On the Complexity Theory of Masked Discrete Diffusion: From $poly(1/\epsilon)$ to Nearly ϵ -Free

Anonymous authors

000

001

002003004

006

007 008 009

010 011

012

013

014

016

018

019

021

023

024

025

026

027

028

029

031

033

035

037

039

040

041

042

043

044 045

046

047

048

050

051

052

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) (Nie et al., 2025), analyzing and mitigating its computational overhead in high-dimensional settings remains challenging. As summarized in Table 1, 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 at least $\tilde{O}(\epsilon^{-1})$ 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. (2025) 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., 2025) 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., 2025), 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 tighting 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., 2025) 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., 2025), 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, ..., K}^d. Specifically, to reach ε-TV convergence, MATU uses an expected number of discrete score calls on the order of

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

Compared to uniformization-based sampler in uniform discrete diffusion (Huang et al., 2025; Liang et al., 2025), 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., 2025). 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 $\mathcal{Y} = \{1, 2, ..., K\}^d$. For any functions $f, g : \mathcal{Y} \to \mathbb{R}$, we define their inner product as

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

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 $R^{\to} \colon \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ which denotes the instantaneous transition rate from y' to y. Formally,

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

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

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

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

$$\frac{\mathrm{d}\,q_t^{\leftarrow}}{\mathrm{d}\,t}(\boldsymbol{y}) = \langle R_t^{\leftarrow}(\boldsymbol{y},\cdot),\,q_t^{\leftarrow}(\cdot)\rangle_{\mathcal{Y}} \quad \text{where} \quad R_t^{\leftarrow}(\boldsymbol{y},\boldsymbol{y}') := R^{\rightarrow}(\boldsymbol{y}',\boldsymbol{y})\,\frac{q_t^{\leftarrow}(\boldsymbol{y})}{q_t^{\leftarrow}(\boldsymbol{y}')}, \tag{3}$$

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

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

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

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 $q_t^{\leftarrow}(y)/q_t^{\leftarrow}(y')$. Since this ratio is generally intractable, it is approximated in practice by a neural network \tilde{v} :

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

$$L_{\text{SE}}(\tilde{v}) = \frac{1}{T} \int_{0}^{T} \mathbb{E}_{\mathbf{y}_{t} \sim q_{t}^{\rightarrow}} \left[\sum_{\mathbf{y} \neq \mathbf{y}_{t}} R^{\rightarrow}(\mathbf{y}_{t}, \mathbf{y}) D_{\phi} \left(v_{T-t, \mathbf{y}_{t}}(\mathbf{y}) \middle\| \tilde{v}_{T-t, \mathbf{y}_{t}}(\mathbf{y}) \right) \right] dt,$$
 (6)

where $D_{\phi}\left(\cdot \| \cdot \right)$ 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 $L_{\text{SE}}(\tilde{v}) \leq \epsilon_{\text{score}}^2$.

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. (2025).

Additional settings. Following Ou et al. (2024), we adopt a diffusion-based language modeling framework. Our vocabulary is $\{1, 2, ..., K\}$, where K denotes the mask token. We aim to generate a length-d sequence (sentence) $\mathbf{y} \in \mathcal{Y} = \{1, 2, ..., K\}^d$. 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

$$\operatorname{numK}\left(\boldsymbol{y}\right) \;\coloneqq\; \sum_{i=1}^{d} \delta_{\mathrm{K}}(\boldsymbol{y}_{i}) \quad \text{and} \quad \operatorname{Ham}(\boldsymbol{y}, \boldsymbol{y}') = \sum_{i=1}^{d} \delta_{\boldsymbol{y}_{i}}(\boldsymbol{y}'_{i})$$

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

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

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

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

Here, DiffIdx (y, y') denotes the single coordinate where y and 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

$$q_{t|s}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}') = \prod_{i=1}^{d} \left[\delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}'_{i}) + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}'_{i})\right) \cdot \delta_{0}(\boldsymbol{y}_{i} - \boldsymbol{y}'_{i}) \cdot e^{-(t-s)} + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}'_{i})\right) \cdot \delta_{K}(\boldsymbol{y}_{i}) \cdot (1 - e^{-(t-s)}) \right],$$
(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

$$\tilde{q}_t(\boldsymbol{y}) \propto \prod_{i=1}^d \exp(-t \cdot [1 - \delta_{K}(\boldsymbol{y}_i)]) = \exp(-t \cdot [d - \text{numK}(\boldsymbol{y})]).$$
 (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,

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

Consequently, to ensure $\mathrm{KL}\left(q_t^{\rightarrow} \| \tilde{q}_t\right) \leq \epsilon$, 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 $\{\hat{y}_t\}_{t=0}^T$ denote the practical reverse process, whose marginal, joint, and conditional distributions satisfy:

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

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

$$\hat{q}_{t+\Delta t|t}(\boldsymbol{y}'|\boldsymbol{y}) \propto \prod_{i=1}^{d} \hat{q}_{t+\Delta t|t}^{(i)}(\boldsymbol{y}[\{i\} \to \{\boldsymbol{y}_{i}'\}]|\boldsymbol{y}), \tag{10}$$

where the token revision function

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

indicates that the coordinates of 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 ith coordinate takes the form:

