dt$ diverges more quickly than $\int e^{-t} dt$, masked diffusion achieves lower inference complexity than uniform diffusion.

6 RELATED WORK

Most recently, the impressive empirical performance of discrete diffusion models (DDMs) has sparked a proliferation of theoretical investigations aiming to elucidate DDMs from various perspectives.

The Sample Complexity. For example, Srikanth et al. (2025) develops a theoretical framework for discrete-state diffusion models and presents the first rigorous sample-complexity bound of $\tilde{O}(\epsilon^{-2})$ under practical assumptions about neural network training. By pursuing a structured error decomposition, the authors illustrate how approximation, statistical, optimization, and clipping constraints

jointly contribute to the total complexity, furnishing dimension-free insights for training discrete-state diffusion models. Meanwhile, Wan et al. (2025) conducts the first non-asymptotic error analysis for discrete flow models on finite state spaces. By proposing a novel Girsanov-type theorem and bounding the KL divergence between two continuous-time Markov chains (CTMCs) with distinct transition rates, they rigorously decompose the transition-rate estimation error (including stochastic, approximation, and early-stopping components). Employing uniformization for sampling, the authors derive an upper bound on the distribution error that avoids any additional discretization error, thereby advancing the theory of discrete flow models beyond existing analyses of discrete diffusion.

The Inference Complexity. In addition to quantifying error tolerance and dimensional dependencies, Liang et al. (2025b) introduces a differential-inequality-based analysis for discrete diffusion models that eliminates the strong regularity assumptions required by Girsanov-based methods, reducing the convergence rates for τ -leaping from quadratic to linear in vocabulary size. Furthermore, Zheng et al. (2024) proposes the first-hitting-sampler (FHS) as a way to exactly simulate the reverse process by analytically sampling both the transition time and position. However, when discrete scores are parameterized by a time-dependent neural network (see Eq. 5), the uniform procedure for selecting the next unmasking position can introduce inference errors beyond those stemming from score estimation alone.

The key issue is that, although each masked position may share the same unmasking probability under the ideal reverse transition q_t^{\leftarrow} , this property may fail once the reverse process is learned. In particular, there can exist $i \neq j$ such that

$$\sum_{y', s.t. \text{ Ham}(y', y)=1, \text{ DiffIdx}(y', y)=i, y_i=K} s_{\theta, t, y}(y') \neq \sum_{y', s.t. \text{ Ham}(y', y)=1, \text{ DiffIdx}(y', y)=j, y_j=K} s_{\theta, t, y}(y').$$

so that uniformly choosing the next position to unmask biases the simulation of the learned reverse process, causing additional inference errors. Although this bias vanishes for time-independent discrete parameterizations (Ou et al., 2024), such as in Devlin et al. (2019); Chang et al. (2022); Ghazvininejad et al. (2019), where

$$q_t^{\leftarrow}(y')/q_t^{\leftarrow}(y) = \frac{e^{-t}}{1 - e^{-t}} \cdot q_0^{\rightarrow}(y'_i || y^{UM}) \approx \frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta, y}(y'), \quad (21)$$

the strong constraints, i.e.,

$$\sum_{y', s.t. \text{ Ham}(y', y)=1, \text{ DiffIdx}(y', y)=i, y_i=K} p_{\theta, y}(y') = 1 = \sum q_0(y'_i || y^{UM}),$$

ensure that every position has identical transition rates. Nevertheless, FHS Zheng et al. (2024) provides no detailed or rigorous proof of its unbiasedness in this setting. In Theorem 7, we close this theoretical gap by coupling the trajectories of FHS and MATU, thereby controlling their differences and formally establishing FHS’s unbiasedness.

7 CONCLUSION

In this paper, we provide a rigorous analysis of masked discrete diffusion. Differ from the analysis of uniform discrete diffusion, we show how to manage the initial KL blow-up and control the reverse-process KL divergence without relying on Girsanov theory. Building on this framework, we prove that Euler-type samplers TV converge in $\tilde{O}(d^2 \epsilon^{-3/2})$. We further introduce a mask-aware truncated uniformization sampler that removes the $\ln(1/\epsilon)$ factor, achieving nearly ϵ -free complexity. This acceleration aligns with the practical observation that masked diffusion denoises each masked token only once, whereas uniform diffusion repeatedly re-denoises already denoised tokens. Our results not only establish the first rigorous foundations for masked discrete diffusion but also explain why masked diffusion significantly reduces overhead in practice, opening avenues for more efficient text generation and advanced masked sampling techniques.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in neural information processing systems*, 34:17981–17993, 2021.
- 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, 2024.
- 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.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11315–11325, 2022.
- Hongrui Chen and Lexing Ying. Convergence analysis of discrete diffusion model: Exact implementation through uniformization. *arXiv preprint arXiv:2402.08095*, 2024.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *International Conference on Learning Representations*, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-predict: Parallel decoding of conditional masked language models. *arXiv preprint arXiv:1904.09324*, 2019.
- Emiel Hooeboom, Alexey A Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg, and Tim Salimans. Autoregressive diffusion models. In *International Conference on Learning Representations*.
- Xunpeng Huang, Yingyu Lin, Nikki Lijing Kuang, Hanze Dong, Difan Zou, Yian Ma, and Tong Zhang. Almost linear convergence under minimal score assumptions: Quantized transition diffusion. *arXiv preprint arXiv:2505.21892*, 2025.
- Yuchen Liang, Renxiang Huang, Lifeng Lai, Ness Shroff, and Yingbin Liang. Absorb and converge: Provable convergence guarantee for absorbing discrete diffusion models. *arXiv preprint arXiv:2506.02318*, 2025a.
- Yuchen Liang, Yingbin Liang, Lifeng Lai, and Ness Shroff. Discrete diffusion models: Novel analysis and new sampler guarantees. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025b.
- 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*, pp. 32819–32848, 2024.
- Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. *arXiv preprint arXiv:2406.03736*, 2024.
- Yinuo Ren, Haoxuan Chen, Grant M Rotskoff, and Lexing Ying. How discrete and continuous diffusion meet: Comprehensive analysis of discrete diffusion models via a stochastic integral framework. *arXiv preprint arXiv:2410.03601*, 2024.

- 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.
- Simo Särkkä and Arno Solin. *Applied stochastic differential equations*, volume 10. Cambridge University Press, 2019.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. pmlr, 2015.
- Aadithya Srikanth, Mudit Gaur, and Vaneet Aggarwal. Discrete state diffusion models: A sample complexity perspective. *arXiv preprint arXiv:2510.10854*, 2025.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 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.
- Zhengyan Wan, Yidong Ouyang, Qiang Yao, Liyan Xie, Fang Fang, Hongyuan Zha, and Guang Cheng. Error analysis of discrete flow with generator matching. *arXiv preprint arXiv:2509.21906*, 2025.
- Zikun Zhang, Zixiang Chen, and Quanquan Gu. Convergence of score-based discrete diffusion models: A discrete-time analysis. *arXiv preprint arXiv:2410.02321*, 2024.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2), 2023.
- Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. *arXiv preprint arXiv:2409.02908*, 2024.

CONTENTS

1	Introduction	1
2	Preliminaries	2
3	The Forward Process of Masked Discrete Diffusion	4
4	Euler Sampler in Masked Discrete Diffusion	5
5	Truncated Uniformization in Masked Discrete Diffusion	6
6	Related Work	9
7	Conclusion	10
A	Notation Summary	14
B	The Markov Processes of Discrete Diffusion Models	14
B.1	The Formulations of the Forward Process	14
B.2	The Proof of Lemma 1	16
B.3	The Proof of Lemma 2	17
C	Euler Discretization Analysis	23
C.1	Proof of Theorem 1	23
D	Truncated Uniformization Inference Analysis	28
D.1	The Proof of Lemma 3	28
D.2	The convergence of Alg. 1	30
D.3	The complexity of Alg. 1	34
E	Technical Lemmas	39
F	FHS convergence under Time-Independent Score Parameterization	41
G	Experiments	44
G.1	Synthetic Experiments.	44
G.2	Real World Experiments	44

NOTATION SUMMARY

We summarize all notations used in the main paper and appendix in Table 2.

Table 2: Summary of key notations used in the paper.

Symbol	Description
q_*	Discrete distribution on $\mathcal{Y} = \{1, 2, \dots, K\}^d$
\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 $(\mathbf{y}_{t'}^\rightarrow, \mathbf{y}_t^\rightarrow)$
\tilde{q}_t^\rightarrow	Approximation of q_t^\rightarrow constructing the reverse initialization, Eq. (9)
$q_{t' t}^\rightarrow(\mathbf{y}' \mathbf{y})$	Conditional transition probability in forward process, Eq. (37)
\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	Marginal distribution of reverse process at time t implemented by Alg. 1
$\hat{q}_{t',t}$	Joint distribution of $(\hat{\mathbf{y}}_{t'}, \hat{\mathbf{y}}_t)$
$\hat{q}_{t' t}(\mathbf{y}' \mathbf{y})$	Conditional transition probability of the ideal reverse process
$R^\rightarrow(\mathbf{y}, \mathbf{y}')$	Forward transition rate, i.e., Eq. (7), from state \mathbf{y}' to \mathbf{y} . 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. (6)
$\hat{R}_t(\cdot, \cdot)$	Truncated version of $\tilde{R}_t(\cdot, \cdot)$ with threshold β_t , Eq. (18)
$R_t^\leftarrow(\mathbf{y}), \tilde{R}_t(\mathbf{y}), \hat{R}_t(\mathbf{y})$	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 $R_t^\leftarrow(\mathbf{y})$, $\beta_t = \text{numK}(\mathbf{y}) \cdot K/(T-t)$, Eq. (17)
$v_{t,\mathbf{y}'}(\mathbf{y})$	Density ratio $q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$
$\tilde{v}_{t,\mathbf{y}'}(\mathbf{y})$	Learned approximation to $v_{t,\mathbf{y}'}(\mathbf{y}) = q_t^\leftarrow(\mathbf{y})/q_t^\leftarrow(\mathbf{y}')$
$\text{numK}(\cdot)$	The number of [MASK] token (or token K) in a vector.
$L_{SE}(\hat{v})$	Score entropy loss used to train \tilde{v} , Eq. (6)
\mathbf{e}_i	One-hot vector with a 1 at position i and 0 elsewhere
$\delta_{\mathbf{y}}(\cdot)$	Indicator function with $\delta_{\mathbf{y}}(\mathbf{y}) = 1$ and $\delta_{\mathbf{y}}(\mathbf{y}') = 0$ ($\mathbf{y}' \neq \mathbf{y}$)

THE MARKOV PROCESSES OF DISCRETE DIFFUSION MODELS

B.1 THE FORMULATIONS OF THE FORWARD PROCESS

Semigroup Formulation. In general, the time-homogeneous CTMC can be described by a Markov semigroup $\mathcal{Q}_t^\rightarrow$ defined as:

$$\mathcal{Q}_t^\rightarrow[f](\mathbf{y}) = \mathbb{E}[f(\mathbf{y}_t)|\mathbf{y}_0 = \mathbf{y}] = \left\langle f, q_{t|0}^\rightarrow(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}} \quad (22)$$

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

$$\mathcal{L}^\rightarrow[f](\mathbf{y}) = \lim_{t \rightarrow 0} \left[\frac{\mathcal{Q}_t^\rightarrow[f] - f}{t} \right](\mathbf{y}) = \left\langle f, \partial_t q_{t|0}^\rightarrow(\cdot|\mathbf{y}) \Big|_{t=0} \right\rangle_{\mathcal{Y}} := \langle f, R^\rightarrow(\cdot, \mathbf{y}) \rangle_{\mathcal{Y}} \quad (23)$$

where

$$R^\rightarrow(\mathbf{y}', \mathbf{y}) := \partial_t q_{t|0}^\rightarrow(\mathbf{y}'|\mathbf{y}) \Big|_{t=0} = \lim_{t \rightarrow 0} \left[\frac{q_{t|0}^\rightarrow(\mathbf{y}'|\mathbf{y}) - \delta_{\mathbf{y}}(\mathbf{y}')}{t} \right]. \quad (24)$$

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

$$R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \geq 0 \text{ when } \mathbf{y}' \neq \mathbf{y} \text{ and } R^{\rightarrow}(\mathbf{y}', \mathbf{y}') = - \sum_{\mathbf{y} \neq \mathbf{y}'} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \leq 0 \quad (25)$$

due to the definition Eq. (24). Under this setting, we can provide the dynamic of $q_{t|0}$ for any t . Specifically, we have

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

where the first inequality follows from the semigroup property. Combined 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. (22), 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

$$\partial_t q_{t|s}^{\rightarrow}(\tilde{\mathbf{y}}|\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} R(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{t|s}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) = \left\langle R(\tilde{\mathbf{y}}, \cdot), q_{t|s}^{\rightarrow}(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}}. \quad (26)$$

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

$$\frac{dq_t^{\rightarrow}}{dt}(\mathbf{y}) = \langle R(\mathbf{y}, \cdot), q_t^{\rightarrow} \rangle_{\mathcal{Y}}. \quad (27)$$

Matrix Formulation. Suppose the support set \mathcal{Y} of q_t^{\rightarrow} be written as $\mathcal{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{|\mathcal{Y}|}\}$, we may consider the marginal distribution q_s^{\rightarrow} to be a vector, i.e.,

$$\mathbf{q}_t^{\rightarrow} = [q_t(\mathbf{y}_1), q_t(\mathbf{y}_2), \dots, q_t(\mathbf{y}_{|\mathcal{Y}|})],$$

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

$$\mathbf{Q}_{t|s}^{\rightarrow} = \begin{bmatrix} q_{t|s}^{\rightarrow}(\mathbf{y}_1|\mathbf{y}_1) & q_{t|s}^{\rightarrow}(\mathbf{y}_1|\mathbf{y}_2) & \dots & q_{t|s}^{\rightarrow}(\mathbf{y}_1|\mathbf{y}_{|\mathcal{Y}|}) \\ q_{t|s}^{\rightarrow}(\mathbf{y}_2|\mathbf{y}_1) & q_{t|s}^{\rightarrow}(\mathbf{y}_2|\mathbf{y}_2) & \dots & q_{t|s}^{\rightarrow}(\mathbf{y}_2|\mathbf{y}_{|\mathcal{Y}|}) \\ \dots & \dots & \dots & \dots \\ q_{t|s}^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}|\mathbf{y}_1) & q_{t|s}^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}|\mathbf{y}_2) & \dots & q_{t|s}^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}|\mathbf{y}_{|\mathcal{Y}|}) \end{bmatrix}.$$

Similarly, the function R can also be presented as

$$\mathbf{R}^{\rightarrow} = \begin{bmatrix} R^{\rightarrow}(\mathbf{y}_1, \mathbf{y}_1) & R^{\rightarrow}(\mathbf{y}_1, \mathbf{y}_2) & \dots & R^{\rightarrow}(\mathbf{y}_1, \mathbf{y}_{|\mathcal{Y}|}) \\ R^{\rightarrow}(\mathbf{y}_2, \mathbf{y}_1) & R^{\rightarrow}(\mathbf{y}_2, \mathbf{y}_2) & \dots & R^{\rightarrow}(\mathbf{y}_2, \mathbf{y}_{|\mathcal{Y}|}) \\ \dots & \dots & \dots & \dots \\ R^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}, \mathbf{y}_1) & R^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}, \mathbf{y}_2) & \dots & R^{\rightarrow}(\mathbf{y}_{|\mathcal{Y}|}, \mathbf{y}_{|\mathcal{Y}|}) \end{bmatrix}. \quad (28)$$

Under this condition, Eq. (27) can be written as

$$d\mathbf{q}_t^{\rightarrow}/dt = \mathbf{R}^{\rightarrow} \cdot \mathbf{q}_t^{\rightarrow} \quad (29)$$

matching the usual presentation shown in Chen & Ying (2024); Zhang et al. (2024).

B.2 THE PROOF OF LEMMA 1

The proof of Lemma 1. For any $t \in [0, T]$, the marginal, joint, and conditional distribution w.r.t. $\{\mathbf{y}_t^{\leftarrow}\}$ are denoted as

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

which have $q_t^{\leftarrow} = q_{T-t}^{\rightarrow}$. Then, we start to check the dynamic of $q_{t|s}^{\leftarrow}$, i.e.,

$$\begin{aligned} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) &= -1 \cdot \partial_{T-t} q_{T-t|T-s}^{\rightarrow}(\mathbf{y}'|\mathbf{y}) = -1 \cdot \partial_{T-t} \left[\frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}') \cdot q_{T-t}^{\rightarrow}(\mathbf{y}')}{q_{T-s}^{\rightarrow}(\mathbf{y})} \right] \\ &= - \underbrace{\partial_{T-t} q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}') \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y}')}{q_{T-s}^{\rightarrow}(\mathbf{y})}}_{\text{Term 1}} - \underbrace{\frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}')}{q_{T-s}^{\rightarrow}(\mathbf{y})} \cdot \partial_{T-t} q_{T-t}^{\rightarrow}(\mathbf{y}')}_{\text{Term 2}}. \end{aligned} \quad (30)$$

For Term 1 of Eq. (30), we have

$$\begin{aligned} \text{Term 1} &= - \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} R^{\rightarrow}(\tilde{\mathbf{y}}, \mathbf{y}') \cdot q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\tilde{\mathbf{y}}) \cdot \frac{q_{T-t}^{\rightarrow}(\tilde{\mathbf{y}})}{q_{T-s}^{\rightarrow}(\mathbf{y})} \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y}')}{q_{T-t}^{\rightarrow}(\tilde{\mathbf{y}})} \\ &= - \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} R^{\rightarrow}(\tilde{\mathbf{y}}, \mathbf{y}') \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y}')}{q_{T-t}^{\rightarrow}(\tilde{\mathbf{y}})} \cdot q_{T-t|T-s}^{\rightarrow}(\tilde{\mathbf{y}}|\mathbf{y}), \end{aligned}$$

where the first equation follows from the Kolmogorov backward theorem (Lemma 14) and Eq. (23):

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

For Term 2 of Eq. (30), we have

$$\begin{aligned} \text{Term 2} &= \frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}')}{q_{T-s}^{\rightarrow}(\mathbf{y})} \cdot \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} R^{\rightarrow}(\mathbf{y}', \tilde{\mathbf{y}}) \cdot q_{T-t}^{\rightarrow}(\tilde{\mathbf{y}}) \\ &= \frac{q_{T-s|T-t}^{\rightarrow}(\mathbf{y}|\mathbf{y}') \cdot q_{T-t}^{\rightarrow}(\mathbf{y}')}{q_{T-s}^{\rightarrow}(\mathbf{y})} \cdot \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} R^{\rightarrow}(\mathbf{y}', \tilde{\mathbf{y}}) \cdot \frac{q_{T-t}^{\rightarrow}(\tilde{\mathbf{y}})}{q_{T-t}^{\rightarrow}(\mathbf{y}')} = 0, \end{aligned}$$

where the first equation follows from Eq. (27) and the last equation follows from the fact

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

Under this condition, by setting

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

then Eq. (30) can be summarized as

$$\partial_t q_{t|s}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) = \left\langle R_t^{\leftarrow}(\mathbf{y}', \cdot), q_{t|s}^{\leftarrow}(\cdot|\mathbf{y}) \right\rangle_{\mathcal{Y}} = \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} R_t^{\leftarrow}(\mathbf{y}', \tilde{\mathbf{y}}) \cdot q_{t|s}^{\leftarrow}(\tilde{\mathbf{y}}|\mathbf{y}). \quad (31)$$

Combining with Bayes' Theorem, we have

$$\frac{dq_t^{\leftarrow}}{dt}(\mathbf{y}) = \langle R_t^{\leftarrow}(\mathbf{y}, \cdot), q_t^{\leftarrow} \rangle_{\mathcal{Y}}. \quad (32)$$

Hence, Eq. (3) establishes.