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

Since $\operatorname{Ham}(\boldsymbol{y},\ \boldsymbol{y}[\{i\}\to\boldsymbol{y}_i')=1$, the definition of R^\to in Eq. (7) ensures that $R^\to(\boldsymbol{y},\ \boldsymbol{y}[\{i\}\to\boldsymbol{y}_i')\neq 0$. Hence, $\hat{q}_{t+h|t}^{(i)}(\boldsymbol{y}[\{i\}\to k]|\boldsymbol{y})$ 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:

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

where $\hat{q}_{t+\Delta t|t}(y'\mid y)\neq 0$ only if $R^{\rightarrow}(y,y')\neq 0$, which implies $\operatorname{Ham}(y,y')=1$ (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

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

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. (2025) hold, implement Euler sampler with Eq. (11), if we require

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

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

Compared to the au-leaping method analyzed in Liang et al. (2025), 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 y,y' can be described by

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

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

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

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 $\Delta t \to 0$:

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

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

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

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

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

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

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 y, the reverse transition rate function satisfies

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

Compared to (16), this bound explicitly depends on $\operatorname{numK}(y)$, the number of mask tokens in y. Since $\operatorname{numK}(y) \leq d$, it is strictly smaller than the uniform bound in (16). Furthermore, $\operatorname{numK}(y)$ 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 $R_t^{\leftarrow}(y',y)$ by a learned neural score $\tilde{v}_{t,y}(y')$, yielding

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

as dictated by Lemma 1 and Eq. (5). Because \tilde{v} is a learned estimator, the outgoing rate $\tilde{R}_t(y)$ 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 $\tilde{R}_t(y)$, restricting it to remain below a fixed constant (Liang et al., 2025) 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 $0 = t_0 < \ldots < t_W = T \delta$, parameters $\beta_{t_1}, \ldots, \beta_{t_W}$ set as Eq. (17), a reverse transition rate function \hat{R}_t^{\leftarrow} obtained by the learnt score function $\tilde{v}_{t,y'}(\cdot)$.
- 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 $\tau_1 < \tau_2 < \ldots < \tau_N$;
 - 7: Set $\mathbf{z}_0 = \hat{\mathbf{y}}_{t_{w-1}}$;
- 387 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 z_{n-1} with

$$\mathbf{z}_n = \begin{cases} \mathbf{z}_{n-1}[\mathbf{z}_i \colon K \to k] & w.p. \ \beta_{t_w}^{-1} \cdot \hat{R}(\mathbf{z}_{n-1}[\mathbf{z}_i \colon K \to k], \mathbf{z}_{n-1}), \\ \mathbf{z}_{n-1}, & w.p. \ 1 - \beta_{t_w}^{-1} \cdot \hat{R}_{\tau_n}^{\leftarrow}(\mathbf{z}_{n-1}). \end{cases}$$

- 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 $1 - \beta^{-1} \tilde{R}_t(y) < 0$, 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 $\tilde{R}_t(y',y)$. This ensures that the non time-discretization property is preserved without additional cost, even when $\tilde{R}_t(y)$ 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 $\hat{\mathbf{y}}_{t_{w-1}}=y_{t_{w-1}}$. Following from the monotonicity of $(e^{T-t}-1)^{-1}$ and $\operatorname{numK}(\hat{\mathbf{y}}_t)$ in Lemma 3, the mask-aware truncation is chosen as $\beta_{t_w}(y_{t_{w-1}})$, then we set

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

and

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

With these truncations, the corrected outgoing rate will be definitely upper bounded by $\beta_{t_w}(y_{t_{w-1}})$. 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

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

and the partition of the reverse process satisfies

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

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

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

to achieve TV $(p_*, \hat{p}) \le 2\epsilon$ 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. (2025) which is a little bit stronger than [A3] due to the time-invariant requirement. All complexities for TV convergence 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]	$\tilde{\mathcal{O}}(d^{5/3}\epsilon^{-2})$
Ren et al. (2024)	Uniformed	au-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], <mark>[A3]</mark> +	$\tilde{O}(d^2\epsilon^{-3/2})$
Liang et al. (2025)	Masked	au-leaping	[A1],[A2],[A3]+	$O(d\epsilon^{-2})$
Liang et al. (2025)	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 $\operatorname{poly}(d/\epsilon)$ 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 1. 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:

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

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 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.
- 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.
- Emiel Hoogeboom, 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, 2025.
- 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.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025.
- 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 highorder 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 learn-ing*, pp. 2256–2265. pmlr, 2015.

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 appli-cation 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. 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.

CONTENTS Introduction **Preliminaries** The Forward Process of Masked Discrete Diffusion 4 Euler Sampler in Masked Discrete Diffusion **Truncated Uniformization in Masked Discrete Diffusion** Conclusion **A Notation Summary B** The Markov Processes of Discrete Diffusion Models C Euler Discretization Analysis **D** Truncated Uniformization Inference Analysis **E** Technical Lemmas The Use of Large Language Models (LLMs)