Moreover, since 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 14 and Eq. (23)

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

Combining with Eq. (3), we have

$$\lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}|\mathbf{y}') - \delta_{\mathbf{y}'}(\mathbf{y})}{\Delta t} \right] = \lim_{s \rightarrow t} \partial_t q_{t|s}^{\leftarrow}(\mathbf{y}|\mathbf{y}') = R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\mathbf{y})}{q_{T-t}^{\rightarrow}(\mathbf{y}')} = R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') \quad (33)$$

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. (33). Hence, Eq. (3) establishes, and the proof is completed. \square

B.3 THE PROOF OF LEMMA 2

Lemma 4. *The close solution of Eq. (29) is*

$$\mathbf{q}_t^{\rightarrow} = \exp(t\mathbf{R}^{\rightarrow}) \cdot \mathbf{q}_0^{\rightarrow} \quad \text{where} \quad \exp(t\mathbf{R}^{\rightarrow}) = \sum_{i=0}^{\infty} \frac{1}{i!} (t\mathbf{R}^{\rightarrow})^i = \mathbf{I} + t\mathbf{R}^{\rightarrow} + \frac{(t\mathbf{R}^{\rightarrow})^2}{2} + \dots$$

Proof. We can easily verify that

$$\frac{d\mathbf{q}_t^{\rightarrow}}{dt} = \frac{d}{dt} [\exp(t\mathbf{R}^{\rightarrow})\mathbf{q}_0^{\rightarrow}] = \frac{d}{dt} [\exp(t\mathbf{R}^{\rightarrow})] \mathbf{q}_0^{\rightarrow}.$$

With the following equation,

$$\frac{d}{dt} [\exp(t\mathbf{R}^{\rightarrow})] = \frac{d}{dt} \left[\sum_{i=0}^{\infty} \frac{(t\mathbf{R}^{\rightarrow})^i}{i!} \right] = \sum_{i=1}^{\infty} \frac{t^{i-1}}{(i-1)!} \cdot (\mathbf{R}^{\rightarrow})^i = \mathbf{R}^{\rightarrow} \cdot \sum_{j=0}^{\infty} \frac{(t\mathbf{R}^{\rightarrow})^j}{j!} = \mathbf{R}^{\rightarrow} \cdot \exp(t\mathbf{R}^{\rightarrow}),$$

we have

$$\frac{d\mathbf{q}_t^{\rightarrow}}{dt} = \mathbf{R}^{\rightarrow} \cdot \exp(t\mathbf{R}^{\rightarrow}) \cdot \mathbf{q}_0^{\rightarrow} = \mathbf{R}^{\rightarrow} \cdot \mathbf{q}_t^{\rightarrow}.$$

Hence, the proof is completed. \square

Lemma 5. Suppose the transition rate matrix \mathbf{R}^\rightarrow shown as Eq. (28) satisfies Eq. (7). It can be decomposed as

$$\mathbf{R}^\rightarrow = \sum_{i=1}^d \mathbf{R}_i^\rightarrow \quad \text{where} \quad \mathbf{R}_i^\rightarrow = \underbrace{\mathbf{I} \otimes \cdots \otimes \mathbf{I}}_{i-1 \text{ terms}} \otimes \mathbf{A} \otimes \cdots \otimes \mathbf{I},$$

where \otimes denotes the Kronecker product, \mathbf{I} denotes the identity matrix on $\mathbb{R}^{K \times K}$, and \mathbf{A} satisfies

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & \cdots & 0 \\ 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 0 \end{bmatrix}. \quad (34)$$

Proof. According to the calculation of the Kronecker product, we have

$$\mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{I}(\mathbf{y}_1, \mathbf{y}'_1) \cdots \mathbf{A}(\mathbf{y}_i, \mathbf{y}'_i) \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d).$$

Under this condition, suppose $\text{Ham}(\mathbf{y}, \mathbf{y}') \geq 2$ and $\text{DiffIdx}(\mathbf{y}, \mathbf{y}') = \{j_1, j_2, \dots\}$ without loss of generality, for any $j \notin \{j_1, j_2\}$, we have

$$\mathbf{R}_j^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{A}(\mathbf{y}_j, \mathbf{y}'_j) \cdot \mathbf{I}(\mathbf{y}_1, \mathbf{y}'_1) \cdots \underbrace{\mathbf{I}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1})}_{=0} \cdots \underbrace{\mathbf{I}(\mathbf{y}_{j_2}, \mathbf{y}'_{j_2})}_{=0} \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d) = 0.$$

Besides, for $j = j_1$, we have

$$\mathbf{R}_{j_1}^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{A}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1}) \cdot \mathbf{I}(\mathbf{y}_1, \mathbf{y}'_1) \cdots \underbrace{\mathbf{I}(\mathbf{y}_{j_2}, \mathbf{y}'_{j_2})}_{=0} \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d) = 0.$$

A similar result will be satisfied for $j = j_2$. Hence, it has

$$\mathbf{R}^\rightarrow(\mathbf{y}, \mathbf{y}') = \sum_{i=1}^d \mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}') = 0 \quad \text{when} \quad \text{Ham}(\mathbf{y}, \mathbf{y}') \geq 2$$

Then, suppose $\text{Ham}(\mathbf{y}, \mathbf{y}') = 1$ and $\text{DiffIdx}(\mathbf{y}, \mathbf{y}') = j_1$, for any $j \neq j_1$, we have

$$\mathbf{R}_j^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{A}(\mathbf{y}_j, \mathbf{y}'_j) \cdot \mathbf{I}(\mathbf{y}_0, \mathbf{y}'_0) \cdots \underbrace{\mathbf{I}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1})}_{=0} \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d) = 0.$$

Otherwise, when $j = j_1$, we have

$$\mathbf{R}_{j_1}^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{A}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1}) \cdot \mathbf{I}(\mathbf{y}_1, \mathbf{y}'_1) \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d) = \mathbf{A}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1})$$

where the second equation establishes since $\text{Ham}(\mathbf{y}, \mathbf{y}') = 1$ and $\mathbf{y}_j = \mathbf{y}'_j$ when $j \neq j_1$. Then, only when $\mathbf{y}_{j_1} = K$, we will have $\mathbf{A}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1}) = 1$ otherwise $\mathbf{A}(\mathbf{y}_{j_1}, \mathbf{y}'_{j_1}) = 0$ due to the definition Eq. (34). That means

$$\mathbf{R}^\rightarrow(\mathbf{y}, \mathbf{y}') = \sum_{i=1}^d \mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}') = 0 \quad \text{when} \quad \text{Ham}(\mathbf{y}, \mathbf{y}') = 1 \text{ and } \mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')} \neq K$$

$$\mathbf{R}^\rightarrow(\mathbf{y}, \mathbf{y}') = \sum_{i=1}^d \mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}') = 1 \quad \text{when} \quad \text{Ham}(\mathbf{y}, \mathbf{y}') = 1 \text{ and } \mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')} = K.$$

Then, suppose $\text{Ham}(\mathbf{y}, \mathbf{y}') = 0$, i.e., $\mathbf{y} = \mathbf{y}'$, for any $j \in \{1, 2, \dots, d\}$, we have

$$\mathbf{R}_j^\rightarrow(\mathbf{y}, \mathbf{y}') = \mathbf{A}(\mathbf{y}_j, \mathbf{y}'_j) \cdot \mathbf{I}(\mathbf{y}_1, \mathbf{y}'_1) \cdots \mathbf{I}(\mathbf{y}_d, \mathbf{y}'_d) = \mathbf{A}(\mathbf{y}_j, \mathbf{y}'_j),$$

and

$$\sum_{i=1}^d \mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}') = \sum_{j=1}^d \mathbf{A}(\mathbf{y}_j, \mathbf{y}_j) = - \sum_{i=1}^d (1 - \delta_K(\mathbf{y}_i)),$$

which implies we have $\mathbf{R}^\rightarrow(\mathbf{y}, \mathbf{y}') = \sum_{i=0}^{d-1} \mathbf{R}_i^\rightarrow(\mathbf{y}, \mathbf{y}')$ when $\mathbf{y} = \mathbf{y}'$. Hence, the proof is completed. \square

Lemma 6. *With the decomposition shown in Lemma 5, i.e.,*

$$\mathbf{R}^{\rightarrow} = \sum_{i=1}^d \mathbf{R}_i^{\rightarrow} \quad \text{where} \quad \mathbf{R}_i^{\rightarrow} = \underbrace{\mathbf{I} \otimes \dots \otimes \mathbf{I}}_{i-1 \text{ terms}} \otimes \mathbf{A} \otimes \underbrace{\mathbf{I} \otimes \dots \otimes \mathbf{I}}_{d-i \text{ terms}}$$

for any $i, j \in \{1, 2, \dots, d\}$, the matrices $\mathbf{R}_i^{\rightarrow}$ and $\mathbf{R}_j^{\rightarrow}$ satisfy

$$\mathbf{R}_i^{\rightarrow} \cdot \mathbf{R}_j^{\rightarrow} = \mathbf{R}_j^{\rightarrow} \cdot \mathbf{R}_i^{\rightarrow},$$

which implies

$$\exp(t\mathbf{R}^{\rightarrow}) = \exp\left(t \sum_{i=1}^d \mathbf{R}_i^{\rightarrow}\right) = \prod_{i=1}^d \exp(t\mathbf{R}_i^{\rightarrow}) = \exp(t\mathbf{A})^{\otimes d}$$

Proof. According to Lemma 5, the matrix \mathbf{R}^{\rightarrow} has the following decomposition, i.e.,

$$\mathbf{R}^{\rightarrow} = \sum_{i=1}^d \mathbf{R}_i^{\rightarrow} \quad \text{where} \quad \mathbf{R}_i^{\rightarrow} = \underbrace{\mathbf{I} \otimes \dots \otimes \mathbf{I}}_{i-1 \text{ terms}} \otimes \mathbf{A} \otimes \underbrace{\mathbf{I} \otimes \dots \otimes \mathbf{I}}_{d-i \text{ terms}}$$

where \otimes denotes the Kronecker product, \mathbf{I} denotes the identity matrix on $\mathbb{R}^{K \times K}$, and \mathbf{A} satisfies

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{bmatrix}.$$

We can easily verify that the matrix \mathbf{A} can be decomposed as

$$\begin{bmatrix} -\mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{1}_{1 \times (K-1)} & 0 \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{I}_{K-1} & \mathbf{0} \\ -\mathbf{1}_{1 \times (K-1)} & 1 \end{bmatrix}}_{\mathbf{U}} \cdot \underbrace{\begin{bmatrix} -\mathbf{I}_{K-1} & 0 \\ \mathbf{0} & 0 \end{bmatrix}}_{\mathbf{\Lambda}} \cdot \underbrace{\begin{bmatrix} \mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{1}_{1 \times (K-1)} & 1 \end{bmatrix}}_{\mathbf{U}^{-1}} \quad \text{where} \quad \mathbf{U}\mathbf{U}^{-1} = \mathbf{U}^{-1}\mathbf{U} = \mathbf{I}_K. \quad (35)$$

Under this condition, $\mathbf{R}_i^{\rightarrow}$ can be reformulated as

$$\begin{aligned} \mathbf{R}_i^{\rightarrow} &= \underbrace{(\mathbf{U}\mathbf{U}^{-1}) \otimes \dots \otimes (\mathbf{U}\mathbf{U}^{-1})}_{i-1 \text{ terms}} \otimes (\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}) \otimes (\mathbf{U}\mathbf{U}^{-1}) \otimes \dots \otimes (\mathbf{U}\mathbf{U}^{-1}) \\ &= (\mathbf{U} \otimes \dots \otimes \mathbf{U}) \cdot \left(\underbrace{\mathbf{I} \otimes \dots \otimes \mathbf{I}}_{i-1 \text{ terms}} \otimes \mathbf{\Lambda} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I} \right) \cdot (\mathbf{U}^{-1} \otimes \dots \otimes \mathbf{U}^{-1}) := \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_i \cdot (\mathbf{U}^{-1})^{\otimes d} \end{aligned}$$

where the last inequality follows from Lemma 13. Under this condition, it has

$$\begin{aligned} \mathbf{R}_i^{\rightarrow} \cdot \mathbf{R}_j^{\rightarrow} &= \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_i \cdot (\mathbf{U}^{-1})^{\otimes d} \cdot \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_j \cdot (\mathbf{U}^{-1})^{\otimes d} = \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_i \cdot \mathbf{\Lambda}_j \cdot (\mathbf{U}^{-1})^{\otimes d} \\ &= \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_j \cdot \mathbf{\Lambda}_i \cdot (\mathbf{U}^{-1})^{\otimes d} = \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_i \cdot (\mathbf{U}^{-1})^{\otimes d} \cdot \mathbf{U}^{\otimes d} \cdot \mathbf{\Lambda}_j \cdot (\mathbf{U}^{-1})^{\otimes d} = \mathbf{R}_j^{\rightarrow} \cdot \mathbf{R}_i^{\rightarrow}, \end{aligned}$$

where the second and forth equations follows from Lemma 13 and Eq. (35).

For the property about the matrix exponential, we start from investigating the case of two commuting matrices, i.e., $\mathbf{R}_1^{\rightarrow}$ and $\mathbf{R}_2^{\rightarrow}$. By definition, we have

$$\exp(\mathbf{R}_1^{\rightarrow} + \mathbf{R}_2^{\rightarrow}) = \sum_{i=0}^{\infty} \frac{1}{i!} (\mathbf{R}_1^{\rightarrow} + \mathbf{R}_2^{\rightarrow})^i = \sum_{i=0}^{\infty} \frac{1}{i!} \sum_{j=0}^i C_i^j \cdot (\mathbf{R}_1^{\rightarrow})^j \cdot (\mathbf{R}_2^{\rightarrow})^{i-j}$$

where the last equation establishes since $\mathbf{R}_1^{\rightarrow}$ and $\mathbf{R}_2^{\rightarrow}$ are commute. Then, we have

$$\begin{aligned} \sum_{i=0}^{\infty} \frac{1}{i!} \sum_{j=0}^i C_i^j \cdot (\mathbf{R}_1^{\rightarrow})^j \cdot (\mathbf{R}_2^{\rightarrow})^{i-j} &= \sum_{i=0}^{\infty} \sum_{j=0}^i \frac{1}{i!} \cdot \frac{i!}{j!(i-j)!} \cdot (\mathbf{R}_1^{\rightarrow})^j \cdot (\mathbf{R}_2^{\rightarrow})^{i-j} \\ &= \sum_{i=0}^{\infty} \sum_{j=0}^i \frac{1}{j!(i-j)!} \cdot (\mathbf{R}_1^{\rightarrow})^j \cdot (\mathbf{R}_2^{\rightarrow})^{i-j} = \left(\sum_{j=0}^{\infty} \frac{(\mathbf{R}_1^{\rightarrow})^j}{j!} \right) \cdot \left(\sum_{i=0}^{\infty} \frac{(\mathbf{R}_2^{\rightarrow})^i}{i!} \right) = \exp(\mathbf{R}_1^{\rightarrow}) \cdot \exp(\mathbf{R}_2^{\rightarrow}). \end{aligned}$$

According to the definition of the matrix exponential, we will have $\exp(\mathbf{A} \otimes \mathbf{B}) = \exp(\mathbf{A}) \otimes \exp(\mathbf{B})$ when one of the factors is the identity. When we multiply all these exponentials, it has

$$\begin{aligned} \exp(\mathbf{R}_1^\rightarrow) \cdot \exp(\mathbf{R}_2^\rightarrow) &= [\exp(\mathbf{A}) \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I}] \cdot [\mathbf{I} \otimes \exp(\mathbf{A}) \otimes \dots \otimes \mathbf{I}] \\ &= [\exp(\mathbf{A}) \cdot \mathbf{I}] \otimes [\mathbf{I} \cdot \exp(\mathbf{A})] \otimes \mathbf{I} \dots \otimes \mathbf{I}. \end{aligned}$$

Then, following a recursive manner, we have

$$\exp\left(t \sum_{i=1}^d \mathbf{R}_i^\rightarrow\right) = \prod_{i=1}^d \exp(t \mathbf{R}_i^\rightarrow) = \exp(t \mathbf{A})^{\otimes d},$$

hence the proof is completed. \square

Lemma 7. Suppose matrix \mathbf{A} is

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{bmatrix},$$

the matrix exponential $\exp(t\mathbf{A})$ becomes

$$\exp(t\mathbf{A}) = \begin{bmatrix} e^{-t} & 0 & \dots & 0 & 0 \\ 0 & e^{-t} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 - e^{-t} & 1 - e^{-t} & \dots & 1 - e^{-t} & 1 \end{bmatrix}.$$

Proof. According to Lemma 4, $\bar{\mathbf{A}}(t) := \exp(t\mathbf{A})$ can be considered as the close solution of the following matrix ODE, i.e.,

$$\frac{d\bar{\mathbf{A}}(t)}{dt} = \mathbf{A} \cdot \bar{\mathbf{A}}(t), \quad \text{where} \quad \bar{\mathbf{A}}(0) = \mathbf{I}. \quad (36)$$

To provide a close form of $\bar{\mathbf{A}}_t$, we first decompose the matrix \mathbf{A} as follows

$$\mathbf{A} = \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{C} & \mathbf{0} \end{bmatrix} \quad \text{where} \quad \mathbf{B} := -\mathbf{I}_{K-1} \in \mathbb{R}^{(K-1) \times (K-1)} \text{ and } \mathbf{C} := [1, 1, \dots, 1] \in \mathbb{R}^{1 \times (K-1)}.$$

Then, the ODE. (36) can be equivalently think column-by-column, the j -th column of $\bar{\mathbf{A}}(t)$ solves

$$\frac{d}{dt} \bar{\mathbf{a}}(t) = \mathbf{A} \bar{\mathbf{a}}(t) \quad \text{where} \quad \bar{\mathbf{a}}(0) = \mathbf{e}_j.$$

We use the block structure to split $\bar{\mathbf{a}}(t) \in \mathbb{R}^K$ into two parts, i.e., $\bar{\mathbf{a}}(t) = [\bar{\mathbf{a}}_1(t), \bar{\mathbf{a}}_K(t)]$ where $\mathbf{q}_1(t) \in \mathbb{R}^{K-1}$ and $\mathbf{a}_K(t) \in \mathbb{R}$ denotes the last coordinate. Under this condition, we have

$$\frac{d}{dt} \bar{\mathbf{a}}_1(t) = \mathbf{B} \bar{\mathbf{a}}_1(t) + \mathbf{0} \cdot \bar{\mathbf{a}}_K(t) = \mathbf{B} \bar{\mathbf{a}}_1(t).$$

According to the definition of $\mathbf{B} = -\mathbf{I}_{K-1}$, we have

$$\frac{d}{dt} \bar{\mathbf{a}}_1(t) = -\bar{\mathbf{a}}_1(t) \quad \Rightarrow \quad \bar{\mathbf{a}}_1(t) = e^{-t} \bar{\mathbf{a}}_1(0).$$

If we consider the solution of $\bar{\mathbf{a}}_K(t)$, it has

$$\frac{d}{dt} \bar{\mathbf{a}}_K(t) = \mathbf{C} \cdot \bar{\mathbf{a}}_1(t) + \mathbf{0} \cdot \bar{\mathbf{a}}_K(t) = \mathbf{C} \cdot e^{-t} \cdot \bar{\mathbf{a}}_1(0).$$

For the initial condition, i.e., $\bar{\mathbf{a}}(0) = \mathbf{e}_j$, where $j \in \{1, 2, \dots, K-1\}$ and $\mathbf{C} \cdot \bar{\mathbf{a}}_1(0) = 1$, then it has

$$\frac{d}{dt} \bar{\mathbf{a}}_K(t) = \mathbf{C} \cdot \bar{\mathbf{a}}_1(t) + \mathbf{0} \cdot \bar{\mathbf{a}}_K(t) = e^{-t},$$

which implies

$$\bar{\mathbf{a}}_K(t) = \bar{\mathbf{a}}_K(0) + 1 - e^{-t} = 1 - e^{-t}.$$

For the initial condition, $\bar{\mathbf{a}}(0) = \mathbf{e}_K$, we have $\mathbf{C} \cdot \bar{\mathbf{a}}_1(0) = 0$ and

$$\bar{\mathbf{a}}_K(t) = \bar{\mathbf{a}}_K(0) + 0 = 1.$$

Therefore, we have

$$\exp(t\mathbf{A}) = \begin{bmatrix} e^{-t} & 0 & \dots & 0 & 0 \\ 0 & e^{-t} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 - e^{-t} & 1 - e^{-t} & \dots & 1 - e^{-t} & 1 \end{bmatrix}.$$

□

Lemma 8 (Forward transition kernel). *Consider the forward CTMC, i.e., $\{\mathbf{y}_t\}_{t=0}^T$ with the infinitesimal operator \mathbf{R}^\rightarrow given in Eq. (7). Then, for any two timestamps $s \leq t$, the forward transition probability satisfies, for any $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$,*

$$\begin{aligned} q_{t|s}^\rightarrow(\mathbf{y}|\mathbf{y}') &= \prod_{i=1}^d \left[\delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i) + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_0(\mathbf{y}_i - \mathbf{y}'_i) \cdot e^{-(t-s)} \right. \\ &\quad \left. + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_K(\mathbf{y}_i) \cdot (1 - e^{-(t-s)}) \right]. \end{aligned} \quad (37)$$

Proof. Under the matrix presentation, Eq. (26) implies the transition matrix $\mathbf{Q}_{t|s}^\rightarrow$ can be considered as the solution of the ODE

$$d\mathbf{Q}_{t|s}^\rightarrow/dt = \mathbf{R}^\rightarrow \cdot \mathbf{Q}_{t|s}^\rightarrow \quad \text{where} \quad \mathbf{Q}_{s|s}^\rightarrow = \mathbf{I}.$$

Combining Lemma 4 and 6, we have

$$\mathbf{Q}_{t|s}^\rightarrow = \exp((t-s)\mathbf{R}^\rightarrow) = \exp((t-s)\mathbf{A})^{\otimes d}, \quad (38)$$

which implies

$$\mathbf{Q}_{t|s}^\rightarrow = \begin{bmatrix} e^{-(t-s)} & 0 & \dots & 0 & 0 \\ 0 & e^{-(t-s)} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 - e^{-(t-s)} & 1 - e^{-(t-s)} & \dots & 1 - e^{-(t-s)} & 1 \end{bmatrix}^{\otimes d}$$

due to the close solution of $\exp((t-s)\mathbf{A})$ shown in Lemma 7. Combining this result with the calculation of the Kronecker product Lemma 12, we have

$$\begin{aligned} q_{t|s}^\rightarrow(\mathbf{y}|\mathbf{y}') &= \prod_{i=1}^d \left[\delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i) + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_0(\mathbf{y}_i - \mathbf{y}'_i) \cdot e^{-(t-s)} \right. \\ &\quad \left. + (1 - \delta_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \delta_K(\mathbf{y}_i) \cdot (1 - e^{-(t-s)}) \right]. \end{aligned}$$

where $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$. Hence, the proof is completed. □

The proof of Lemma 2. According to Eq. (29), the solution of \mathbf{q}_t^\rightarrow can be calculated as

$$\mathbf{q}_t^\rightarrow = \exp(t\mathbf{R}^\rightarrow) \cdot \mathbf{q}_0^\rightarrow = \exp(t\mathbf{A})^{\otimes d} \cdot \mathbf{q}_0^\rightarrow = \begin{bmatrix} e^{-t} & 0 & \dots & 0 & 0 \\ 0 & e^{-t} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & e^{-t} & 0 \\ 1 - e^{-t} & 1 - e^{-t} & \dots & 1 - e^{-t} & 1 \end{bmatrix}^{\otimes d} \cdot \mathbf{q}_0^\rightarrow$$

where the first equation follows from Lemma 4, the second equation follows from Lemma 6, and the last equation follows from Lemma 7. With the calculation of the Kronecker product Lemma 12, we have

$$q_t^{\rightarrow}(\mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} \exp(t\mathbf{A})^{\otimes d}(\mathbf{y}, \mathbf{y}') \cdot q_0^{\rightarrow}(\mathbf{y}') = \sum_{\mathbf{y}'} \left[\prod_{i=1}^d \exp(t\mathbf{A})(\mathbf{y}_i, \mathbf{y}'_i) \right] \cdot q_0^{\rightarrow}(\mathbf{y}'). \quad (39)$$

Under this condition, for any \mathbf{y} , we denote the coordinate set of token K as \mathcal{K} satisfying $\mathbf{y}_i = K \quad \forall i \in \mathcal{K}(\mathbf{y})$, and

$$\mathbf{y}_{\mathcal{K}^c(\mathbf{y})} = \mathbf{y}'_{\mathcal{K}^c(\mathbf{y})} \quad \Leftrightarrow \quad \mathbf{y}_i = \mathbf{y}'_i \quad \forall i \notin \mathcal{K}(\mathbf{y}).$$

Then, Eq. (39) can be rewritten as

$$\begin{aligned} q_t^{\rightarrow}(\mathbf{y}) &= \sum_{\mathbf{y}'_{\mathcal{K}^c(\mathbf{y})} = \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}} \left[\prod_{j \notin \mathcal{K}} \exp(t\mathbf{A})(\mathbf{y}_j, \mathbf{y}'_j) \cdot \prod_{j \in \mathcal{K}} \exp(t\mathbf{A})(K, \mathbf{y}'_j) \right] \cdot q_0^{\rightarrow}(\mathbf{y}') \\ &\quad + \sum_{\mathbf{y}'_{\mathcal{K}^c(\mathbf{y})} \neq \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}} \left[\prod_{j=1}^d \exp(t\mathbf{A})(\mathbf{y}_j, \mathbf{y}'_j) \right] \cdot q_0^{\rightarrow}(\mathbf{y}') \\ &= \sum_{\mathbf{y}'_{\mathcal{K}^c(\mathbf{y})} = \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}} \left[e^{-t \cdot |\mathcal{K}^c(\mathbf{y})|} \cdot (1 - e^{-t})^{|\mathcal{K}(\mathbf{y})|} \right] \cdot q_0^{\rightarrow}(\mathbf{y}') \\ &\leq e^{-t \cdot (d - \text{numK}(\mathbf{y}))} \cdot \sum_{\mathbf{y}'_{\mathcal{K}^c(\mathbf{y})} = \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}} q_0^{\rightarrow}(\mathbf{y}') \leq \exp(-t \cdot (d - \text{numK}(\mathbf{y}))), \end{aligned}$$

where the second equation establishes since we have

$$\exp(t\mathbf{A})(\mathbf{y}_j, \mathbf{y}'_j) = \begin{cases} e^{-t} & \mathbf{y}_j = \mathbf{y}'_j \quad \text{and} \quad \mathbf{y}_j \neq K \\ \mathbf{1}_K(\mathbf{y}'_j) \cdot (1 - e^{-t}) + (1 - \mathbf{1}_K(\mathbf{y}'_j)) & \mathbf{y}_j = K \\ 0 & \text{otherwise} \end{cases}.$$

According to the definition of $\tilde{q}(\mathbf{y})$, we can calculate the normalizing constant of \tilde{q} as

$$\tilde{Z}_t = \sum_{\mathbf{y}} \exp(-t \cdot (d - \text{numK}(\mathbf{y}))) = \sum_{i=0}^d \sum_{\text{numK}(\mathbf{y})=i} \exp(-t \cdot (d - i)) = \sum_{i=1}^d C_d^i \cdot e^{-t \cdot i} = (1 + e^{-t})^d.$$

Therefore, the KL divergence between q_t^{\rightarrow} and \tilde{q}_t can be written as

$$\begin{aligned} \text{KL}(q_t^{\rightarrow} \| \tilde{q}_t) &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\rightarrow}(\mathbf{y}) \cdot \ln \frac{q_t^{\rightarrow}(\mathbf{y})}{\tilde{q}_t(\mathbf{y})} = q_t^{\rightarrow}([K, \dots, K]) \cdot \ln \frac{q_t^{\rightarrow}([K, \dots, K])}{\tilde{q}_t([K, \dots, K])} + \sum_{\mathbf{y} \neq [K, \dots, K]} q_t^{\rightarrow}(\mathbf{y}) \cdot \ln \frac{q_t^{\rightarrow}(\mathbf{y})}{\tilde{q}_t(\mathbf{y})} \\ &\leq \ln \tilde{Z}_t + \sum_{\mathbf{y} \neq [K, \dots, K]} q_t^{\rightarrow}(\mathbf{y}) \ln \frac{q_t^{\rightarrow}(\mathbf{y})}{\exp(-t \cdot (d - \text{numK}(\mathbf{y}))) / \tilde{Z}_t} = \ln \tilde{Z}_t + \sum_{\mathbf{y} \neq [K, \dots, K]} q_t^{\rightarrow}(\mathbf{y}) \ln \tilde{Z}_t \\ &\leq 2 \ln \tilde{Z}_t = 2 \ln [1 + (1 + e^{-t})^d - 1] \leq 2 \cdot (1 + e^{-t})^d - 2. \end{aligned}$$

Suppose we require the TV distance to be small enough, e.g.,

$$\text{KL}(q_t^{\rightarrow} \| \tilde{q}_t) \leq \epsilon \quad \Leftrightarrow \quad (1 + e^{-t})^d - 1 \leq \epsilon/2 \quad \Leftrightarrow \quad d \ln(1 + e^{-t}) \leq \ln(1 + \epsilon/2),$$

then, since $\ln(1 + c) \leq c$ when $c > 0$, the sufficient condition for the establishment of the above equation is to require

$$d \cdot e^{-t} \leq \ln(1 + \epsilon/2) \quad \Leftrightarrow \quad t \geq \ln(d / \ln(1 + \epsilon/2)) \quad \Leftarrow \quad t \geq \ln(4d/\epsilon),$$

where the last derivation establishes since $\epsilon/4 \leq \ln(1 + \epsilon/2)$ when $\epsilon \leq 1$ without loss of generality. Hence, the proof is completed. \square

C EULER DISCRETIZATION ANALYSIS

By Assumption 2 of Liang et al. (2025a), $\tilde{v}_{t,y}(\mathbf{y}') \leq M$.

[A1]- Score approximation error assumption The discrete score \tilde{v}_t obtained from Eq. (6) is well-trained, and its estimation error satisfies for the chosen discretization step size h , and $T = nh + \delta$:

$$\frac{1}{T - \delta} \sum_{k=0}^{n-1} \int_{kh}^{(k+1)h} \mathbb{E}_{\mathbf{y}_t \sim q_t^{\leftarrow}} \left[\sum_{\mathbf{y} \neq \mathbf{y}_t} R^{\rightarrow}(\mathbf{y}_t, \mathbf{y}) D_{\phi}(v_{kh, \mathbf{y}_t}(\mathbf{y}) || \tilde{v}_{kh, \mathbf{y}_t}(\mathbf{y})) \right] dt \leq \epsilon_{score}^2.$$

C.1 PROOF OF THEOREM 1

Consider the Euler-discretization update in Eq. (11):

$$q_{t+\Delta t|t}^{Eu}(\mathbf{y}'|\mathbf{y}) \propto \delta_{\mathbf{y}}(\mathbf{y}') + \Delta t \cdot \tilde{R}_t(\mathbf{y}', \mathbf{y}) = \delta_{\mathbf{y}}(\mathbf{y}') + \Delta t \cdot R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot \tilde{v}_{t,y}(\mathbf{y}')$$

Without loss of generality, assume that $\tilde{R}_t(\mathbf{y}', \mathbf{y})^{\top}$ satisfies the two sufficient conditions of the transition rate matrix: its off-diagonal entries are non-negative, and each row sums to zero¹. In this way, both $e^{h\tilde{R}_t}$ and $I + h\tilde{R}_t$ are the transpose of valid transition matrices. The probability transition matrix of the Euler discretization can then be written as $Q_{t,t+h}^{Eu} = I + h\tilde{R}_t^{\top}$, where each element can be written as

$$Q_{t,t+h}^{Eu}(\mathbf{y}, \mathbf{y}') = q_{t+h|t}^{Eu}(\mathbf{y}'|\mathbf{y}) = \delta_{\mathbf{y}}(\mathbf{y}') + h \cdot \tilde{R}_t(\mathbf{y}', \mathbf{y}) \quad (40)$$

To prove the convergence bound for TV $(q_{\delta}^{\leftarrow}, q_{T-\delta}^{Eu})$, we introduce an auxiliary process q^{EI} using the exponential integrator update $Q_{t,t+h}^{Eu} = e^{h\tilde{R}_t^{\top}}$ (Zhang et al., 2024). We first prove the bound for TV $(q_{T-\delta}^{Eu}, q_{T-\delta}^{EI})$ and TV $(q_{\delta}^{\leftarrow}, q_{T-\delta}^{EI})$ separately, and use the triangle inequality to conclude the proof. Take $T = nh + \delta$.

Bound for TV $(q_{T-\delta}^{Eu}, q_{T-\delta}^{EI})$. For time interval $[kh, (k+1)h]$, by the chain rule of TV distance (Lemma 16), we have

$$\text{TV} \left(q_{(k+1)h}^{Eu}, q_{(k+1)h}^{EI} \right) \leq \text{TV} \left(q_{kh}^{Eu}, q_{kh}^{EI} \right) + \mathbb{E}_{\mathbf{y} \sim q_{kh}^{Eu}} \text{TV} \left(q_{(k+1)h|kh}^{Eu}(\cdot | \mathbf{y}), q_{(k+1)h|kh}^{EI}(\cdot | \mathbf{y}) \right) \quad (41)$$

By the definition of total variation distance, we have

$$\begin{aligned} \text{TV} \left(q_{(k+1)h|kh}^{Eu}(\cdot | \mathbf{y}), q_{(k+1)h|kh}^{EI}(\cdot | \mathbf{y}) \right) &= \sum_{\mathbf{y}'} \left| q_{(k+1)h|kh}^{Eu}(\mathbf{y}' | \mathbf{y}) - q_{(k+1)h|kh}^{EI}(\mathbf{y}' | \mathbf{y}) \right| \\ &= \sum_{\mathbf{y}'} \left| Q_{kh, (k+1)h}^{Eu}(\mathbf{y}, \mathbf{y}') - Q_{kh, (k+1)h}^{EI}(\mathbf{y}, \mathbf{y}') \right| \end{aligned} \quad (42)$$

Writing out the difference between $Q_{kh, (k+1)h}^{Eu} = I + h\tilde{R}_{kh}^{\top}$ and $Q_{kh, (k+1)h}^{EI} = e^{h\tilde{R}_{kh}^{\top}}$ using the Taylor series expansion for the matrix exponential:

$$Q_{kh, (k+1)h}^{EI} = e^{h\tilde{R}_{kh}^{\top}} = \sum_{i=0}^{\infty} \frac{1}{i!} (h\tilde{R}_{kh}^{\top})^i = I + h\tilde{R}_{kh}^{\top} + \frac{1}{2!} h^2 (\tilde{R}_{kh}^{\top})^2 + \frac{1}{3!} h^3 (\tilde{R}_{kh}^{\top})^3 + \dots,$$

we have

$$Q_{kh, (k+1)h}^{EI} - Q_{kh, (k+1)h}^{Eu} = e^{h\tilde{R}_{kh}^{\top}} - \left(I + h\tilde{R}_{kh}^{\top} \right) = \sum_{i=2}^{\infty} \frac{1}{i!} (h\tilde{R}_{kh}^{\top})^i.$$

¹Notice that our notation of R is the *transpose* of the convention used in some other works.

Thus, by the triangle inequality, we have

$$\begin{aligned}
& \sum_{\mathbf{y}' \in \mathcal{Y}} \left| Q_{kh,(k+1)h}^{EI}(\mathbf{y}, \mathbf{y}') - Q_{kh,(k+1)h}^{Eu}(\mathbf{y}, \mathbf{y}') \right| = \sum_{\mathbf{y}' \in \mathcal{Y}} \left| \sum_{i=2}^{\infty} \frac{1}{i!} \left((h \tilde{R}_{kh}^\top)^i \right) (\mathbf{y}, \mathbf{y}') \right| \\
& \leq \sum_{\mathbf{y}' \in \mathcal{Y}} \sum_{i=2}^{\infty} \frac{h^i}{i!} \left| \left((\tilde{R}_{kh}^\top)^i \right) (\mathbf{y}, \mathbf{y}') \right| \\
& = \sum_{i=2}^{\infty} \frac{h^i}{i!} \sum_{\mathbf{y}' \in \mathcal{Y}} \left| \left((\tilde{R}_{kh}^\top)^i \right) (\mathbf{y}, \mathbf{y}') \right| \quad (\text{Tonelli's theorem for series}) \\
& = \sum_{i=2}^{\infty} \frac{h^i}{i!} \sum_{\mathbf{y}' \in \mathcal{Y}} \left| \left((\tilde{R}_{kh})^i \right) (\mathbf{y}', \mathbf{y}) \right| \leq \sum_{i=2}^{\infty} \frac{h^i}{i!} \left\| (\tilde{R}_{kh})^i \right\|_1 \leq \sum_{i=2}^{\infty} \frac{h^i}{i!} \left\| \tilde{R}_{kh} \right\|_1^i,
\end{aligned}$$

where $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{i,j}| = \max_{x \neq \mathbf{0}} \|Ax\|_1 / \|x\|_1$ denotes the 1-norm of the matrix. And the last inequality is due to the multiplicative property of this matrix norm.