A NOTATION SUMMARY

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

652		
653		

	Table 2: Summary of key notations used in the paper.		
Symbol	Description		
$egin{array}{c} q_{*} & & & & & & & & & & & & & & & & & & &$	Discrete distribution on $\mathcal{Y} = \{1, 2, \dots, K\}^d$		
$\mathbf{y}_t^{\rightarrow}$	Forward-time CTMC on ${\cal Y}$		
$q_t^{ ightarrow}$	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 $(\mathbf{y}_{t'}^{\rightarrow}, \mathbf{y}_{t}^{\rightarrow})$		
$ ilde{q}_t$	Aapproximation of q_t^{\rightarrow} constructing the reverse initialization, Eq. (9)		
$q_{t' t}^{ ightarrow}(oldsymbol{y}' oldsymbol{y})$	Conditional transition probability in forward process, Eq. (36)		
$\mathbf{y}_t^{\leftarrow}$	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 $(\mathbf{y}_{t'}^{\leftarrow}, \mathbf{y}_{t}^{\leftarrow})$		
$q_{t' t}^\leftarrow(oldsymbol{y}' oldsymbol{y})$	Conditional transition probability of the ideal reverse process		
\hat{q}_t	Marginal distribution of reverse process at time t implemented by Alg. 1		
qt',t	Joint distribution of $(\hat{\mathbf{y}}_{t'}, \hat{\mathbf{y}}_t)$		
$\hat{q}_{t' t}(oldsymbol{y}' oldsymbol{y})$	Conditional transition probability of the ideal reverse process		
$\overline{R^{ ightarrow}(oldsymbol{y},oldsymbol{y}')}$	Forward transition rate, i.e., Eq. (7), from state y' to y . This follows the ordering of the conditional distribution $p(y y')$, which is the <i>transpose</i> of the convention used in some other works.		
$R_t^{\leftarrow}(oldsymbol{y},oldsymbol{y}')$	Reverse transition rate at time t from state y' to y , $R_t^{\leftarrow}(y,y') \coloneqq R^{\rightarrow}(y',y) \cdot \frac{q_t^{\leftarrow}(y)}{q_t^{\leftarrow}(y')}$, Eq. (3)		
$ ilde{R}_t(oldsymbol{y},oldsymbol{y}')$	Estimated reverse transition rate using the learned density ratio, $\tilde{R}_t(\boldsymbol{y}, \boldsymbol{y}') = R^{\rightarrow}(\boldsymbol{y}', \boldsymbol{y}) \cdot \tilde{v}_{t, \boldsymbol{y}'}(\boldsymbol{y})$, Eq. (6)		
$\hat{R}_t(\cdot,\cdot)$	Truncated version of $R_t(\cdot, \cdot)$ with threshold β_t , Eq. (18)		
$R_t^{\leftarrow}(\boldsymbol{y}), \; \tilde{R}_t(\boldsymbol{y}), \; \hat{R}_t(\boldsymbol{y})$	Total reverse transition rate out of state ${\boldsymbol y}$ for each rate type, defined as $R({\boldsymbol y}) \coloneqq$		
	$\sum_{m{y}' eq m{y}} R(m{y}', m{y}) ext{ with } R \in \{R_t^{\leftarrow}, \ ilde{R}_t, \ \hat{R}_t\}$		
β_t	Upper bound on $R_t^{\leftarrow}(\boldsymbol{y})$, $\beta_t = \text{numK}(\boldsymbol{y}) \cdot K/(T-t)$, Eq. (17)		

B THE MARKOV PROCESSES OF DISCRETE DIFFUSION MODELS

Score entropy loss used to train \tilde{v} , Eq. (6)

Density ratio $q_t^{\leftarrow}(\boldsymbol{y})/q_t^{\leftarrow}(\boldsymbol{y}')$

B.1 THE FORMULATIONS OF THE FORWARD PROCESS

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

Learned approximation to $v_{t, \mathbf{y}'}(\mathbf{y}) = q_t^{\leftarrow}(\mathbf{y})/q_t^{\leftarrow}(\mathbf{y}')$

The number of [MASK] token (or token K) in a vector.

One-hot vector with a 1 at position i and 0 elsewhere

Indicator function with $\delta_{y}(y) = 1$ and $\delta_{y}(y') = 0$ $(y' \neq y)$

$$Q_t^{\rightarrow}[f](\boldsymbol{y}) = \mathbb{E}\left[f(\mathbf{y}_t)|\mathbf{y}_0 = \boldsymbol{y}\right] = \left\langle f, q_{t|0}^{\rightarrow}(\cdot|\boldsymbol{y}) \right\rangle_{\mathcal{Y}}$$
(21)

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

$$\mathcal{L}^{\rightarrow}[f](\boldsymbol{y}) = \lim_{t \to 0} \left[\frac{\mathcal{Q}_{t}^{\rightarrow}[f] - f}{t} \right](\boldsymbol{y}) = \left\langle f, \partial_{t} q_{t|0}^{\rightarrow}(\cdot|\boldsymbol{y}) \Big|_{t=0} \right\rangle_{\mathcal{Y}} := \left\langle f, R^{\rightarrow}(\cdot, \boldsymbol{y}) \right\rangle_{\mathcal{Y}}$$
(22)

where

 $v_{t, \boldsymbol{y}'}(\boldsymbol{y})$

 $\tilde{v}_{t, \boldsymbol{y}'}(\boldsymbol{y})$

 $L_{\rm SE}(\hat{v})$

 \boldsymbol{e}_i

 $\delta_{\boldsymbol{y}}(\cdot)$

 $numK(\cdot)$

$$R^{\rightarrow}(\boldsymbol{y}',\boldsymbol{y}) := \partial_t q_{t|0}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y})\Big|_{t=0} = \lim_{t \to 0} \left[\frac{q_{t|0}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y}) - \delta_{\boldsymbol{y}}(\boldsymbol{y}')}{t} \right].$$
 (23)