Therefore,

$$\begin{aligned}
\sum_{\mathbf{y}' \in \mathcal{Y}} \left| Q_{kh,(k+1)h}^{EI}(\mathbf{y}, \mathbf{y}') - Q_{kh,(k+1)h}^{Eu}(\mathbf{y}, \mathbf{y}') \right| & \leq \sum_{i=2}^{\infty} \frac{h^i}{i!} \left\| \tilde{R}_{kh} \right\|_1^i = e^{h \left\| \tilde{R}_{kh} \right\|_1} - 1 - h \left\| \tilde{R}_{kh} \right\|_1 \\
& \leq \left(h \left\| \tilde{R}_{kh} \right\|_1 \right)^2,
\end{aligned}$$

when $h \left\| \tilde{R}_{kh} \right\|_1 \leq 1$. Plugging this into Eq. (41) and (42), we have

$$\text{TV} \left(q_{(k+1)h}^{Eu}, q_{(k+1)h}^{EI} \right) \leq \text{TV} \left(q_{kh}^{Eu}, q_{kh}^{EI} \right) + \left(h \left\| \tilde{R}_{kh} \right\|_1 \right)^2, \quad (43)$$

when $h \left\| \tilde{R}_{kh} \right\|_1 \leq 1$.

By Assumption 2 of Liang et al. (2025a), $\tilde{v}_{t,\mathbf{y}}(\mathbf{y}') \leq M$. We have

$$\begin{aligned}
\left\| \tilde{R}_t \right\|_1 & = \max_{\mathbf{y}} \sum_{\mathbf{y}' \in \mathcal{Y}} \left| \tilde{R}_t(\mathbf{y}', \mathbf{y}) \right| = \max_{\mathbf{y}} \sum_{\mathbf{y}' \in \mathcal{Y}} |R^\rightarrow(\mathbf{y}, \mathbf{y}') \cdot \tilde{v}_{t,\mathbf{y}}(\mathbf{y}')| \\
& = \max_{\mathbf{y}} \left((d - \text{numK}(\mathbf{y})) + \sum_{\text{Ham}(\mathbf{y}, \mathbf{y}')=1 \text{ and } \mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')=K}} |\tilde{v}_{t,\mathbf{y}}(\mathbf{y}')| \right) \quad (\text{By Eq. 7}) \\
& \leq \max_{\mathbf{y}} (d - \text{numK}(\mathbf{y}) + KdM) \\
& \leq 2KdM.
\end{aligned}$$

Thus $\left\| \tilde{R}_{kh} \right\|_1 \leq 2KdM$. By (43) we have

$$\begin{aligned}
\text{TV} \left(q_{nh}^{Eu}, q_{nh}^{EI} \right) & \leq \text{TV} \left(q_0^{Eu}, q_0^{EI} \right) + \sum_{k=1}^n \left(h \left\| \tilde{R}_{kh} \right\|_1 \right)^2 \\
& \leq K^2 d^2 \sum_{k=1}^n h^2 M^2 \leq K^2 d^2 n h^2 M^2 \leq K^2 (T - \delta) h d^2 M^2.
\end{aligned} \quad (44)$$

By taking $h \leq \frac{\varepsilon}{K^2 d^2 M^2 \log(d/\varepsilon)}$, then $\text{TV} \left(q_{nh}^{Eu}, q_{nh}^{EI} \right) \leq \varepsilon$.

Bound for $\text{TV} \left(q_{T-\delta}^{EI}, q_{T-\delta}^{\leftarrow} \right)$. We first prove $\text{KL} \left(q_{T-\delta}^{\leftarrow} \| q_{T-\delta}^{EI} \right)$, then use Pinsker's inequality to derive the bound for $\text{TV} \left(q_{T-\delta}^{EI}, q_{T-\delta}^{\leftarrow} \right)$.

For time interval $[kh, (k+1)h]$, we have

$$\text{KL} \left(q_{(k+1)h}^{\leftarrow} \| q_{(k+1)h}^{EI} \right) = \text{KL} \left(q_{kh}^{\leftarrow} \| q_{kh}^{EI} \right) + \int_{kh}^{(k+1)h} \frac{d\text{KL} \left(q_t^{\leftarrow} \| q_t^{EI} \right)}{dt} dt. \quad (45)$$

By the chain rule of KL divergence (Lemma 15)

$$\begin{aligned} \frac{d}{dt} \text{KL} \left(q_t^{\leftarrow} \| q_t^{EI} \right) &= \lim_{\Delta t \rightarrow 0} \frac{\text{KL} \left(q_{t+\Delta t}^{\leftarrow} \| q_{t+\Delta t}^{EI} \right) - \text{KL} \left(q_t^{\leftarrow} \| q_t^{EI} \right)}{\Delta t} \\ &\leq \lim_{\Delta t \rightarrow 0} \mathbb{E}_{\mathbf{y} \sim q_t^{\leftarrow}} \frac{\text{KL} \left(q_{t+\Delta t|t}^{\leftarrow}(\cdot | \mathbf{y}) \| q_{t+\Delta t|t}^{EI}(\cdot | \mathbf{y}) \right)}{\Delta t} \\ &= \mathbb{E}_{\mathbf{y} \sim q_t^{\leftarrow}} \underbrace{\lim_{\Delta t \rightarrow 0} \frac{\text{KL} \left(q_{t+\Delta t|t}^{\leftarrow}(\cdot | \mathbf{y}) \| q_{t+\Delta t|t}^{EI}(\cdot | \mathbf{y}) \right)}{\Delta t}}_{\text{Term 1}} \end{aligned} \quad (46)$$

For each $\mathbf{y} \in \mathcal{Y}$, we focus on Term 1 of Eq. (46), and have

$$\begin{aligned} \text{Term 1} &= \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y}) \cdot \ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right] \\ &= \lim_{\Delta t \rightarrow 0} \underbrace{\left[\sum_{\mathbf{y}' \neq \mathbf{y}} \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right]}_{\text{Term 1.1}} + \\ &\quad \underbrace{\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y}) \right) \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right]}_{\text{Term 1.2}}. \end{aligned} \quad (47)$$

For Term 1.1, we have

$$\begin{aligned} \text{Term 1.1} &= \sum_{\mathbf{y}' \neq \mathbf{y}} \lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{\Delta t} \right] \cdot \lim_{\Delta t \rightarrow 0} \left[\ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \left[\lim_{\Delta t \rightarrow 0} \left(\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{\Delta t} \cdot \frac{\Delta t}{q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right) \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_{kh}(\mathbf{y}', \mathbf{y})}, \end{aligned} \quad (48)$$

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

$$\begin{aligned} \text{Term 1.2} &= \lim_{\Delta t \rightarrow 0} \left[1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y}) \right] \cdot \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}' | \mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{EI}(\mathbf{y}' | \mathbf{y})} \right] \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} \left(\tilde{R}_{kh}(\mathbf{y}', \mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right) = \tilde{R}_{kh}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \end{aligned} \quad (49)$$

where the first inequality follows from Lemma 9. Plugging Eq. (48), Eq. (49) and Eq. (47), into Eq. (46) we have

$$\begin{aligned}
\frac{d\text{KL}(q_t^{\leftarrow} \| q_t^{EI})}{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}_{kh}(\mathbf{y}', \mathbf{y})} + \tilde{R}_{kh}(\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}_{kh}(\mathbf{y}', \mathbf{y})} + \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_{kh}(\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})} + \tilde{v}_{kh, \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}) \tilde{v}_{kh, \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[-v_{t, \mathbf{y}}(\mathbf{y}') + \tilde{v}_{kh, \mathbf{y}}(\mathbf{y}') + v_{t, \mathbf{y}}(\mathbf{y}') \ln \frac{v_{t, \mathbf{y}}(\mathbf{y}')}{\tilde{v}_{kh, \mathbf{y}}(\mathbf{y}')} \right] \\
&= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \underbrace{\sum_{\text{Ham}(\mathbf{y}, \mathbf{y}')=1 \text{ and } \mathbf{y}_{\text{DiffIdx}}(\mathbf{y}, \mathbf{y}')=K} \left[-v_{t, \mathbf{y}}(\mathbf{y}') + \tilde{v}_{kh, \mathbf{y}}(\mathbf{y}') + v_{t, \mathbf{y}}(\mathbf{y}') \ln \frac{v_{t, \mathbf{y}}(\mathbf{y}')}{\tilde{v}_{kh, \mathbf{y}}(\mathbf{y}')} \right]}_{\text{Term 2}}
\end{aligned} \tag{50}$$

For $\mathbf{y}' = \mathbf{y}[y_i \rightarrow k]$, by Eq. (55) we have $v_{t, \mathbf{y}}(\mathbf{y}') = \frac{q_t^{\leftarrow}(\mathbf{y}[y_i \rightarrow k])}{q_t^{\leftarrow}(\mathbf{y})} \leq \frac{1}{e^{(T-t)} - 1}$. By (Liang et al., 2025a, Lemma 2), there exist $c > 0$ such that $v_{t, \mathbf{y}}(\mathbf{y}') \geq \frac{1}{c} e^{-(T-t)}$. Therefore, by (Zhang et al., 2024, Proposition 3), letting $C = \max\{M, ce^T\}$, Term 2 satisfies

$$\text{Term 2} \leq \sum_{\text{Ham}(\mathbf{y}, \mathbf{y}')=1 \text{ and } \mathbf{y}_{\text{DiffIdx}}(\mathbf{y}, \mathbf{y}')=K} \left(C \|v_{t, \mathbf{y}}(\mathbf{y}') - v_{kh, \mathbf{y}}(\mathbf{y}')\|^2 + 2C^2 D_\phi(v_{kh, \mathbf{y}}(\mathbf{y}') \| \tilde{v}_{kh, \mathbf{y}}(\mathbf{y}')) \right)$$

where D_ϕ is the Bregman divergence with $\phi(x) = x \ln x$ (as Eq. (6)), i.e.,

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

By (Liang et al., 2025a, Lemma 7), we have $\|v_{t, \mathbf{y}}(\mathbf{y}') - v_{kh, \mathbf{y}}(\mathbf{y}')\| \lesssim \gamma^{-1}(t - kh) \lesssim h$, where γ is defined in (Liang et al., 2025a, Assumption 4). We therefore have

$$\begin{aligned}
\text{Term 2} &\lesssim CK \text{numK}(\mathbf{y}) h^2 + C^2 \sum_{\text{Ham}(\mathbf{y}, \mathbf{y}')=1 \text{ and } \mathbf{y}_{\text{DiffIdx}}(\mathbf{y}, \mathbf{y}')=K} (D_\phi(v_{kh, \mathbf{y}}(\mathbf{y}') \| \tilde{v}_{kh, \mathbf{y}}(\mathbf{y}'))) \\
&\lesssim CK dh^2 + C^2 \sum_{\text{Ham}(\mathbf{y}, \mathbf{y}')=1 \text{ and } \mathbf{y}_{\text{DiffIdx}}(\mathbf{y}, \mathbf{y}')=K} (D_\phi(v_{kh, \mathbf{y}}(\mathbf{y}') \| \tilde{v}_{kh, \mathbf{y}}(\mathbf{y}')))
\end{aligned} \tag{51}$$

Since by Eq. (45), we have

$$\text{KL}(q_{nh}^{\leftarrow} \| q_{nh}^{EI}) = \text{KL}(q_0^{\leftarrow} \| q_0^{EI}) + \sum_{k=0}^{n-1} \int_{kh}^{(k+1)h} \frac{d\text{KL}(q_t^{\leftarrow} \| q_t^{EI})}{dt} dt.$$

Then, by Eq. (50), Eq. (51), Eq. (6) and Assumption [A1], we have

$$\text{KL}(q_{T-\delta}^{\leftarrow} \| q_{T-\delta}^{EI}) \lesssim (T - \delta) C^2 \epsilon_{\text{score}}^2 + C(T - \delta) K dh^2.$$

By Pinsker's inequality, we have

$$\text{TV}(q_{T-\delta}^{\leftarrow}, q_{T-\delta}^{EI}) \leq \sqrt{\frac{1}{2} \text{KL}(q_{T-\delta}^{\leftarrow} \| q_{T-\delta}^{EI})} \lesssim \sqrt{\frac{1}{2} \sqrt{(T - \delta) C^2 \epsilon_{\text{score}}^2 + C(T - \delta) K dh^2}}$$

By taking $\epsilon_{\text{score}} \lesssim \varepsilon/(\sqrt{TC})$, and $h \lesssim \varepsilon/\sqrt{CdT}$, we have $\text{TV}(q_{T-\delta}^{\leftarrow}, q_{T-\delta}^{EI}) \leq \varepsilon$.

Therefore, taking $h \lesssim \min\{\frac{\varepsilon}{K^2 d^2 M^2 \log(d/\varepsilon)}, \frac{\varepsilon}{\sqrt{Cd \log(d/\varepsilon)}}\}$, by the triangle inequality, we have

$$\text{TV}(q_{\delta}^{\leftarrow}, q_{T-\delta}^{Eu}) \leq \text{TV}(q_{T-\delta}^{Eu}, q_{T-\delta}^{EI}) + \text{TV}(q_{\delta}^{\leftarrow}, q_{T-\delta}^{EI}) \lesssim \varepsilon.$$

Plugging in $C = \Theta(d/\varepsilon)$, we have for $h \lesssim \min\{\frac{\varepsilon}{K^2 d^2 M^2 \log(d/\varepsilon)}, \frac{\varepsilon^{\frac{3}{2}}}{d\sqrt{\log(d/\varepsilon)}}\}$, we have $\text{TV}(q_{\delta}^{\leftarrow}, q_{T-\delta}^{Eu}) \lesssim \varepsilon$.

Hence, the proof is completed.

Lemma 9. *Following the notations shown in Section 2, for $t \in [kh, (k+1)h]$, we have*

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

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

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

which automatically makes

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

Under this condition, we have

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

which implies (with the dominated convergence theorem)

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

Only when $i = 1$, we have

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

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

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

Hence, the proof is completed. \square

D TRUNCATED UNIFORMIZATION INFERENCE ANALYSIS

D.1 THE PROOF OF LEMMA 3

The proof of Lemma 3. According to the definition, we have

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

Since the definition of the transition rate matrix, i.e., Eq. (7), for any \mathbf{y}' with $\text{Ham}(\mathbf{y}', \mathbf{y}) > 1$, it has $R^{\rightarrow}(\mathbf{y}, \mathbf{y}') = 0$. Moreover, even when $\text{Ham}(\mathbf{y}', \mathbf{y}) = 1$, it has

$$R^{\rightarrow}(\mathbf{y}, \mathbf{y}') = 0 \quad \text{when} \quad \mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')} \neq K.$$

Define the function to transfer the i -th element of \mathbf{y} (\mathbf{y}_i) from k' to k as

$$\mathbf{y}[\mathbf{y}_i: k' \rightarrow k] = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{i-1}, k, \mathbf{y}_{i+1}, \dots, \mathbf{y}_d].$$

That means $R_t^{\leftarrow}(\mathbf{y})$ can be rewritten as

$$R_t^{\leftarrow}(\mathbf{y}) = \sum_{i, \mathbf{y}_i=K} \left[\sum_{k=1}^{K-1} R^{\rightarrow}(\mathbf{y}, \mathbf{y}[\mathbf{y}_i: K \rightarrow k]) \cdot \frac{q_t^{\leftarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k])}{q_t^{\leftarrow}(\mathbf{y})} \right]. \quad (52)$$

To upper bound the RHS of the above equation, we consider controlling

$$\begin{aligned} \frac{q_t^{\leftarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k])}{q_t^{\leftarrow}(\mathbf{y})} &= \frac{q_{T-t}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k])}{q_{T-t}^{\rightarrow}(\mathbf{y})} = \frac{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k]|\mathbf{y}_0)}{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)} \\ &= \frac{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0) \cdot \frac{q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k]|\mathbf{y}_0)}{q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)}}{\sum_{\mathbf{y}_0 \in \mathcal{Y}} q_0^{\rightarrow}(\mathbf{y}_0) \cdot q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)} = \mathbb{E}_{\mathbf{y}_0 \sim q_0^{\rightarrow}(\cdot|\mathbf{y})} \left[\frac{q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k]|\mathbf{y}_0)}{q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)} \right], \end{aligned} \quad (53)$$

where the last equation follows from Bayes' Theorem, i.e.,

$$q_0^{\rightarrow|T-t}(\mathbf{y}_0|\mathbf{y}) \cdot q_{T-t}^{\rightarrow}(\mathbf{y}) = q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0) \cdot q_0^{\rightarrow}(\mathbf{y}_0) \Leftrightarrow q_0^{\rightarrow|T-t}(\mathbf{y}_0|\mathbf{y}) \propto q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0) \cdot q_0^{\rightarrow}(\mathbf{y}_0).$$

Then, we only need to control $q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i \rightarrow k]|\mathbf{y}_0)/q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)$ where both the denominator and the numerator can be calculated accurately by Lemma 8. Specifically, we have

$$\begin{aligned} q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0) &= \prod_{j \in \{1, \dots, i-1, i+1, \dots, d\}} \left[\mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j}) + (1 - \mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j})) \cdot \mathbf{1}_0(\mathbf{y}_j - \mathbf{y}_{0,j}) \cdot e^{-(T-t)} \right. \\ &\quad \left. + (1 - \mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j})) \cdot \mathbf{1}_K(\mathbf{y}_j) \cdot (1 - e^{-(T-t)}) \right] \\ &\quad \left[\mathbf{1}_{(K,K)}(K, \mathbf{y}_{0,i}) + (1 - \mathbf{1}_{(K,K)}(K, \mathbf{y}_{0,i})) \cdot (1 - e^{-(T-t)}) \right] \end{aligned}$$

and

$$\begin{aligned} q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i: K \rightarrow k]|\mathbf{y}_0) &= \prod_{j \in \{1, \dots, i-1, i+1, \dots, d\}} \left[\mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j}) + (1 - \mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j})) \cdot \mathbf{1}_0(\mathbf{y}_j - \mathbf{y}_{0,j}) \cdot e^{-(T-t)} \right. \\ &\quad \left. + (1 - \mathbf{1}_{(K,K)}(\mathbf{y}_j, \mathbf{y}_{0,j})) \cdot \mathbf{1}_K(\mathbf{y}_j) \cdot (1 - e^{-(T-t)}) \right] \\ &\quad \left[(1 - \mathbf{1}_{(K,K)}(k, \mathbf{y}_{0,i})) \cdot \mathbf{1}_0(k - \mathbf{y}_{0,i}) \cdot e^{-(T-t)} \right]. \end{aligned}$$

Since the factor except for the i -th term will be canceled, we have

$$\begin{aligned} \frac{q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i \rightarrow k]|\mathbf{y}_0)}{q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)} &= \frac{(1 - \mathbf{1}_{(K,K)}(k, \mathbf{y}_{0,i})) \cdot \mathbf{1}_0(k - \mathbf{y}_{0,i}) \cdot e^{-(T-t)}}{\mathbf{1}_{(K,K)}(K, \mathbf{y}_{0,i}) + (1 - \mathbf{1}_{(K,K)}(K, \mathbf{y}_{0,i})) \cdot (1 - e^{-(T-t)})} \\ &= \frac{\mathbf{1}_0(k - \mathbf{y}_{0,i}) \cdot e^{-(T-t)}}{1 - e^{-(T-t)}} \leq \frac{e^{-(T-t)}}{1 - e^{-(T-t)}} = \frac{1}{e^{(T-t)} - 1}. \end{aligned} \quad (54)$$

Plugging this result into Eq. (53), the density ratio of the reverse process will have

$$\frac{q_t^{\leftarrow}(\mathbf{y}[\mathbf{y}_i \rightarrow k])}{q_t^{\leftarrow}(\mathbf{y})} = \mathbb{E}_{\mathbf{y}_0 \sim q_{0|T-t}^{\rightarrow}(\cdot|\mathbf{y})} \left[\frac{q_{T-t|0}^{\rightarrow}(\mathbf{y}[\mathbf{y}_i : K \rightarrow k]|\mathbf{y}_0)}{q_{T-t|0}^{\rightarrow}(\mathbf{y}|\mathbf{y}_0)} \right] \leq \frac{1}{e^{(T-t)} - 1}. \quad (55)$$

Combining with the fact, i.e.,

$$R^{\rightarrow}(\mathbf{y}, \mathbf{y}[\mathbf{y}_i : K \rightarrow k]) = 1$$

from Eq. (7), Eq. (52) can be upper bounded as

$$R_t^{\leftarrow}(\mathbf{y}) = \sum_{i, \mathbf{y}_i = K} \left[\sum_{k=1}^{K-1} \frac{q_t^{\leftarrow}(\mathbf{y}[\mathbf{y}_i : K \rightarrow k])}{q_t^{\leftarrow}(\mathbf{y})} \right] \leq \frac{\text{numK}(\mathbf{y}) \cdot K}{e^{(T-t)} - 1}.$$

Hence, the proof is completed. \square

Remark 1. Here, an interesting property is that compared with the upper bound of $\beta_t(\mathbf{y})$ in the uniform forward process Chen & Ying (2024), i.e.,

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \leq K \cdot d \cdot \frac{1 + e^{-2(T-t)}}{1 - e^{-2(T-t)}} \leq K \cdot d \cdot (1 + (T-t)^{-1}).$$

the upper bound of $\beta_t(\mathbf{y})$ in absorbing forward process will only be

$$\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \leq K \cdot \text{numK}(\mathbf{y}) \cdot \frac{e^{-(T-t)}}{1 - e^{-(T-t)}}.$$

The latter upper bound is strictly better compared with the former one, since the number of mask tokens, i.e., $\text{numK}(\mathbf{y}) \leq d$. Besides, with the time growth (from 0 to T), $\text{numK}(\mathbf{y})$ will be monotonic decrease for $R_t^{\leftarrow}(\mathbf{y})$ (from d to 0). Since the dominating term in the complexity analysis of truncated uniformization is β_t , the discrete diffusion models with absorbing forward process are expected to have a better result. The mechanism of the acceleration can be explained in one sentence, i.e.,

At each uniformization step, absorbing the discrete diffusion model knows the token needs (masked token)/ or does not need (unmasked token) to denoise, and an unmasked token will not be denoised twice.

Rigorously, this property can be summarized by Lemma 10.

Lemma 10. Suppose Assumption [A2] hold, and $0 < t_0 \leq t$, we have $q_{t|t_0}^{\leftarrow}(\mathbf{y}|\mathbf{y}_0) \neq 0$ if and only if

$$\mathbf{y} \in \mathcal{Y}^{\leftarrow}(\mathbf{y}_0) = \{\mathbf{y}' | \forall i, \quad \mathbf{y}_{0,i} = K \text{ or } \mathbf{y}'_i = \mathbf{y}_{0,i}\}.$$

Proof. According to the Bayes' theorem, for any $t \geq t_0$, it has

$$\begin{aligned} q_{t,t_0}^{\leftarrow}(\mathbf{y}, \mathbf{y}_0) &= q_{t|t_0}^{\leftarrow}(\mathbf{y}|\mathbf{y}_0) \cdot q_{t_0}^{\leftarrow}(\mathbf{y}_0) = q_{T-t,T-t_0}^{\rightarrow}(\mathbf{y}, \mathbf{y}_0) \\ &= q_{T-t_0,T-t}^{\rightarrow}(\mathbf{y}_0, \mathbf{y}) = q_{T-t_0|T-t}^{\rightarrow}(\mathbf{y}_0|\mathbf{y}) \cdot q_{T-t}^{\rightarrow}(\mathbf{y}), \end{aligned} \quad (56)$$

where the third equation follows from the reversibility of the absorbing forward process shown in Campbell et al. (2022). Following from the forward transition kernel shown in Lemma 8, we know that

$$q_{T-t_0|T-t}^{\rightarrow}(\mathbf{y}_0|\mathbf{y}) \neq 0 \quad \Leftrightarrow \quad \mathbf{y}_0 \in \mathcal{Y}^{\rightarrow}(\mathbf{y}) = \{\mathbf{y}' | \forall i, \quad \mathbf{y}'_i = \mathbf{y}_i \text{ or } \mathbf{y}'_i = K\}. \quad (57)$$

Combining Assumption [A2] and Lemma 8, we have $q_{\tau}^{\rightarrow}(\mathbf{y}) > 0$ for all $\mathbf{y} \in \mathcal{Y}$, which implies

$$q_{t_0}^{\leftarrow}(\mathbf{y}_0) = q_{T-t_0}^{\rightarrow}(\mathbf{y}_0) > 0 \quad \text{and} \quad q_t^{\rightarrow}(\mathbf{y}) > 0. \quad (58)$$

Then, we can summarize

$$q_{t|t_0}^{\leftarrow}(\mathbf{y}|\mathbf{y}_0) \neq 0 \quad \Leftrightarrow \quad \mathbf{y} \in \mathcal{Y}^{\leftarrow}(\mathbf{y}_0) = \{\mathbf{y}' | \forall i, \quad \mathbf{y}_{0,i} = K \text{ or } \mathbf{y}'_i = \mathbf{y}_{0,i}\}.$$

Hence, the proof is completed. \square

D.2 THE CONVERGENCE OF ALG. 1

Suppose, with the infinitesimal reverse transition rate, the particles in Alg. 1 during the reverse process are denoted as random variables $\{\hat{\mathbf{y}}_t\}_{t=0}^{T-\delta}$, whose underlying distributions are \hat{q}_t . Then, the implementation will be equivalent to the following Poisson process. For $t \in (t_{w-1}, t_w]$, $\hat{\mathbf{y}}_{t_{w-1}} = \mathbf{y}_0$ and $\hat{\mathbf{y}}_t = \mathbf{y}$,

1. With probability $\Delta t \cdot \beta_{t_w}(\mathbf{y}_0)$, allow a state transition.
2. Conditioning on an allowed transition, move from \mathbf{y} to \mathbf{y}' with probability

$$\hat{M}_{t|t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) = \begin{cases} \beta_{t_w}^{-1}(\mathbf{y}_0) \cdot \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y}) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \beta_{t_w}^{-1}(\mathbf{y}_0) \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) & \text{otherwise} \end{cases}.$$

Here we should note that

$$\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) \leq \beta_t(\mathbf{y}) = K \cdot \text{numK}(\mathbf{y}) \cdot \frac{1}{e^{T-t} - 1} \leq K \cdot \text{numK}(\mathbf{y}_0) \cdot \frac{1}{e^{T-t_w} - 1} = \beta_{t_w}(\mathbf{y}_0),$$

where the second inequality established since $\text{numK}(\hat{\mathbf{y}}_t) \leq \text{numK}(\hat{\mathbf{y}}_{t_{w-1}})$ and $(e^{T-t} - 1)^{-1}$ is monotonic increasing. Under these two steps, the practical conditional probability satisfies

$$\begin{aligned} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) &= \begin{cases} \Delta t \cdot \beta_{t_w}(\mathbf{y}_0) \cdot \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y}) \cdot \beta_{t_w}^{-1}(\mathbf{y}_0) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \Delta t \cdot \beta_{t_w}(\mathbf{y}_0) + \Delta t \cdot \beta_{t_w}(\mathbf{y}_0) \cdot (1 - \beta_{t_w}(\mathbf{y}_0)^{-1} \cdot \hat{R}_{t, \mathbf{y}_0}(\mathbf{y})) & \mathbf{y}' = \mathbf{y}, \end{cases} \\ &= \begin{cases} \Delta t \cdot \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y}) & \mathbf{y}' \neq \mathbf{y} \\ 1 - \Delta t \cdot \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) & \mathbf{y}' = \mathbf{y} \end{cases}. \end{aligned} \quad (59)$$

Lemma 11. *Following the notations shown in Section A, we have*

$$\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] = \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}).$$

Proof. Since we have required $\Delta t \rightarrow 0$, for any $\mathbf{y}' \neq \mathbf{y}$, it has

$$\begin{aligned} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) &\rightarrow \hat{q}_{t|t}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) = 0 \\ \text{and } q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) &= q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \rightarrow q_{t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) = 0, \end{aligned}$$

where the first row follows from Eq. (59) and the second row follows from Lemma. 1. This automatically makes

$$\left| \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right| \leq \frac{1}{2} < 1.$$

Under this condition, we have

$$\begin{aligned} \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} &= \ln \left[1 + \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \\ &= \sum_{i=1}^{\infty} \frac{(-1)^{i+1}}{i} \cdot \left[\frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right]^i, \end{aligned}$$

which implies (with the dominated convergence theorem)

$$\begin{aligned} & \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \\ &= \sum_{i=1}^{\infty} \frac{(-1)^{i+1}}{i} \cdot \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)}{\Delta t} \\ & \quad \cdot \lim_{\Delta t \rightarrow 0} \frac{\left(\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right) \right)^{i-1}}{\left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)^i}. \end{aligned}$$

Only when $i = 1$, we have

$$\lim_{\Delta t \rightarrow 0} \frac{\left(\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right) \right)^{i-1}}{\left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)^i} = 1,$$

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

$$\begin{aligned} & \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \\ &= \lim_{\Delta t \rightarrow 0} \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) - q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right)}{\Delta t} \\ &= \sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right) = \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}), \end{aligned}$$

where the second equation follows from Eq. (59) and the second row follows from Lemma. 1. Hence, the proof is completed. \square

Theorem 3 (The convergence of Alg. 1). *Suppose Assumption [A1] and [A2] hold, if Alg. 1 has*

$t_0 = 0$, $t_W = T - \delta$, and $\epsilon_{score} \leq T^{-1/2} \cdot \epsilon$ where $T = \ln(4d/\epsilon^2)$ and $\delta \leq d^{-1}\epsilon$, the TV distance between the target discrete distribution q_ and the underlying distribution of the output particle $\hat{q}_{T-\delta}$ will satisfy $\text{TV}(q_*, \hat{q}_{T-\delta}) \leq 2\epsilon$.*

Proof. Here we provide the upper bound of TV distance accumulation in a specific segment, e.g., from t_{w-1} to t_w . According to the chain rule of KL divergence, i.e., Lemma 15, we have

$$\begin{aligned} \text{KL}(q_{t_w}^{\leftarrow} \| \hat{q}_{t_w}) &\leq \text{KL}(q_{t_{w-1}}^{\leftarrow} \| \hat{q}_{t_{w-1}}) + \mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL}(q_{t_w|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t_w|t_{w-1}}(\cdot|\mathbf{y}_0)) \right] \\ &= \text{KL}(q_{t_{w-1}}^{\leftarrow} \| \hat{q}_{t_{w-1}}) + \int_{t_{w-1}}^{t_w} d\mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL}(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_0)) \right] \end{aligned} \quad (60)$$

Then, it has

$$\begin{aligned} & d\mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL}(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_0)) \right] / dt \\ &= \lim_{\Delta t \rightarrow 0} (\Delta t)^{-1} \cdot \mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL}(q_{t+\Delta t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t+\Delta t|t_{w-1}}(\cdot|\mathbf{y}_0)) - \text{KL}(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_0)) \right] \\ &\leq \lim_{\Delta t \rightarrow 0} (\Delta t)^{-1} \cdot \mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\mathbb{E}_{\mathbf{y} \sim q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0)} \left(\text{KL}(q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}, \mathbf{y}_0) \| \hat{q}_{t+\Delta t|t, t_{w-1}}(\cdot|\mathbf{y}, \mathbf{y}_0)) \right) \right] \end{aligned}$$

where the inequality follows from the chain rule of the KL divergence, i.e., Lemma 15. Then, it has

$$\begin{aligned} & d\mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL}(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_0)) \right] / dt \\ &\leq \sum_{\mathbf{y} \in \mathcal{Y}, \mathbf{y}_0 \in \mathcal{Y} \rightarrow (\mathbf{y})} q_{t, t_{w-1}}^{\leftarrow}(\mathbf{y}, \mathbf{y}_0) \cdot \underbrace{\lim_{\Delta t \rightarrow 0} \left[\frac{\text{KL}(q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}, \mathbf{y}_0) \| \hat{q}_{t+\Delta t|t, t_{w-1}}(\cdot|\mathbf{y}, \mathbf{y}_0))}{\Delta t} \right]}_{\text{Term 1}} \end{aligned} \quad (61)$$

where the inequality and the notation $\mathcal{Y}^{\rightarrow}(\cdot)$ follows from Lemma 10. For each $\mathbf{y}^{\leftarrow} \in \mathcal{Y}$, $\mathbf{y}_0^{\leftarrow} \in \mathcal{Y}^{\rightarrow}(\mathbf{y})$, we focus on Term 1 of Eq. (61), and have

$$\begin{aligned}
 \text{Term 1} &= \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \cdot \ln \frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \\
 &= \lim_{\Delta t \rightarrow 0} \underbrace{\left[\sum_{\mathbf{y}' \neq \mathbf{y}} \frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right]}_{\text{Term 1.1}} + \\
 &\quad \underbrace{\lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right) \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right]}_{\text{Term 1.2}}. \tag{62}
 \end{aligned}$$

For Term 1.1, we have

$$\begin{aligned}
 \text{Term 1.1} &= \sum_{\mathbf{y}' \neq \mathbf{y}} \lim_{\Delta t \rightarrow 0} \left[\frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\Delta t} \right] \cdot \lim_{\Delta t \rightarrow 0} \left[\ln \frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \\
 &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \left[\lim_{\Delta t \rightarrow 0} \left(\frac{q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{\Delta t} \cdot \frac{\Delta t}{\hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right) \right] \\
 &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y})}, \tag{63}
 \end{aligned}$$

where the last equation follows from Lemma 1 and Eq. (59). For Term 1.2, we have

$$\begin{aligned}
 \text{Term 1.2} &= \lim_{\Delta t \rightarrow 0} \left[1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0) \right] \\
 &\quad \cdot \lim_{\Delta t \rightarrow 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t, t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t, t_{w-1}}(\mathbf{y}'|\mathbf{y}, \mathbf{y}_0)} \right] \leq 1 \cdot (\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y})) \tag{64}
 \end{aligned}$$

where the first inequality follows from Lemma 11. Plugging Eq. (63), Eq. (64) and Eq. (62), into Eq. (61) we have

$$\begin{aligned}
 &\mathbb{E}_{\mathbf{y}_0 \sim q_{t_{w-1}}^{\leftarrow}} \left[\text{KL} \left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_0) \parallel \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_0) \right) \right] / dt \\
 &\leq \sum_{\mathbf{y} \in \mathcal{Y}, \mathbf{y}_0 \in \mathcal{Y}^{\rightarrow}(\mathbf{y})} q_{t, t_{w-1}}^{\leftarrow}(\mathbf{y}, \mathbf{y}_0) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y})} + \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \right). \tag{65}
 \end{aligned}$$

Then, for any $\mathbf{y} \in \mathcal{Y}$ and $\mathbf{y}_0 \in \mathcal{Y}^{\rightarrow}(\mathbf{y})$, we have

$$\begin{aligned}
 &\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y})} + \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \\
 &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \tilde{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \\
 &\quad + \underbrace{\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{\tilde{R}_t(\mathbf{y}', \mathbf{y})}{\hat{R}_{t, \mathbf{y}_0}(\mathbf{y}', \mathbf{y})} + \hat{R}_{t, \mathbf{y}_0}(\mathbf{y}) - \tilde{R}_t(\mathbf{y})}_{\text{Term 2}}. \tag{66}
 \end{aligned}$$

When $\tilde{R}_t(\mathbf{y}) \leq \beta_{t_w}(\mathbf{y}_0)$, due to Eq. (18), we have

$$\hat{R}_{t,\mathbf{y}_0}(\mathbf{y}', \mathbf{y}) = \tilde{R}_t(\mathbf{y}', \mathbf{y}) \quad \text{and} \quad \hat{R}_{t,\mathbf{y}_0}(\mathbf{y}) = \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{R}_{t,\mathbf{y}_0}(\mathbf{y}', \mathbf{y}) = \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_t(\mathbf{y}', \mathbf{y}) = \tilde{R}_t(\mathbf{y})$$

which implies Term 2 = 0 in Eq. (66). Otherwise, we have

$$\frac{\hat{R}_{t,\mathbf{y}_0}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} = \frac{\beta_{t_w}(\mathbf{y}_0)}{\tilde{R}_t(\mathbf{y})} \quad \text{and} \quad \frac{\hat{R}_{t,\mathbf{y}_0}(\mathbf{y})}{\tilde{R}_t(\mathbf{y})} = \frac{\beta_{t_w}(\mathbf{y}_0)}{\tilde{R}_t(\mathbf{y})},$$

which implies

$$\begin{aligned} \text{Term 2} &= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{\tilde{R}_t(\mathbf{y})}{\beta_{t_w}(\mathbf{y}_0)} + \beta_{t_w}(\mathbf{y}_0) - \tilde{R}_t(\mathbf{y}) \\ &= R_t^{\leftarrow}(\mathbf{y}) \cdot \ln \left[1 + \frac{\tilde{R}_t(\mathbf{y}) - \beta_{t_w}(\mathbf{y}_0)}{\beta_{t_w}(\mathbf{y}_0)} \right] + \beta_{t_w}(\mathbf{y}_0) - \tilde{R}_t(\mathbf{y}) \\ &\leq \beta_{t_w}(\mathbf{y}_0) \cdot \left[\frac{\tilde{R}_t(\mathbf{y}) - \beta_{t_w}(\mathbf{y}_0)}{\beta_{t_w}(\mathbf{y}_0)} \right] + \beta_{t_w}(\mathbf{y}_0) - \tilde{R}_t(\mathbf{y}) = 0, \end{aligned}$$

where the last inequality follows from

$$\mathbf{y}_0 \in \mathcal{Y}^{\rightarrow}(\mathbf{y}) \Rightarrow \text{numK}(\mathbf{y}) \leq \text{numK}(\mathbf{y}_0) \Rightarrow R_t^{\leftarrow}(\mathbf{y}) \leq \beta_{t_w}(\mathbf{y}_0).$$

Combining with Eq. (66) and Eq. (65), we have

$$\begin{aligned} &\text{d}\mathbb{E}_{\mathbf{y}_0 \sim q_{t_w-1}^{\leftarrow}} \left[\text{KL} \left(q_{t_w-1}^{\leftarrow}(\cdot | \mathbf{y}_0) \| \hat{q}_{t_w-1}(\cdot | \mathbf{y}_0) \right) \right] / \text{d}t \\ &\leq \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \tilde{R}_t(\mathbf{y}) - R_t^{\leftarrow}(\mathbf{y}) \right) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_t(\mathbf{y}', \mathbf{y})} + \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_t(\mathbf{y}', \mathbf{y}) - \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right) \quad (67) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot \left[-\frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} + \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') + \frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} \ln \frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y}) \tilde{v}_{t,\mathbf{y}}(\mathbf{y}')} \right] \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} q_t^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') D_{\phi} \left(\frac{q_t^{\leftarrow}(\mathbf{y}')}{q_t^{\leftarrow}(\mathbf{y})} \| \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') \right), \end{aligned}$$

where D_{ϕ} is the Bregman divergence with $\phi(c) = c \ln c$ (as Eq. (6)), and the last equation follows from the definition of Bregman divergence:

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

Therefore, Eq. (60) can be rewritten as

$$\text{KL} \left(q_{t_w}^{\leftarrow} \| \hat{q}_{t_w} \right) \leq \text{KL} \left(q_{t_w-1}^{\leftarrow} \| \hat{q}_{t_w-1} \right) + \int_{t_w-1}^{t_w} \mathbb{E}_{\mathbf{y} \sim q_{T-t}^{\rightarrow}} \left[\sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot D_{\phi} \left(v_{t,\mathbf{y}}(\mathbf{y}') \| \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') \right) \right] \text{d}t.$$

With a recursive manner, we have

$$\text{KL} \left(q_{T-\delta}^{\leftarrow} \| \hat{q}_{T-\delta} \right) \leq \text{KL} \left(q_0^{\leftarrow} \| \hat{q}_0 \right) + L_{\text{SE}}(\tilde{v}) = \text{KL} \left(q_T^{\rightarrow} \| \hat{q}_0 \right) + L_{\text{SE}}(\tilde{v}) \leq (1 + e^{-T})^d - 1 + T\epsilon_{\text{score}}^2,$$

where the last inequality follows from Lemma 2 and Assumption [A1]

$$\hat{q}_0(\mathbf{y}) = \tilde{q}_T(\mathbf{y}) \propto \exp(-T \cdot (d - \text{numK}(\mathbf{y}))).$$

If we set

$$T \geq \ln(4d/\epsilon^2) \quad \text{and} \quad \epsilon_{\text{score}} \leq T^{-1/2} \cdot \epsilon,$$

it has $(1 + e^{-T})^d - 1 \leq \epsilon^2$ and $T\epsilon_{\text{score}}^2 \leq \epsilon^2$, which means $\text{KL} \left(q_{T-\delta}^{\leftarrow} \| \hat{q}_{T-\delta} \right) \leq 2\epsilon^2$.