According to the time-homogeneous property, we have

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

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

$$R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') \ge 0 \text{ when } \boldsymbol{y}' \ne \boldsymbol{y} \quad \text{and} \quad R^{\rightarrow}(\boldsymbol{y}', \boldsymbol{y}') = -\sum_{\boldsymbol{y} \ne \boldsymbol{y}'} R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') \le 0$$
 (24)

due to the definition Eq. (23). Under this setting, we can provide the dynamic of $q_{t|0}$ for any t. Specifically, we have

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

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

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

derived from Eq. (21), we have

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

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

$$\partial_t q_{t|s}^{\rightarrow}(\tilde{\boldsymbol{y}}|\boldsymbol{y}) = \sum_{\boldsymbol{y}' \in \mathcal{Y}} R(\tilde{\boldsymbol{y}}, \boldsymbol{y}') \cdot q_{t|s}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y}) = \left\langle R(\tilde{\boldsymbol{y}}, \cdot), q_{t|s}^{\rightarrow}(\cdot|\boldsymbol{y}) \right\rangle_{\mathcal{Y}}.$$
 (25)

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

$$\frac{\mathrm{d}q_t^{\rightarrow}}{\mathrm{d}t}(\boldsymbol{y}) = \langle R(\boldsymbol{y}, \cdot), q_t^{\rightarrow} \rangle_{\mathcal{Y}}.$$
 (26)

Matrix Formulation. Suppose the support set \mathcal{Y} of q_t^{\rightarrow} be written as $\mathcal{Y} = \{y_1, y_2, \dots, y_{|\mathcal{Y}|}\}$, we may consider the marginal distribution q_s^{\rightarrow} to be a vector, i.e.,

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

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

$$\boldsymbol{Q}_{t|s}^{\rightarrow} = \begin{bmatrix} q_{t|s}^{\rightarrow}(\boldsymbol{y}_{1}|\boldsymbol{y}_{1}) & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{1}|\boldsymbol{y}_{2}) & \dots & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{1}|\boldsymbol{y}_{|\mathcal{Y}|}) \\ q_{t|s}^{\rightarrow}(\boldsymbol{y}_{2}|\boldsymbol{y}_{1}) & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{2}|\boldsymbol{y}_{2}) & \dots & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{2}|\boldsymbol{y}_{|\mathcal{Y}|}) \\ \dots & \dots & \dots & \dots \\ q_{t|s}^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}|\boldsymbol{y}_{1}) & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}|\boldsymbol{y}_{2}) & \dots & q_{t|s}^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}|\boldsymbol{y}_{|\mathcal{Y}|}) \end{bmatrix}.$$

Similarly, the function R can also be presented as

$$\boldsymbol{R}^{\rightarrow} = \begin{bmatrix} R^{\rightarrow}(\boldsymbol{y}_{1}, \boldsymbol{y}_{1}) & R^{\rightarrow}(\boldsymbol{y}_{1}, \boldsymbol{y}_{2}) & \dots & R^{\rightarrow}(\boldsymbol{y}_{1}, \boldsymbol{y}_{|\mathcal{Y}|}) \\ R^{\rightarrow}(\boldsymbol{y}_{2}, \boldsymbol{y}_{1}) & R^{\rightarrow}(\boldsymbol{y}_{2}, \boldsymbol{y}_{2}) & \dots & R^{\rightarrow}(\boldsymbol{y}_{2}, \boldsymbol{y}_{|\mathcal{Y}|}) \\ \dots & \dots & \dots & \dots \\ R^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}, \boldsymbol{y}_{1}) & R^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}, \boldsymbol{y}_{2}) & \dots & R^{\rightarrow}(\boldsymbol{y}_{|\mathcal{Y}|}, \boldsymbol{y}_{|\mathcal{Y}|}) \end{bmatrix}.$$
(27)

Under this condition, Eq. (26) can be written as

$$\mathrm{d}\boldsymbol{q}_t^{\rightarrow}/\mathrm{d}t = \boldsymbol{R}^{\rightarrow} \cdot \boldsymbol{q}_t^{\rightarrow} \tag{28}$$

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.,

$$\partial_{t}q_{t|s}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y}) = -1 \cdot \partial_{T-t}q_{T-t|T-s}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y}) = -1 \cdot \partial_{T-t} \left[\frac{q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}') \cdot q_{T-t}^{\rightarrow}(\boldsymbol{y}')}{q_{T-s}^{\rightarrow}(\boldsymbol{y})} \right]$$

$$= -\underbrace{\partial_{T-t}q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}') \cdot \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y}')}{q_{T-s}^{\rightarrow}(\boldsymbol{y})}}_{\text{Term 1}} - \underbrace{\frac{q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}')}{q_{T-s}^{\rightarrow}(\boldsymbol{y})}}_{\text{Term 2}} \cdot \partial_{T-t}q_{T-t}^{\rightarrow}(\boldsymbol{y}')}_{\text{Term 2}}.$$
(29)

For Term 1 of Eq. (29), we have

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

= $-\sum_{\tilde{\boldsymbol{y}} \in \mathcal{Y}} R^{\rightarrow}(\tilde{\boldsymbol{y}}, \boldsymbol{y}') \cdot \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y}')}{q_{T-t}^{\rightarrow}(\tilde{\boldsymbol{y}})} \cdot q_{T-t|T-s}^{\rightarrow}(\tilde{\boldsymbol{y}}|\boldsymbol{y}),$

where the first equation follows from the Kolmogorov backward theorem (Lemma 14) and Eq. (22):