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

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

where $\Gamma(q_*, q_\delta^\rightarrow)$ is the set of all couplings of $(q_*, q_\delta^\rightarrow)$, and the inequality holds because $(\mathbf{y}, \mathbf{y}')$ gives a coupling of $(q_*, q_\delta^\rightarrow)$. Without loss of generality, we suppose $q_0^\rightarrow(\mathbf{y}) > 0$ for all $\text{numK}(\mathbf{y}) = 0$, then, combining the transition kernel given Lemma 8 and Assumption [A2], we have

$$\mathbb{P}(\mathbf{y} = \mathbf{y}') = \sum_{\mathbf{y} \in \mathcal{Y}, \text{numK}(\mathbf{y})=0} q_0^\rightarrow[\mathbf{y}] \cdot q_{\delta|0}^\rightarrow(\mathbf{y}|\mathbf{y}) = \sum_{\mathbf{y} \in \mathcal{Y}, \text{numK}(\mathbf{y})=0} q_0^\rightarrow(\mathbf{y}) \cdot e^{-\delta d} = e^{-\delta d}.$$

Thus, by choosing $\delta \leq \epsilon/d$, we have

$$\delta \leq d^{-1}\epsilon \leq d^{-1} \cdot \ln\left(\frac{1}{1-\epsilon}\right) \Rightarrow e^{\delta d} \leq \frac{1}{1-\epsilon} \Rightarrow \text{TV}(q_*, q_\delta^\rightarrow) \leq 1 - e^{-\delta d} \leq \epsilon. \quad (68)$$

Finally, we have

$$\text{TV}(q_0^\rightarrow, \hat{q}_{T-\delta}^\leftarrow) \leq \text{TV}(q_0^\rightarrow, q_\delta^\rightarrow) + \text{TV}(q_{T-\delta}^\leftarrow, \hat{q}_{T-\delta}^\leftarrow) \leq \epsilon + \sqrt{\frac{1}{2} \text{KL}(q_{T-\delta}^\leftarrow \parallel \hat{q}_{T-\delta}^\leftarrow)} \leq 2\epsilon.$$

Hence the proof is completed. \square

D.3 THE COMPLEXITY OF ALG. 1

Theorem 4 (The complexity of Alg. 1). *Suppose Assumption [A1] and [A2] hold, following from the settings shown in Theorem 3, if we implement Alg. 1 with*

$t_w - t_{w-1} = \eta$ where $w \in \{1, 2, \dots, W\}$, $W = (T - \delta)/\eta$, $\eta = \epsilon/2d$, and $\epsilon < 1$ the expectation of iteration/score estimation complexity of Alg. 1 will be upper bounded by

$$2K(d - \epsilon^2/4) + 12Kd \ln d$$

to achieve $\text{TV}(q_, \hat{q}) \leq 2\epsilon$ where \hat{p} denotes the underlying distribution of generated samples.*

Proof. We denote $\{\hat{\mathbf{y}}_t\}_{t=0}^{T-\delta}$ to present the reverse process. For a specific trajectory, e.g., $\{\hat{\mathbf{y}}_t\}_{t=0}^{T-\delta} = \{\hat{\mathbf{y}}\}_{t=0}^{T-\delta}$, the total expected iteration number will be equivalent to the summation of Poisson expectations of W segments, i.e.,

$$\sum_{i=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) = \frac{K \cdot \text{numK}(\hat{\mathbf{y}}_{t_{w-1}})}{e^{T-t_w} - 1} \cdot (t_w - t_{w-1}),$$

which means the expected iteration number of the reverse process can be written as

$$\mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] = \sum_{w=1}^W \mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] \cdot \frac{K}{e^{(T-t_w)} - 1} \cdot (t_w - t_{w-1}). \quad (69)$$

Although $\mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})]$ is respect to the practical distribution $\hat{\mathbf{y}}_{t_{w-1}} \sim \hat{q}_{t_{w-1}}$, we can approximate it by the forward marginal distribution, i.e.,

$$\mathbb{E}[\text{numK}(\mathbf{y}_{t_{w-1}}^\leftarrow)] = \mathbb{E}[\text{numK}(\mathbf{y}_{T-t_{w-1}}^\rightarrow)] \quad \text{where} \quad \mathbf{y}_{t_{w-1}}^\leftarrow \sim q_{t_{w-1}}^\leftarrow \quad \text{and} \quad \mathbf{y}_{T-t_{w-1}}^\rightarrow \sim q_{t_{w-1}}^\rightarrow$$

Specifically, with Assumption [A2], we have $\mathbb{E}[\text{numK}(\mathbf{y}_0^\rightarrow)] = 0$. Under this condition, the transition kernel becomes

$$q_{t|0}^\rightarrow(\mathbf{y}|\mathbf{y}') = \prod_{i=1}^d \left[\underbrace{(1 - \mathbf{1}_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \mathbf{1}_0(\mathbf{y}_i - \mathbf{y}'_i) \cdot e^{-t}}_{\text{remain non-mask token}} + \underbrace{(1 - \mathbf{1}_{(K,K)}(\mathbf{y}_i, \mathbf{y}'_i)) \cdot \mathbf{1}_K(\mathbf{y}_i) \cdot (1 - e^{-t})}_{\text{turn into mask token}} \right].$$

due to Lemma 8. Let $\mathbb{P}[\text{numK}(\mathbf{y}_t^\rightarrow) = k]$ be the probability that exactly k out of the d coordinates are mask tokens (K) at time t . Because each of the d coordinates evolves independently (and identically, each with probability $1 - e^{-t}$ of being the mask token at time t), we get a standard Binomial random variable:

- Each coordinate is K with probability $1 - e^{-t}$.
- Each coordinate is non- K with probability e^{-t} .

Hence, we have

$$\mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k] = C_d^k \cdot (1 - e^{-t})^k \cdot (e^{-t})^{d-k} \quad \text{and} \quad \mathbb{E}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k] = d \cdot (1 - e^{-t}).$$

Then, for any w , we have $\mathbb{E}[\text{numK}(\mathbf{y}_{T-t_{w-1}}^{\rightarrow})] = d \cdot (1 - e^{-(T-t_{w-1})})$. Under the settings shown in Theorem 3, we have

$$\text{TV}(\mathbf{q}_{T-t_{w-1}}^{\rightarrow}, \hat{\mathbf{q}}_{t_{w-1}}) \leq \text{TV}(\mathbf{q}_{T-t_W}^{\rightarrow}, \hat{\mathbf{q}}_{t_W}) \leq 2\epsilon,$$

which implies

$$\left| \mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] - \mathbb{E}[\text{numK}(\mathbf{y}_{t_{w-1}}^{\leftarrow})] \right| \leq d \cdot \text{TV}(\mathbf{q}_{T-t_{w-1}}^{\rightarrow}, \hat{\mathbf{q}}_{t_{w-1}}) \leq 2d\epsilon.$$

Then, we have

$$\begin{aligned} \mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] &\leq \mathbb{E}[\text{numK}(\mathbf{y}_{t_{w-1}}^{\leftarrow})] + 2d\epsilon = \mathbb{E}[\text{numK}(\mathbf{y}_{T-t_{w-1}}^{\rightarrow})] + 2d\epsilon \\ &= d \cdot (1 - e^{-(T-t_{w-1})}) + 2d\epsilon. \end{aligned} \quad (70)$$

Plugging Eq. (70) into Eq. (69), we have

$$\begin{aligned} &\mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] \\ &\leq \sum_{w=1}^W d \cdot (1 - e^{-(T-t_{w-1})}) \cdot \frac{K}{e^{(T-t_w)} - 1} \cdot (t_w - t_{w-1}) \\ &\quad + \sum_{w=1}^W 2d\epsilon \cdot \frac{K}{e^{(T-t_w)} - 1} \cdot (t_w - t_{w-1}) \\ &= \underbrace{Kd \cdot \sum_{w=1}^W e^{-(T-t_w)} \cdot (t_w - t_{w-1}) \cdot \frac{1 - e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}}}_{\text{Term 1}} \\ &\quad + \underbrace{2Kd\epsilon \cdot \sum_{w=1}^W (e^{T-t_w} - 1)^{-1} \cdot (t_w - t_{w-1})}_{\text{Term 2}} \end{aligned} \quad (71)$$

Then, we suppose the segments share the same length η , i.e.,

$$t_w - t_{w-1} = \eta \quad \text{where} \quad w \in \{1, 2, \dots, W\}, \quad W = (T - \delta)/\eta, \quad \text{and} \quad \eta = \epsilon/2d.$$

Under these conditions, we have

$$\begin{aligned} \eta \leq \frac{\delta}{2} &\leq \ln\left(\frac{1}{2} + \frac{e^\delta}{2}\right) \Rightarrow e^\eta \leq \frac{e^\delta}{2} + \frac{1}{2} \leq \frac{e^{(T-t_{w-1})}}{2} + \frac{1}{2} \quad \forall w \in \{1, \dots, W\} \\ \Rightarrow e^\eta &\leq \frac{1 + e^{-(T-t_{w-1})}}{2e^{-(T-t_{w-1})}} \Rightarrow 2 \cdot e^{-(T-t_{w-1}-\eta)} \leq 1 + e^{-(T-t_{w-1})} \\ \Rightarrow 1 - e^{-(T-t_{w-1})} &\leq 2 - 2e^{-(T-t_{w-1}-\eta)} \Rightarrow \frac{1 - e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}} \leq 2. \end{aligned} \quad (72)$$

Plugging these results into Term 1 of Eq. (71), we have

$$\begin{aligned} \text{Term 1} &= 2Kd \cdot \sum_{w=1}^W e^{-(T-t_w)} \cdot \eta = 2Kd \cdot \sum_{w=1}^W e^{-(T-w\eta)} \cdot \eta \\ &= 2Kd \cdot \eta \cdot e^{-T} \cdot \frac{e^{(W+1)\eta} - e^\eta}{e^\eta - 1} \leq 2Kd \cdot e^\eta \cdot (e^{-\delta} - e^{-T}) \leq 2Kd \cdot (1 - e^{-T}) \\ &= 2Kd \cdot \left(1 - \frac{\epsilon^2}{4d}\right). \end{aligned} \quad (73)$$

Moreover, we have

$$\frac{e^{T-t_{w-1}} - 1}{e^{T-t_w} - 1} = \frac{e^{T-t_{w-1}}}{e^{T-t_w}} \cdot \frac{1 - e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}} \leq e^\eta \cdot 2 \leq 2e,$$

where the first inequality follows from Eq. (72) and the last inequality is established when $\eta \leq 1$. Then, Term 2 of Eq. (71) can be upper bounded as

$$\begin{aligned} \text{Term 2} &= 2Kd\epsilon \cdot \sum_{w=1}^W \frac{\eta}{e^{T-t_w} - 1} \leq 4e \cdot Kd\epsilon \cdot \sum_{w=1}^W \frac{\eta}{e^{T-t_{w-1}} - 1} \leq 4e \cdot Kd\epsilon \cdot \sum_{w=1}^W \frac{\eta}{T - t_{w-1}} \\ &\leq 4e \cdot dK\epsilon \cdot \int_0^{T-\delta} \frac{1}{T-t} dt = 4e \cdot dK\epsilon \cdot \ln \frac{T}{\delta} \leq 4e \cdot dK\epsilon \cdot \ln \frac{4d^2}{\epsilon^3} \leq 12e \cdot Kd \ln d \cdot \epsilon \ln \frac{1}{\epsilon} \end{aligned} \quad (74)$$

where the last inequality follows from

$$4 \leq d \quad \text{and} \quad \ln \frac{d^3}{\epsilon^3} = 3 \ln \frac{d}{\epsilon} \leq 3 \ln d \ln \frac{1}{\epsilon}$$

without loss of generality. Moreover, when $\epsilon < 1$, we have

$$\epsilon \ln \frac{1}{\epsilon} \leq e^{-1},$$

which follows from the monotonicity of the function $x \ln x$. Under this condition, the RHS of Eq. (74) has the following bound

$$\text{Term 2} \leq 12 \cdot Kd \ln d. \quad (75)$$

Finally, plugging Eq. (73) and Eq. (75) into Eq. (71), the expected calls of discrete scores will be

$$\mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] \leq 2K(d - \epsilon^2) + 12Kd \ln d.$$

Hence, the proof is completed. \square

Corollary 5. Suppose Assumption [AI] hold, following from the settings shown in Theorem 3, if we implement Alg. 1 with

$t_w - t_{w-1} = \eta$ where $w \in \{1, 2, \dots, W\}$, $W = (T - \delta)/\eta$, $\eta = \epsilon/2d$, and $\epsilon < 1$ the expectation of iteration/score estimation complexity of Alg. 1 will be upper bounded by

$$\min \left\{ O(Kd \ln(d/\epsilon)), O \left(Kd \cdot \frac{\mathbb{E}[\text{numK}(\mathbf{y}_0^{\rightarrow})]}{\epsilon} \right) \right\} + O(Kd \ln d)$$

to achieve $\text{TV}(q_*, \hat{q}) \leq 2\epsilon$ where \hat{p} denotes the underlying distribution of generated samples.

Proof. Similar to the proof shown in Theorem 4, the expected iteration number of the reverse process can be written as

$$\mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] = \sum_{w=1}^W \mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] \cdot \frac{K}{e^{(T-t_w)} - 1} \cdot (t_w - t_{w-1}).$$

Although $\mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})]$ is respect to the practical distribution $\hat{\mathbf{y}}_{t_{w-1}} \sim \hat{q}_{t_{w-1}}$, we can approximate it by the forward marginal distribution, i.e.,

$$\mathbb{E}[\text{numK}(\mathbf{y}_{t_{w-1}}^{\leftarrow})] = \mathbb{E}[\text{numK}(\mathbf{y}_{T-t_{w-1}}^{\rightarrow})] \quad \text{where} \quad \mathbf{y}_{t_{w-1}}^{\leftarrow} \sim q_{t_{w-1}}^{\leftarrow} \quad \text{and} \quad \mathbf{y}_{T-t_{w-1}}^{\rightarrow} \sim q_{t_{w-1}}^{\rightarrow}.$$

Let $\mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k]$ be the probability that exactly k out of the d coordinates are mask tokens (K) at time t presented as

$$\mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k] = \sum_{i=0}^k \mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k | \text{numK}(\mathbf{y}_0^{\rightarrow}) = i] \cdot \Pr[\text{numK}(\mathbf{y}_0^{\rightarrow}) = i].$$

Because each of the d coordinates evolves independently (and identically, each with probability $1 - e^{-t}$ of being the mask token at time t), we get a standard Binomial random variable:

- Each coordinate is K with probability $1 - e^{-t}$.
- Each coordinate is non- K with probability e^{-t} .

Hence, we have

$$\mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k | \text{numK}(\mathbf{y}_0^{\rightarrow}) = i] = \underbrace{C_{d-i}^{k-i}}_{\text{unmask} \rightarrow \text{mask count}} \cdot \underbrace{(1 - e^{-t})^{k-i}}_{\text{prob of mask transition}} \cdot \underbrace{(e^{-t})^{d-k}}_{\text{prob of unmask kept}}.$$

Under this condition, the expected number of MASK token at forward time t will become

$$\begin{aligned} \mathbb{E}[\text{numK}(\mathbf{y}_t^{\rightarrow})] &= \sum_{k=0}^d k \cdot \mathbb{P}[\text{numK}(\mathbf{y}_t^{\rightarrow}) = k] \\ &= \sum_{k=0}^d \sum_{i=0}^k C_{d-i}^{k-i} \cdot (1 - e^{-t})^{k-i} \cdot (e^{-t})^{d-k} \cdot \mathbb{P}[\text{numK}(\mathbf{y}_0^{\rightarrow}) = i] \\ &= \sum_{i=0}^d \mathbb{P}[\text{numK}(\mathbf{y}_0^{\rightarrow}) = i] \cdot \underbrace{\sum_{k=i}^d C_{d-i}^{k-i} \cdot (1 - e^{-t})^{k-i} \cdot (e^{-t})^{d-k}}_{\text{Term 1}}. \end{aligned} \quad (76)$$

For Term 1 in Eq. (76), suppose $j = k - i$, we have

$$\begin{aligned} \text{Term 1} &= \sum_{j=0}^{d-i} (j + i) \cdot C_{d-i}^j \cdot (1 - e^{-t})^j \cdot (e^{-t})^{(d-i-j)} \\ &= \sum_{j=0}^{d-i} j \cdot C_{d-i}^j \cdot (1 - e^{-t})^j \cdot (e^{-t})^{(d-i-j)} + i \cdot \sum_{j=0}^{d-i} C_{d-i}^j \cdot (1 - e^{-t})^j \cdot (e^{-t})^{(d-i-j)} \\ &= \sum_{j=0}^{d-i} j \cdot C_{d-i}^j \cdot (1 - e^{-t})^j \cdot (e^{-t})^{(d-i-j)} + i \cdot (1 - e^{-t} + e^{-t})^{d-i} = d - (d - i)e^{-t} \end{aligned}$$

where the last equation follows from the expectation of binomial distributions. Then, Eq. (76) can be written as

$$\begin{aligned} \mathbb{E}[\text{numK}(\mathbf{y}_t^{\rightarrow})] &= \sum_{i=0}^d \mathbb{P}[\text{numK}(\mathbf{y}_0^{\rightarrow}) = i] \cdot (d - d \cdot e^{-t} + i \cdot e^{-t}) \\ &= d \cdot (1 - (1 - \mathbb{E}[\text{numK}(\mathbf{y}_0^{\rightarrow})]/d) \cdot e^{-t}). \end{aligned}$$

Without loss of generality, we suppose

$$r_0 := 1 - \mathbb{E}[\text{numK}(\mathbf{y}_0^{\rightarrow})]/d > 0.$$

Then, following from Eq. (70), we have

$$\mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] \leq d \cdot (1 - r_0 \cdot e^{-(T-t_{w-1})}) + 2d\epsilon,$$

and

$$\begin{aligned} \mathbb{E} \left[\sum_{w=1}^W \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1}) \right] &\leq \underbrace{Kd \cdot \sum_{w=1}^W e^{-(T-t_w)} \cdot (t_w - t_{w-1}) \cdot \frac{1 - r_0 \cdot e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}}}_{\text{Term 1}} \\ &\quad + \underbrace{2Kd\epsilon \cdot \sum_{w=1}^W (e^{T-t_w} - 1)^{-1} \cdot (t_w - t_{w-1})}_{\text{Term 2}} \end{aligned} \quad (77)$$

Here the second term can be upper bounded as

$$\text{Term 2} \leq 12 \cdot Kd \ln d$$

by choosing the mixing time T and early stopping time δ as Theorem 3, which follows from Eq. (74). For Term 1 of Eq. (77), we will discuss it in categories. Suppose the expected number of mask token satisfies

$$\mathbb{E}[\text{numK}(\mathbf{y}_0^\rightarrow)] \leq C_0 \cdot \epsilon \Leftrightarrow r_0 \geq 1 - C_0 \cdot \epsilon/d,$$

and the segments share the same length η , i.e.,

$$t_w - t_{w-1} = \eta \quad \text{where} \quad w \in \{1, 2, \dots, W\}, \quad W = (T - \delta)/\eta, \quad \text{and} \quad \eta = \epsilon/2d,$$

following from Eq. (72), we have

$$\eta \leq \frac{\delta}{2} \Rightarrow \frac{1 - e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}} \leq 2.$$

Combining with the following fact, i.e.,

$$\begin{aligned} \frac{(1 - r_0) \cdot e^{-(T-t_{w-1})}}{1 - e^{-(T-t_w)}} &= \frac{1 - r_0}{e^{(T-t_{w-1})} - e^\eta} = (1 - r_0) \cdot e^{-\eta} \cdot (e^{T-t_{w-1}-\eta} - 1)^{-1} \\ &\leq (1 - r_0) \cdot \delta^{-1} = C_0, \end{aligned}$$

Eq. (73) demonstrates that

$$\text{Term 1} \leq (2 + C_0) \cdot Kd \cdot \left(1 - \frac{\epsilon^2}{4d}\right).$$

On the other hand, we have

$$\begin{aligned} \text{Term 1} &\leq Kd \cdot \sum_{w=1}^W \frac{\eta}{e^{T-t_w} - 1} \leq Kd \sum_{w=1}^W \frac{\eta}{T - t_w} \leq 1.5Kd \sum_{w=1}^W \frac{\eta}{T - t_{w-1}} \\ &\lesssim 1.5Kd \cdot \int_{\delta}^1 t^{-1} dt \leq 1.5Kd \ln(1/\delta) = 1.5Kd \ln(d/\epsilon), \end{aligned}$$

where the forth inequality follows from the choice of η , i.e.,

$$\eta \leq \delta/2 \Rightarrow (T - t_{w-1}) - (T - t_w) = \eta \leq \frac{\delta}{2} \leq \frac{T - t_w}{2}.$$

Hence, the total complexity will be

$$\min \left\{ O(Kd \ln(d/\epsilon)), O \left(Kd \cdot \frac{\mathbb{E}[\text{numK}(\mathbf{y}_0^\rightarrow)]}{\epsilon} \right) \right\} + O(Kd \ln d).$$

Hence, the proof is completed. \square

Corollary 6. Suppose Assumption [A1] and [A2] hold, if Alg. 1 has

$t_0 = 0, \quad t_W = T - \delta, \quad \text{and} \quad \epsilon_{\text{score}} \leq T^{-1/2} \cdot \epsilon \quad \text{where} \quad T = \ln(4d/\epsilon^2) \quad \text{and} \quad \delta \leq d^{-1}\epsilon,$
and draw initial $\bar{\mathbf{y}}_0 \sim \delta_{[K, \dots, K]}(\cdot)$, the TV distance between the target discrete distribution q_* and the underlying distribution of the output particle $\bar{q}_{T-\delta}$ will satisfy $\text{TV}(q_*, \bar{q}_{T-\delta}) \leq 2.5\epsilon$.

Proof. We consider a stochastic process $\{\bar{\mathbf{y}}_t\}_{t=0}^{T-\delta}$ which satisfies $\bar{\mathbf{y}}_t \sim \bar{q}_t$. The initial distribution is $\bar{\mathbf{y}}_0 \sim \bar{q}_0 = \delta_{[K, K, \dots, K]}(\mathbf{y})$. Suppose the joint and conditional distribution are

$$(\bar{\mathbf{y}}_{t'}, \bar{\mathbf{y}}_t) \sim \bar{q}_{t',t} \quad \text{and} \quad \bar{q}_{t|t'}(\bar{\mathbf{y}}_t | \bar{\mathbf{y}}_{t'}) = \bar{q}_{t,t'}(\bar{\mathbf{y}}_t, \bar{\mathbf{y}}_{t'}) / \bar{q}_{t'}(\bar{\mathbf{y}}_{t'}) \quad \text{where} \quad t > t'.$$

Specifically, we suppose the random variables $\{\bar{\mathbf{y}}_t\}_{t=0}^{T-\delta}$ share the same transition as that in $\{\hat{\mathbf{y}}_t\}_{t=0}^{T-\delta}$ shown in Theorem 3, which means $\bar{q}_{t|t'} = \hat{q}_{t|t'}$ for any $t > t'$, which implies $\{\bar{\mathbf{y}}_t\}_{t=0}^{T-\delta}$ can be implemented by

1. Initialize the particles as $\bar{\mathbf{y}}_0 \sim \bar{q}_0 = \delta_{[K, K, \dots, K]}(\mathbf{y})$
2. Update $\{\bar{\mathbf{y}}_t\}_{t>0}^{T-\delta}$ with Alg. 1

Then, due to the chain rule of TV distance, i.e., Lemma 16, we have

$$\begin{aligned} \text{TV}(\hat{q}_{T-\delta}, \bar{q}_{T-\delta}) &\leq \text{TV}(\hat{q}_{T-\delta,0}, \bar{q}_{T-\delta,0}) \\ &\leq \text{TV}(\hat{q}_0, \bar{q}_0) + \mathbb{E}_{\hat{y}_0 \sim \hat{q}_0} \left[\text{TV}(\hat{q}_{T-\delta|0}, \bar{q}_{T-\delta|0}) \right] = \text{TV}(\hat{q}_0, \bar{q}_0). \end{aligned} \quad (78)$$

Since \bar{q}_0 is the mask token dirac measure, we have

$$\text{TV}(\hat{q}_0, \bar{q}_0) = 1 - \hat{q}_0([K, K, \dots, K]).$$

According to the proof of Lemma 2, we can easily find that

$$\hat{q}_0([K, \dots, K]) = \frac{1}{(1 + e^{-T})^d}.$$

By requiring $T \geq \ln(4d/\epsilon)$ and $\epsilon \leq 1$, we have

$$\begin{aligned} T \geq \ln(4d/\epsilon) &\Rightarrow t \geq \ln(d/\ln(1 + \epsilon/2)) \Leftrightarrow d \cdot e^{-T} \leq \ln(1 + \epsilon/2) \\ &\Rightarrow d \ln(1 + e^{-T}) \leq \ln(1 + \epsilon/2) \Leftrightarrow (1 + e^{-T})^d - 1 \leq \epsilon/2. \end{aligned}$$

That means

$$\text{TV}(\hat{q}_0, \bar{q}_0) = 1 - \frac{1}{(1 + e^{-T})^d} \leq 1 - \frac{1}{1 + \epsilon/2} \leq \epsilon/2.$$

Plugging this inequality into Eq. 78, we have $\text{TV}(\hat{q}_{T-\delta}, \bar{q}_{T-\delta}) \leq \epsilon/2$. Then combining it with Theorem 3, i.e., $\text{TV}(q_*, \hat{q}_{T-\delta}, \leq) 2\epsilon$, we have $\text{TV}(q_*, \bar{q}_{T-\delta}, \leq) 2.5\epsilon$. Hence, the proof is completed. \square

E TECHNICAL LEMMAS

Lemma 12 (Basic Kronecker product). *Suppose the Kronecker product for n matrices defined on $\mathbb{R}^{d \times d}$, i.e.,*

$$\bar{A} := A_1 \otimes A_2 \otimes \dots \otimes A_n,$$

then we have

$$\bar{A}_{[a_{1,i}, a_{2,i}, \dots, a_{n,i}], [a_{1,j}, a_{2,j}, \dots, a_{n,j}]} := \bar{A}_{\sum_{k=1}^n a_{k,i} \cdot d^{n-k}, \sum_{k=1}^n a_{k,j} \cdot d^{n-k}} = \prod_{k=1}^n [A_k]_{a_{k,i}, a_{k,j}}.$$

Proof. This lemma can easily be proved by the definition of Kronecker product. \square

Lemma 13 (Mixed-product property of Kronecker product). *Suppose the matrices $A, B, C, D \in \mathbb{R}^{d \times d}$, then, the products AC and BD are well-defined. We have*

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD).$$

Proof. We prove this by examining the product on the left-hand side, $(A \otimes B)(C \otimes D)$, and showing it coincides block-by-block with $(AC) \otimes (BD)$.

We start from the definition of Kronecker products in blocks. By definition, the Kronecker product $A \otimes B$ can be seen as an $(d \times d)$ block matrix in which the (i, j) -th block is $a_{ij} B$. Hence,

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{pmatrix}.$$

Similarly,

$$C \otimes D = \begin{pmatrix} c_{11}D & c_{12}D & \cdots & c_{1r}D \\ c_{21}D & c_{22}D & \cdots & c_{2r}D \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1}D & c_{n2}D & \cdots & c_{nr}D \end{pmatrix}.$$

Then, we form the Product $(A \otimes B)(C \otimes D)$. When multiplying two block matrices, we sum over the matching inner block dimensions. Specifically, the (i, k) -block of $(A \otimes B)(C \otimes D)$ is given by

$$\sum_{j=1}^n \left((a_{ij} B) (c_{jk} D) \right).$$

Inside each term, we treat $a_{ij} B$ and $c_{jk} D$ as scalar-matrix products. We can rewrite the expression as:

$$\sum_{j=1}^n a_{ij} c_{jk} (BD) = \left(\sum_{j=1}^n a_{ij} c_{jk} \right) BD.$$

Notice that the factor $\sum_{j=1}^n a_{ij} c_{jk}$ is precisely $(AC)_{ik}$, the (i, k) -th entry of the matrix product AC . Thus, each (i, k) -block of $(A \otimes B)(C \otimes D)$ simplifies to

$$(AC)_{ik} (BD).$$

Now observe that the Kronecker product $(AC) \otimes (BD)$ can also be viewed as an $(m \times r)$ block matrix whose (i, k) -th block is

$$(AC)_{ik} (BD).$$

Hence, the (i, k) -th block of $(AC) \otimes (BD)$ matches exactly with the (i, k) -th block we computed for $(A \otimes B)(C \otimes D)$. Since these two matrices agree in every block of a $d^2 \times d^2$ partition, we conclude

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD),$$

as desired. \square

Lemma 14 (Kolmogorov backward theorem, adapted from Theorem 5.11 in Särkkä & Solin (2019)). *For a specific SDE, if we denote the transition density from $\mathbf{x}(s)$ to $\mathbf{y}(t)$ as $p(\mathbf{y}, t | \mathbf{x}, s)$, then it solves the backward Kolmogorov equation*

$$-\frac{\partial p(\mathbf{y}, t | \mathbf{x}, s)}{\partial s} = \mathcal{L}p(\mathbf{y}, t | \mathbf{x}, s)$$

where \mathcal{L} denotes the infinitesimal operator of the SDE.

Lemma 15 (The chain rule of KL divergence). *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

$$\begin{aligned} \text{KL}(p_{x,z} \| q_{x,z}) &= \text{KL}(p_z \| q_z) + \mathbb{E}_{\mathbf{z} \sim p_z} [\text{KL}(p_{x|z}(\cdot | \mathbf{z}) \| q_{x|z}(\cdot | \mathbf{z}))] \\ &= \text{KL}(p_x \| q_x) + \mathbb{E}_{\mathbf{x} \sim p_x} [\text{KL}(p_{z|x}(\cdot | \mathbf{x}) \| q_{z|x}(\cdot | \mathbf{x}))] \end{aligned}$$

where the latter equation implies

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

Lemma 16 (The chain rule of TV distance). *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

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

Algorithm 2 FIRST HITTING SAMPLING

```

1: Input: The sequence length  $d$ , the vocabulary  $\mathcal{V} = \{1, 2, \dots, K\}$  where  $K$  is the mask token,
   the noise schedule  $\alpha_t$  and its inverse function  $\alpha^{-1}$ , the pretrained masked diffusion model  $p_\theta$ 
2:  $\tilde{\mathbf{y}}_0 = [K, K, \dots, K]$ .
3:  $\tilde{\tau}_0 = 1$ .
4: for  $n = 0$  to  $d - 1$  do
5:   Sample  $u_n \sim \text{Uniform}(0, 1)$ 
6:    $\tilde{\tau}_{n+1} = \alpha^{-1}(1 - u_n^{d-n}(1 - \alpha_{\tilde{\tau}_n}))$ 
7:   Randomly and uniformly select an index  $l$  from  $\{i : \tilde{\mathbf{y}}_n^{(i)} = K\}$  (i.e., masked positions in  $\tilde{\mathbf{y}}_n$ )

8:    $\mathbf{p}_n = p_{\theta, l}(\cdot | \tilde{\mathbf{y}}_n, \tilde{\tau}_{n+1}) \in \mathbb{R}_+^K$ 
9:    $\tilde{\mathbf{y}}_{n+1} = \tilde{\mathbf{y}}_n, \tilde{\mathbf{y}}_{n+1}^{(l)} \sim \text{Cat}(\mathbf{p}_n(\cdot, l))$ 
10: end for
11: return  $\tilde{\mathbf{y}}_d$ .

```

F FHS CONVERGENCE UNDER TIME-INDEPENDENT SCORE PARAMETERIZATION

In the following, we will prove that the distribution generated by first hitting sampling (Zheng et al., 2024) approaches to the target data distribution p_* in TV distance. The core step is to introduce our MATU as the reference probability path.

We start from some additional notations. Specifically, suppose following two elements $\mathbf{y}, \hat{\mathbf{y}} \in \mathcal{Y} = \{1, 2, \dots, K\}^d$ satisfying

$$\mathbf{y} = [\mathbf{y}_1, \dots, \mathbf{y}_i, \dots, \mathbf{y}_d] \quad \hat{\mathbf{y}} = [\mathbf{y}_1, \dots, \hat{\mathbf{y}}_i, \dots, \mathbf{y}_d],$$

which means the Hamming distance between \mathbf{y} and $\hat{\mathbf{y}}$ is 1 and they are only different at i -th coordinate. Suppose $\mathbf{y}_i = K$ and $\hat{\mathbf{y}}_i \neq K$, then we can define the conditional distribution at specific coordinate, e.g., i , given unmask tokens $\mathbf{y}_{\mathcal{K}^c(\mathbf{y})}$

$$q_{0,i}(\hat{\mathbf{y}}_i | \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}) = \frac{\sum_{\tilde{\mathbf{y}} \in \mathcal{Y}^+, \tilde{\mathbf{y}}_{\mathcal{K}^c(\mathbf{y})} = \hat{\mathbf{y}}_{\mathcal{K}^c(\mathbf{y})} q_0(\tilde{\mathbf{y}})}{\sum_{\tilde{\mathbf{y}} \in \mathcal{Y}^+, \tilde{\mathbf{y}}_{\mathcal{K}^c(\mathbf{y})} = \mathbf{y}_{\mathcal{K}^c(\mathbf{y})} q_0(\tilde{\mathbf{y}})}$$

For the completeness of the analysis, we first show the FHS in Alg. 2.

Bridge the discrete score estimation error and the pretrained masked diffusion models in FHS.