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

For Term 2 of Eq. (29), we have

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

where the first equation follows from Eq. (26) and the last equation follows from the fact

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

Under this condition, by setting

$$R_t^{\leftarrow}(\boldsymbol{y}', \tilde{\boldsymbol{y}}) \coloneqq R(\tilde{\boldsymbol{y}}, \boldsymbol{y}') \cdot \frac{q_t^{\leftarrow}(\boldsymbol{y}')}{q_t^{\leftarrow}(\tilde{\boldsymbol{y}})}$$

then Eq. (29) can be summarized as

$$\partial_t q_{t|s}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y}) = \left\langle R_t^{\leftarrow}(\boldsymbol{y}',\cdot), q_{t|s}^{\leftarrow}(\cdot|\boldsymbol{y}) \right\rangle_{\mathcal{Y}} = \sum_{\tilde{\boldsymbol{y}} \in \mathcal{Y}} R_t^{\leftarrow}(\boldsymbol{y}', \tilde{\boldsymbol{y}}) \cdot q_{t|s}^{\leftarrow}(\tilde{\boldsymbol{y}}|\boldsymbol{y}). \tag{30}$$

Combining with Bayes' Theorem, we have

$$\frac{\mathrm{d}q_t^{\leftarrow}}{\mathrm{d}t}(\boldsymbol{y}) = \langle R_t^{\leftarrow}(\boldsymbol{y},\cdot), q_t^{\leftarrow} \rangle_{\mathcal{Y}}.$$
 (31)

Hence, Eq. (3) establishes.

Moreover, since the RHS of Eq. (4) satisfies

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

Besides, we have

$$\lim_{s \to t} \partial_t q_{t|s}^{\leftarrow}(\boldsymbol{y}|\boldsymbol{y}') = \lim_{s \to t} \partial_t \left[q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y})}{q_{T-s}^{\rightarrow}(\boldsymbol{y}')} \right]$$

$$= \lim_{s \to t} \left[\partial_t (q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y})) \cdot \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y})}{q_{T-s}^{\rightarrow}(\boldsymbol{y}')} + q_{T-s|T-t}^{\rightarrow}(\boldsymbol{y}'|\boldsymbol{y}) \cdot \frac{\partial_t q_{T-t}^{\rightarrow}(\boldsymbol{y})}{q_{T-s}^{\rightarrow}(\boldsymbol{y}')} \right].$$

When $y \neq y'$, we have

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

which implies

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

The last equation follows from the Kolmogorov backward theorem, i.e., Lemma 14 and Eq. (22)

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

Combining with Eq. (3), we have

$$\lim_{\Delta t \to 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\boldsymbol{y}|\boldsymbol{y}') - \delta_{\boldsymbol{y}'}(\boldsymbol{y})}{\Delta t} \right] = \lim_{s \to t} \partial_t q_{t|s}^{\leftarrow}(\boldsymbol{y}|\boldsymbol{y}') = R^{\rightarrow}(\boldsymbol{y}', \boldsymbol{y}) \cdot \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y})}{q_{T-t}^{\rightarrow}(\boldsymbol{y}')} = R_t^{\leftarrow}(\boldsymbol{y}, \boldsymbol{y}')$$
(32)

when $y' \neq y$. Besides, we have

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

which means

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

where the last inequality follows from Eq. (32). Hence, Eq. (3) establishes, and the proof is completed. \Box

B.3 THE PROOF OF LEMMA 2

Lemma 4. The close solution of Eq. (28) is

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

Proof. We can easily verify that

$$\frac{\mathrm{d}\boldsymbol{q}_t^{\rightarrow}}{\mathrm{d}t} = \frac{\mathrm{d}}{\mathrm{d}t} \left[\exp(t\boldsymbol{R}^{\rightarrow}) \boldsymbol{q}_0^{\rightarrow} \right] = \frac{\mathrm{d}}{\mathrm{d}t} \left[\exp(t\boldsymbol{R}^{\rightarrow}) \right] \boldsymbol{q}_0^{\rightarrow}.$$

With the following equation,

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

we have

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

Hence, the proof is completed.

Lemma 5. Suppose the transition rate matrix \mathbb{R}^{\to} shown as Eq. (27) satisfies Eq. (7). It can be decomposed as

$$m{R}^{
ightarrow} = \sum_{i=1}^d m{R}_i^{
ightarrow} \quad ext{where} \quad m{R}_i^{
ightarrow} = oldsymbol{I} \otimes \cdots \otimes m{A} \otimes \cdots \otimes m{I},$$

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

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{bmatrix}. \tag{33}$$

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

$$R_i^{\rightarrow}(y,y') = I(y_1,y_1') \cdot \ldots \cdot A(y_i,y_i') \cdot \ldots \cdot I(y_d,y_d').$$

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

$$\boldsymbol{R}_{j}^{\rightarrow}(\boldsymbol{y},\boldsymbol{y}') = \boldsymbol{A}(\boldsymbol{y}_{j},\boldsymbol{y}'_{j}) \cdot \boldsymbol{I}(\boldsymbol{y}_{1},\boldsymbol{y}'_{1}) \cdot \ldots \cdot \underbrace{\boldsymbol{I}(\boldsymbol{y}_{j_{1}},\boldsymbol{y}'_{j_{1}})}_{-0} \cdot \ldots \cdot \underbrace{\boldsymbol{I}(\boldsymbol{y}_{j_{2}},\boldsymbol{y}'_{j_{2}})}_{-0} \cdot \ldots \cdot \boldsymbol{I}(\boldsymbol{y}_{d},\boldsymbol{y}'_{d}) = 0.$$