We need to note that the output of pretrained masked diffusion model satisfies

$$p_{\theta, i}(\cdot | \mathbf{y}, \tilde{\tau}_{n+1}) = p_{\theta, i}(\cdot | \mathbf{y}) \approx q_{0,i}(\hat{\mathbf{y}}_i | \mathbf{y}_{\mathcal{K}^c(\mathbf{y})}),$$

where first equation comes from the time-independent parameterization, and the second approximation comes from the training objective, i.e.,

$$\mathcal{L}_w^d(\mathbf{y}_0) = \sum_{i=1}^d \mathbf{w}_i \cdot \mathbb{E}_{\mathbb{P}_{\mathbf{y}}[\text{numK}(\mathbf{y})=i | \mathbf{y}_0]} \left[\sum_{\mathbf{y}_i=K} -\log p_{\theta, i}(\mathbf{y}_{0,i} | \mathbf{y}) \right]. \quad (79)$$

With proper settings on \mathbf{w} and the change of summation order, the above training loss of FHS will be equivalent to the λ -DCE loss shown in Ou et al. (2024), i.e.,

$$\mathcal{L}_{\lambda\text{-DCE}}(\mathbf{y}_0) = \mathbb{E}_{\lambda \sim \text{Uniform}(0,1)} \frac{1}{\lambda} \cdot \mathbb{E}_{\mathbf{y}_{\lambda} \sim q_{\lambda|0}^{\rightarrow}(\cdot | \mathbf{y}_0)} \left[\sum_{\mathbf{y}_{\lambda,i}=K} -\log p_{\theta}(\mathbf{y}_{0,i} | \mathbf{y}) \right].$$

Then, following from Appendix C.1 and Appendix C.2 in Ou et al. (2024), by choosing $\lambda(t) = 1 - e^{-t}$, with change of variable, the λ -DCE loss will be equivalent to the denoising score entropy loss, i.e.,

$$\begin{aligned} \mathcal{L}_{\text{DSE}}(\mathbf{y}_0) = & \int_0^T \mathbb{E}_{\mathbf{y}_t \sim q_{t|0}^{\rightarrow}(\cdot | \mathbf{y}_0)} \left[\sum_{\mathbf{y}' \neq \mathbf{y}_t} R^{\rightarrow}(\mathbf{y}_t, \mathbf{y}') \cdot \left(\frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t) \right. \right. \\ & \left. \left. - \frac{e^{-t}}{1 - e^{-t}} \cdot \delta_{\mathbf{y}_0, \text{DiffIdx}(\mathbf{y}_t, \mathbf{y}')}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}_t, \mathbf{y}')} \cdot \log \left(\frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t) \right) \right) \right] dt. \end{aligned}$$

Following from Theorem 3.4 in Lou et al. (2024), we note that DSE and SE share the same minimum, i.e.,

$$\begin{aligned} & \arg \min_{\theta} \mathbb{E}_{\mathbf{y}_0 \sim q_*} [\mathcal{L}_{\text{DSE}}(\mathbf{y}_0)] \\ &= \arg \min_{\theta} \int_0^T \mathbb{E}_{\mathbf{y}_t \sim q_t^{\rightarrow}} \left[\sum_{\mathbf{y}' \neq \mathbf{y}_t} R^{\rightarrow}(\mathbf{y}_t, \mathbf{y}') \cdot \left(\frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t) \right. \right. \\ & \quad \left. \left. \frac{q_t^{\rightarrow}(\mathbf{y}')}{q_t^{\rightarrow}(\mathbf{y}_t)} \cdot \log \left(\frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t) \right) \right) \right] := \arg \min_{\theta} \mathcal{L}_{\text{SE}}(\theta) \end{aligned} \quad (80)$$

By supposing

$$\tilde{v}_{t, \mathbf{y}_t}(\mathbf{y}') := \frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t) \quad \text{where} \quad \text{Ham}(\mathbf{y}', \mathbf{y}_t) = 1 \quad \text{and} \quad \mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} \neq \mathbf{K},$$

we know the Eq. 80 exactly matches Eq. 6.

Therefore, optimizing Eq. 79 in FHS is equivalent to parameterize the discrete score as

$$\begin{aligned} \frac{q_t^{\rightarrow}(\mathbf{y}')}{q_t^{\rightarrow}(\mathbf{y}_t)} &= \frac{e^{-t}}{1 - e^{-t}} \cdot q_{0, \text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t, \mathcal{K}(\mathbf{y}_t)) = v_{t, \mathbf{y}_t}(\mathbf{y}') \\ &\approx \tilde{v}_{t, \mathbf{y}_t}(\mathbf{y}') := \frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}_t)} | \mathbf{y}_t), \end{aligned}$$

and optimize Eq. 6. Following the analysis paradigm in this paper, we assume Assumption [A1] is also satisfies for this parametrization.

Bridge the trajectories between FHS and MATU. We have the following theorem.

Theorem 7 (The convergence of Alg. 1). *Suppose Assumption [A1] and [A2] hold, if the discrete scores are parameterized by time-independent neural network as Ou et al. (2024), the TV distance between the target discrete distribution q_* and the underlying distribution of the output particle \bar{q}_0 of Alg. 2 will satisfy $\text{TV}(q_*, \hat{q}_{T-\delta}) \leq 2\epsilon$.*

Proof. Under this time-independent parameterization, we suppose the trajectory of MATU as $\{\hat{\mathbf{y}}_t\}_{t=0}^T$ whose underlying distribution is denoted as $\hat{\mathbf{y}}_t \sim \bar{q}_t$. For FHS, we consider a sequence of random variables $\{\bar{\mathbf{y}}_k\}_{k \in \{0, 1, \dots, d\}}$ where $\bar{\mathbf{y}}_k$ denotes the random variables after $(d - k)$ -step update of FHS. We have $\text{numK}(\bar{\mathbf{y}}_k) = k$. To investigate the TV distance between $\hat{\mathbf{y}}_{T-\delta}$ and $\bar{\mathbf{y}}_0$, we have

$$\begin{aligned} \text{TV}(\hat{q}_{T-\delta}, \bar{q}_0) &= \frac{1}{2} \cdot \sum_{\mathbf{y}, \text{numK}(\mathbf{y})=0} |\bar{q}_0(\mathbf{y}) - \hat{q}_{T-\delta}(\mathbf{y})| + \frac{1}{2} \cdot \sum_{\mathbf{y}, \text{numK}(\mathbf{y}) \neq 0} \hat{q}_{T-\delta}(\mathbf{y}) \\ &= \frac{1}{2} \cdot \sum_{\mathbf{y}, \text{numK}(\mathbf{y})=0} |\bar{q}_0(\mathbf{y}) - \hat{q}_{T-\delta}(\mathbf{y})| + \bar{q}_0(\mathbf{y}) - \hat{q}_{T-\delta}(\mathbf{y}) \leq \sum_{\mathbf{y}, \text{numK}(\mathbf{y}) \neq 0} |\hat{q}_{T-\delta}(\mathbf{y}) - \bar{q}_0(\mathbf{y})| \end{aligned} \quad (81)$$

Currently, we define a distribution sequence

$$\{p_k\}_{k \in \{0, 1, \dots, d\}} \quad \text{where} \quad p_k(t) = \Pr[\text{the } k\text{-th transition happens at time } t].$$

Besides, suppose that at the transition time t the particle is \mathbf{y}' , MATU implies the transition from \mathbf{y}' to \mathbf{y} follows

$$\Pr[\mathbf{y} | \text{transition time} = t \text{ and particle is } \mathbf{y}'] = \hat{R}_t(\mathbf{y}, \mathbf{y}') / \hat{R}_t(\mathbf{y}')$$

Under this setting, we have

$$\begin{aligned}
\hat{q}_{T-\delta}(\mathbf{y}) &= \int_0^{T-\delta} p_{\text{numK}(\mathbf{y})}(t) \cdot \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} \hat{q}_t(\mathbf{y}') \cdot \Pr[\mathbf{y} | \text{transition time} = t \text{ and particle is } \mathbf{y}'] dt \\
&= \int_0^{T-\delta} p_{\text{numK}(\mathbf{y})}(t) \cdot \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} \hat{q}_t(\mathbf{y}') \cdot \frac{\hat{R}_t(\mathbf{y}, \mathbf{y}')}{\hat{R}_t(\mathbf{y}')} dt \\
&= \int_0^{T-\delta} p_{\text{numK}(\mathbf{y})}(t) \cdot \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} \hat{q}_t(\mathbf{y}') \cdot \frac{\tilde{R}_t(\mathbf{y}, \mathbf{y}')}{\tilde{R}_t(\mathbf{y}')} dt
\end{aligned} \tag{82}$$

Due to the time-independent parametrization of the discrete score, we have

$$\tilde{R}_t(\mathbf{y}, \mathbf{y}') = R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \tilde{v}_{t, \mathbf{y}'}(\mathbf{y}) = R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot \frac{e^{-t}}{1 - e^{-t}} \cdot p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}'),$$

which implies it has

$$\frac{\tilde{R}_t(\mathbf{y}, \mathbf{y}')}{\tilde{R}_t(\mathbf{y}')} = \frac{p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}'))}{\sum_{\mathbf{y} \neq \mathbf{y}', \text{Ham}(\mathbf{y}, \mathbf{y}')=1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}'))} = p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')).$$

Plugging this equation into Eq. 82, we have

$$\begin{aligned}
\hat{q}_{T-\delta}(\mathbf{y}) &= \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')) \int_0^{T-\delta} p_{\text{numK}(\mathbf{y})}(t) \cdot \hat{q}_t(\mathbf{y}') dt \\
&= \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')) \cdot \sum_{\mathbf{y}'', \text{numK}(\mathbf{y}'')=\text{numK}(\mathbf{y}')-1} p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}'')}) | \mathbf{y}'') \\
&\quad \cdot \dots \cdot \int_{t_{\text{numK}(\mathbf{y})}, t_{\text{numK}(\mathbf{y})-1}, \dots, t_1} p_{\text{numK}(\mathbf{y}), \text{numK}(\mathbf{y}), \dots, 1}(t_{\text{numK}(\mathbf{y})-1}, \dots, t_1) \cdot \hat{q}_{\tau}([K, \dots, K]) d \\
&\leq \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')) \cdot \sum_{\mathbf{y}'', \text{numK}(\mathbf{y}'')=\text{numK}(\mathbf{y}')-1} p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}'')}) | \mathbf{y}'') \\
&\quad \cdot \dots \cdot \sum_{\mathbf{y}^{(1)}, \text{numK}(\mathbf{y}^{(1)})=1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}^{(1)}, [K, \dots, K])}) | [K, \dots, K]) \hat{q}_0([K, \dots, K]),
\end{aligned}$$

where the last inequality follows from the fact

$$\hat{q}_{\tau}([K, \dots, K]) \leq \hat{q}_0([K, \dots, K]) \quad \forall \tau > 0.$$

According to the update of FHS, we can easily find that

$$\begin{aligned}
\bar{q}_0(\mathbf{y}) &= \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')) \cdot \sum_{\mathbf{y}'', \text{numK}(\mathbf{y}'')=\text{numK}(\mathbf{y}')-1} p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}'')}) | \mathbf{y}'') \\
&\quad \cdot \dots \cdot \sum_{\mathbf{y}^{(1)}, \text{numK}(\mathbf{y}^{(1)})=1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}^{(1)}, [K, \dots, K])}) | [K, \dots, K]) \underbrace{\bar{q}_d([K, \dots, K])}_{=1}.
\end{aligned}$$

Suppose the conditional distribution as

$$\begin{aligned}
\bar{p}_{\theta}(\mathbf{y} | [K, \dots, K]) &= \sum_{\mathbf{y}', \text{numK}(\mathbf{y}')=\text{numK}(\mathbf{y})-1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}, \mathbf{y}')}) | \mathbf{y}')) \cdot \sum_{\mathbf{y}'', \text{numK}(\mathbf{y}'')=\text{numK}(\mathbf{y}')-1} p_{\theta}(\mathbf{y}'_{\text{DiffIdx}(\mathbf{y}', \mathbf{y}'')}) | \mathbf{y}'') \\
&\quad \cdot \dots \cdot \sum_{\mathbf{y}^{(1)}, \text{numK}(\mathbf{y}^{(1)})=1} p_{\theta}(\mathbf{y}_{\text{DiffIdx}(\mathbf{y}^{(1)}, [K, \dots, K])}) | [K, \dots, K]),
\end{aligned}$$

then we have

$$\begin{aligned}
\bar{q}_0(\mathbf{y}) - \hat{q}_{T-\delta}(\mathbf{y}) &\leq \bar{p}_{\theta}(\mathbf{y} | [K, \dots, K]) \cdot (\bar{q}_d([K, \dots, K]) - \hat{q}_0([K, \dots, K])) \\
&= \bar{p}_{\theta}(\mathbf{y} | [K, \dots, K]) \cdot (1 - \hat{q}_0([K, \dots, K]))
\end{aligned} \tag{83}$$

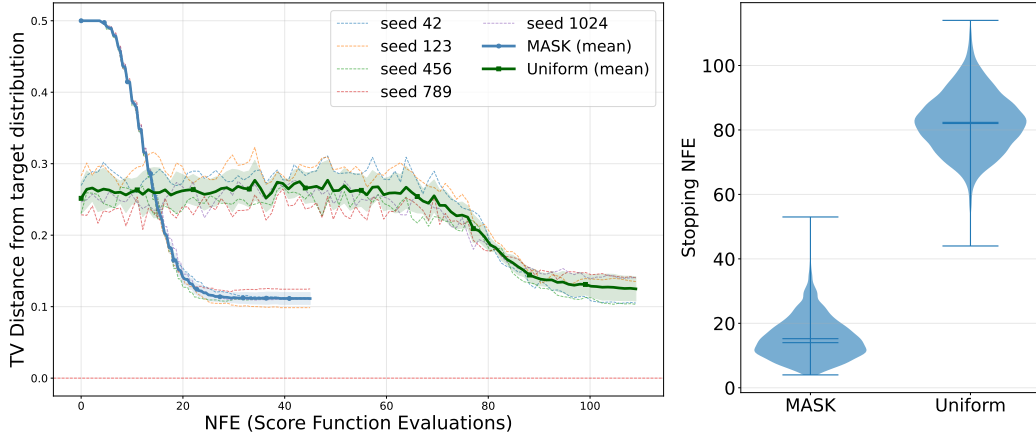


Figure 1: **Synthetic experiment results on sampling efficiency.** We compare our proposed Masked Discrete Diffusion (MASK) against the Uniform baseline with vocabulary size $K = 3$ and sequence length $d = 4$. **Left:** The Total Variation (TV) distance between the empirical and ground truth distributions as a function of the Number of (Score) Function Evaluations (NFE). The solid lines represent the mean over 5 seeds, and shaded regions indicate the standard deviations. Our method achieves faster convergence to the target distribution. **Right:** Violin plots illustrating the distribution of Stopping NFE. The MASK method requires significantly fewer evaluations to terminate compared to the Uniform baseline.

According to the proof of Lemma 2, we know that

$$\begin{aligned} \hat{q}_0([K, \dots, K]) &= (1 + e^{-T})^{-d} \quad \text{and} \quad (1 + e^{-T})^d - 1 \leq \epsilon/2 \\ \Rightarrow 0 &\leq 1 - (1 + e^{-T})^{-d} = 1 - \hat{q}_0([K, \dots, K]) \leq 1 - 1/(1 + \epsilon/2) \leq \epsilon/2. \end{aligned} \quad (84)$$

Combining Eq. 81, Eq. 83 and Eq. 84, we have

$$\text{TV}(\hat{q}_{T-\delta}, \bar{q}_0) \leq \epsilon/2 \quad \text{and} \quad \text{TV}(q_*, \bar{q}_0) \leq \text{TV}(\hat{q}_{T-\delta}, \bar{q}_0) + \text{TV}(q_*, \hat{q}_{T-\delta}) \leq \epsilon$$

where last inequality follows from Theorem 3. Hence, the proof is completed. \square

G EXPERIMENTS

G.1 SYNTHETIC EXPERIMENTS.

We conduct synthetic experiments to validate our theoretical findings and compare the sampling efficiency of our Masked Discrete Diffusion model against the uniform baseline.

Experiment Setup. We utilize a state space defined by vocabulary size $K = 3$ and sequence length $d = 4$. The ground truth distribution, p^* , is constructed by assigning a random mass sampled uniformly from $(0, 1)$ to each of the K^d possible sequences and normalizing the distribution. We report results averaged over 5 independent random seeds. For each seed, we generate 1000 trajectories using our method (Algorithm 1, MATU) and the truncated uniformization baseline with a uniform stationary distribution (adapted from Huang et al. (2025)). Performance is evaluated via the Total Variation (TV) distance between the empirical marginal distribution and p^* , plotted as a function of the Number of (Score) Function Evaluations (NFE). Quantitative results are shown in Figure 1, and illustrative sampling trajectories are visualized in Figure 2.

G.2 REAL WORLD EXPERIMENTS

We consider to introduce our Alg. 1 (MATU) into the text generation task.

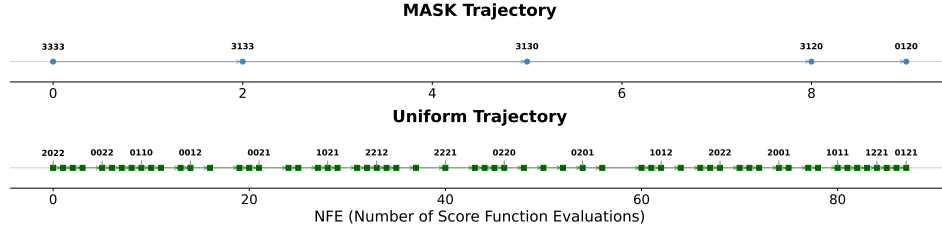


Figure 2: Visualization of individual sampling trajectories. The plots show single sampling paths, with labels indicating the intermediate discrete states. The MASK method (top) navigates the state space efficiently with few steps. In contrast, the Uniform baseline (bottom) exhibits diffusive behavior with many small steps—often reverting previous changes—resulting in a high NFE cost.

Experimental Settings In this paragraph, we follow the problem setting as SEDD shown in Lou et al. (2024), and consider the unconditional text generation task with the small pretrained SEDD Absorbing model. The sequence length of generated sample is constrained as $d = 1024$, and the vocabulary size will be $K = 50258$, including the mask token. We choose the typical Euler and Tweedie’s τ -leaping (analytic samples in Lou et al. (2024)’s implementation) as our baselines. For the step number choice, we only consider $\{1024, 2048\}$. Because MATU does not consider the conditional independent assumption for the reverse process. Under this condition, it requires at least d steps to generate one no-mask sample.

The inexact adaptation from MATU. In SEDD experiments, The exact implementation of Alg. 1 will require the inference complexity to be $K \times d = 50258 \times 1024$, which is far beyond an acceptable inference complexity. Since the choice of K can be used to control the inference complexity, in the following experiment we will choose

$$K = \text{required steps/generated sequence length},$$

which is an inexact implementation of Alg. 1 (MATU), while makes it to be possible to be tuned via the choice of the step number. Moreover, the implementation of Euler and Tweedie’s τ -leaping is based on log-linear noise schedule, which means the transition rate matrix of the forward process satisfies

$$R_t^{\rightarrow}(\mathbf{y}, \mathbf{y}') = \sigma(t) R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \quad \text{where} \quad \sigma(t) = \frac{1 - \epsilon}{1 - (1 - \epsilon) \cdot t}$$

and R^{\rightarrow} follows from Eq. 7. Under this condition, the reverse transition rate matrix will become

$$\begin{aligned} R_t^{\leftarrow}(\mathbf{y}, \mathbf{y}') &:= \sigma(1 - t) \cdot R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \frac{q_t^{\leftarrow}(\mathbf{y})}{q_t^{\leftarrow}(\mathbf{y}')} \\ &= \sigma(1 - t) \cdot R^{\rightarrow}(\mathbf{y}', \mathbf{y}) \cdot s_{\theta, 1-t, \mathbf{y}'}(\mathbf{y}). \end{aligned}$$

Empirical Results. We use PPL and entropy as two criteria to measure the generation quality for different samplers. The results are summarized as the following tables. We will release the detailed code and implementation after the acceptance of this paper.

Table 3: Comparison of the inference generation performance, we calculate the average perplexity and entropy for 32 samples generated by Euler, Analytic and MATU. The experiments show even with an inexact implementation, MATU still outperform then other samplers consistently.

Samplers	Steps	Avg Perplexity	Std Perplexity	Avg Entropy	Std Entropy	Wall-clock time
Euler	1024	41.42	11.68	7.588	0.301	27.35s/sample
Analytic	1024	41.81	11.57	7.597	0.286	24.15s/sample
MATU	1024	40.54	11.20	7.554	0.230	32.23s/sample
Euler	2048	33.32	7.141	7.492	0.258	53.43s/sample
Analytic	2048	32.50	6.952	7.489	0.250	46.88s/sample
MATU	2048	31.82	6.717	7.394	0.332	60.05s/sample