Besides, for $j = j_1$, we have

$$R_{j_1}^{\rightarrow}(y,y') = A(y_{j_1},y'_{j_1}) \cdot I(y_1,y'_1) \cdot \dots \cdot \underbrace{I(y_{j_2},y'_{j_2})}_{=0} \cdot \dots \cdot I(y_d,y'_d) = 0.$$

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

$$oldsymbol{R}^{
ightarrow}(oldsymbol{y},oldsymbol{y}') = \sum_{i=1}^d oldsymbol{R}_i^{
ightarrow}(oldsymbol{y},oldsymbol{y}') = 0 \quad ext{when} \quad ext{Ham}(oldsymbol{y},oldsymbol{y}') \geq 2$$

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

$$R_j^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') = A(\boldsymbol{y}_j, \boldsymbol{y}_j') \cdot I(\boldsymbol{y}_0, \boldsymbol{y}_0') \cdot \dots \cdot \underbrace{I(\boldsymbol{y}_{j_1}, \boldsymbol{y}_{j_1}')}_{=0} \cdot \dots \cdot I(\boldsymbol{y}_d, \boldsymbol{y}_d') = 0.$$

Otherwise, when $j = j_1$, we have

$$m{R}_{j_1}^{
ightarrow}(m{y},m{y}') = m{A}(m{y}_{j_1},m{y}_{j_1}')\cdotm{I}(m{y}_1,m{y}_1')\cdot\ldots\cdotm{I}(m{y}_d,m{y}_d') = m{A}(m{y}_{j_1},m{y}_{j_1}')$$

where the second equation establishes since $\operatorname{Ham}(\boldsymbol{y},\boldsymbol{y}')=1$ and $\boldsymbol{y}_j=\boldsymbol{y}'_j$ when $j\neq j_1$. Then, only when $\boldsymbol{y}_{j_1}=K$, we will have $\boldsymbol{A}(\boldsymbol{y}_{j_1},\boldsymbol{y}'_{j_1})=1$ otherwise $\boldsymbol{A}(\boldsymbol{y}_{j_1},\boldsymbol{y}'_{j_1})=0$ due to the definition Eq. (33). That means

$$m{R}^{
ightarrow}(m{y},m{y}') = \sum_{i=1}^d m{R}_i^{
ightarrow}(m{y},m{y}') = 0 \quad ext{when} \quad ext{Ham}(m{y},m{y}') = 1 \text{ and } m{y}_{ ext{DiffIdx}(m{y},m{y}')}
eq K$$

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

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

$$m{R}_i^{
ightarrow}(m{y},m{y}') = m{A}(m{y}_j,m{y}_j') \cdot m{I}(m{y}_1,m{y}_1') \cdot \ldots \cdot m{I}(m{y}_d,m{y}_d') = m{A}(m{y}_j,m{y}_j'),$$

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 $R^{\rightarrow}(y,y')=\sum_{i=0}^{d-1}R_i^{\rightarrow}(y,y')$ when y=y'. Hence, the proof is completed.

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

$$m{R}^{
ightarrow} = \sum_{i=1}^d m{R}_i^{
ightarrow} \quad where \quad m{R}_i^{
ightarrow} = \underbrace{m{I} \otimes \ldots \otimes m{I}}_{i-1 \; terms} \otimes m{A} \otimes \underbrace{m{I} \otimes \ldots \otimes m{I}}_{d-i \; terms},$$

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

$$oldsymbol{R}_i^{
ightarrow} \cdot oldsymbol{R}_i^{
ightarrow} = oldsymbol{R}_i^{
ightarrow} \cdot oldsymbol{R}_i^{
ightarrow},$$

which implies

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

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

$$m{R}^{
ightarrow} = \sum_{i=1}^d m{R}_i^{
ightarrow} \quad ext{where} \quad m{R}_i^{
ightarrow} = \underbrace{m{I} \otimes \ldots \otimes m{I}}_{i-1 ext{ terms}} \otimes m{A} \otimes \underbrace{m{I} \otimes \ldots \otimes m{I}}_{d-i ext{ terms}},$$

where \otimes denotes the Kronecker product, I denotes the identity matrix on $\mathbb{R}^{K \times K}$, and 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 A can be decomposed as

$$\begin{bmatrix} -\boldsymbol{I}_{K-1} & \boldsymbol{0} \\ \boldsymbol{1}_{1\times(K-1)} & 0 \end{bmatrix} = \underbrace{\begin{bmatrix} \boldsymbol{I}_{K-1} & \boldsymbol{0} \\ -\boldsymbol{1}_{1\times(K-1)} & 1 \end{bmatrix}}_{\boldsymbol{U}} \cdot \underbrace{\begin{bmatrix} -\boldsymbol{I}_{K-1} & \boldsymbol{0} \\ \boldsymbol{0} & 0 \end{bmatrix}}_{\boldsymbol{\Lambda}} \cdot \underbrace{\begin{bmatrix} \boldsymbol{I}_{K-1} & \boldsymbol{0} \\ \boldsymbol{1}_{1\times(K-1)} & 1 \end{bmatrix}}_{\boldsymbol{U}^{-1}} \quad \text{where} \quad \boldsymbol{U}\boldsymbol{U}^{-1} = \boldsymbol{U}^{-1}\boldsymbol{U} = \boldsymbol{I}_{K}.$$
(34)

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

$$\boldsymbol{R}_i^{\rightarrow} = \underbrace{(\boldsymbol{U}\boldsymbol{U}^{-1}) \otimes \ldots \otimes (\boldsymbol{U}\boldsymbol{U}^{-1})}_{i-1 \text{ terms}} \otimes (\boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^{-1}) \otimes (\boldsymbol{U}\boldsymbol{U}^{-1}) \otimes \ldots (\boldsymbol{U}\boldsymbol{U}^{-1})$$

$$= (\boldsymbol{U} \otimes \ldots \otimes \boldsymbol{U}) \cdot \left(\underbrace{\boldsymbol{I} \otimes \ldots \otimes \boldsymbol{I}}_{i-1 \text{ terms}} \otimes \Lambda \otimes \boldsymbol{I} \ldots \otimes \boldsymbol{I} \right) \cdot \left(\boldsymbol{U}^{-1} \otimes \ldots \otimes \boldsymbol{U}^{-1} \right) \coloneqq \boldsymbol{U}^{\otimes d} \cdot \Lambda_i \cdot (\boldsymbol{U}^{-1})^{\otimes d}$$

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

$$\begin{split} \boldsymbol{R}_{i}^{\rightarrow} \cdot \boldsymbol{R}_{j}^{\rightarrow} &= \boldsymbol{U}^{\otimes d} \cdot \Lambda_{i} \cdot (\boldsymbol{U}^{-1})^{\otimes d} \cdot \boldsymbol{U}^{\otimes d} \cdot \Lambda_{j} \cdot (\boldsymbol{U}^{-1})^{\otimes d} = \boldsymbol{U}^{\otimes d} \cdot \Lambda_{i} \cdot \Lambda_{j} \cdot (\boldsymbol{U}^{-1})^{\otimes d} \\ &= \boldsymbol{U}^{\otimes d} \cdot \Lambda_{i} \cdot \Lambda_{i} \cdot (\boldsymbol{U}^{-1})^{\otimes d} = \boldsymbol{U}^{\otimes d} \cdot \Lambda_{i} \cdot (\boldsymbol{U}^{-1})^{\otimes d} \cdot \boldsymbol{U}^{\otimes d} \cdot \Lambda_{j} \cdot (\boldsymbol{U}^{-1})^{\otimes d} = \boldsymbol{R}_{i}^{\rightarrow} \cdot \boldsymbol{R}_{i}^{\rightarrow}. \end{split}$$

where the second and forth equations follows from Lemma 13 and Eq. (34).

For the property about the matrix exponential, we start from investigating the case of two commuting matrices, i.e., R_1^{\rightarrow} and R_2^{\rightarrow} . By definition, we have

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

where the last equation establishes since R_1^{\rightarrow} and R_2^{\rightarrow} are commute. Then, we have

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

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

$$\exp(\mathbf{R}_{1}^{\rightarrow}) \cdot \exp(\mathbf{R}_{2}^{\rightarrow}) = [\exp(\mathbf{A}) \otimes \mathbf{I} \otimes \ldots \otimes \mathbf{I}] \cdot [\mathbf{I} \otimes \exp(\mathbf{A}) \otimes \ldots \otimes \mathbf{I}]$$
$$= [\exp(\mathbf{A}) \cdot \mathbf{I}] \otimes [\mathbf{I} \cdot \exp(\mathbf{A})] \otimes \mathbf{I} \ldots \otimes \mathbf{I}.$$

Then, following a recursive manner, we have

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

hence the proof is completed.

 Lemma 7. Suppose matrix **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, $\bar{A}(t) := \exp(tA)$ can be considered as the close solution of the following matrix ODE, i.e.,

$$\frac{\mathrm{d}\bar{\boldsymbol{A}}(t)}{\mathrm{d}t} = \boldsymbol{A} \cdot \bar{\boldsymbol{A}}(t), \quad \text{where} \quad \bar{\boldsymbol{A}}(0) = \boldsymbol{I}. \tag{35}$$

To provide a close form of A_t , we first decompose the matrix A as follows

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

Then, the ODE. (35) can be equivalently think column-by-column, the j-th column of $\bar{A}(t)$ solves

$$\frac{\mathrm{d}}{\mathrm{d}t}\bar{a}(t) = A\bar{a}(t)$$
 where $a(0) = e_j$.

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

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

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

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

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

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

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

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

which implies

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

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

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

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 R^{\to} given in Eq. (7). Then, for any two timestamps $s \leq t$, the forward transition probability satisfies, for any $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$,

$$q_{t|s}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}') = \prod_{i=1}^{d} \left[\delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}') + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}')\right) \cdot \delta_{0}(\boldsymbol{y}_{i} - \boldsymbol{y}_{i}') \cdot e^{-(t-s)} + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}')\right) \cdot \delta_{K}(\boldsymbol{y}_{i}) \cdot (1 - e^{-(t-s)}) \right].$$
(36)

Proof. Under the matrix presentation, Eq. (25) implies the transition matrix $Q_{t|s}^{\rightarrow}$ can be considered as the solution of the ODE

$$\mathrm{d} oldsymbol{Q}_{t|s}^{
ightarrow}/\mathrm{d} t = oldsymbol{R}^{
ightarrow} \cdot oldsymbol{Q}_{t|s}^{
ightarrow} \quad ext{where} \quad oldsymbol{Q}_{s|s}^{
ightarrow} = oldsymbol{I}.$$

Combining Lemma 4 and 6, we have

$$\mathbf{Q}_{t|s}^{\rightarrow} = \exp\left((t-s)\mathbf{R}^{\rightarrow}\right) = \exp\left((t-s)\mathbf{A}\right)^{\otimes d},\tag{37}$$

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 $\exp((t-s)A)$ shown in Lemma 7. Combining this result with the calculation of the Kronecker product Lemma 12, we have

$$q_{t|s}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}') = \prod_{i=1}^{d} \left[\delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}') + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}')\right) \cdot \delta_{0}(\boldsymbol{y}_{i} - \boldsymbol{y}_{i}') \cdot e^{-(t-s)} + \left(1 - \delta_{(K,K)}(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}')\right) \cdot \delta_{K}(\boldsymbol{y}_{i}) \cdot (1 - e^{-(t-s)}) \right].$$

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

The proof of Lemma 2. According to Eq. (28), the solution of q_t^{\rightarrow} can be calculated as

$$m{q}_t^{
ightarrow} = \exp(tm{R}^{
ightarrow}) \cdot m{q}_0^{
ightarrow} = \exp(tm{A})^{\otimes d} \cdot m{q}_0^{
ightarrow} = egin{bmatrix} e^{-t} & 0 & \dots & 0 & 0 \ 0 & e^{-t} & \dots & 0 & 0 \ dots & dots & \ddots & dots & dots \ 0 & 0 & \dots & e^{-t} & 0 \ 1 - e^{-t} & 1 - e^{-t} & \dots & 1 - e^{-t} & 1 \end{bmatrix}^{\otimes d}$$

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

$$\boldsymbol{q}_{t}^{\rightarrow}(\boldsymbol{y}) = \sum_{\boldsymbol{y}' \in \mathcal{Y}} \exp(t\boldsymbol{A})^{\otimes d}(\boldsymbol{y}, \boldsymbol{y}') \cdot \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}') = \sum_{\boldsymbol{y}'} \left[\prod_{i=1}^{d} \exp(t\boldsymbol{A})(\boldsymbol{y}_{i}, \boldsymbol{y}'_{i}) \right] \cdot \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}'). \tag{38}$$

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

$$oldsymbol{y}_{\mathcal{K}^c(oldsymbol{y})} = oldsymbol{y}_{\mathcal{K}^c(oldsymbol{y})}' \quad \Leftrightarrow \quad oldsymbol{y}_i = oldsymbol{y}_i' \ orall \ i
ot\in \mathcal{K}(oldsymbol{y}).$$

Then, Eq. (38) can be rewritten as

$$\begin{split} \boldsymbol{q}_{t}^{\rightarrow}(\boldsymbol{y}) &= \sum_{\boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}^{\rightarrow} = \boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}} \left[\prod_{j \notin \mathcal{K}} \exp(t\boldsymbol{A})(\boldsymbol{y}_{j}, \boldsymbol{y}_{j}^{\prime}) \cdot \prod_{j \neq i}^{d} \exp(t\boldsymbol{A})(K, \boldsymbol{y}_{j}^{\prime}) \right] \cdot \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}^{\prime}) \\ &+ \sum_{\boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}^{\prime} \neq \boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}} \left[\prod_{j=1}^{d} \exp(t\boldsymbol{A})(\boldsymbol{y}_{j}, \boldsymbol{y}_{j}^{\prime}) \right] \cdot \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}^{\prime}) \\ &= \sum_{\boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}^{\prime} = \boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}} \left[e^{-t \cdot |\mathcal{K}^{c}(\boldsymbol{y})|} \cdot (1 - e^{-t})^{|\mathcal{K}(\boldsymbol{y})|} \right] \cdot \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}^{\prime}) \\ &\leq e^{-t \cdot (d - \operatorname{numK}(\boldsymbol{y}))} \cdot \sum_{\boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}^{\prime} = \boldsymbol{y}_{\mathcal{K}^{c}(\boldsymbol{y})}} \boldsymbol{q}_{0}^{\rightarrow}(\boldsymbol{y}^{\prime}) \leq \exp(-t \cdot (d - \operatorname{numK}(\boldsymbol{y}))), \end{split}$$

where the second equation establishes since we have

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

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

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

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

$$\operatorname{KL}\left(q_{t}^{\rightarrow} \| \tilde{q}_{t}\right) = \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\rightarrow}(\boldsymbol{y}) \cdot \ln \frac{q_{t}^{\rightarrow}(\boldsymbol{y})}{\tilde{q}_{t}(\boldsymbol{y})} = q_{t}^{\rightarrow}([K, \dots, K]) \cdot \ln \frac{q_{t}^{\rightarrow}([K, \dots, K])}{\tilde{q}_{t}([K, \dots, K])} + \sum_{\boldsymbol{y} \neq [K, \dots, K]} q_{t}^{\rightarrow}(\boldsymbol{y}) \cdot \ln \frac{q_{t}^{\rightarrow}(\boldsymbol{y})}{\tilde{q}_{t}(\boldsymbol{y})}$$

$$\leq \ln \tilde{Z}_{t} + \sum_{\boldsymbol{y} \neq [K, \dots, K]} q_{t}^{\rightarrow}(\boldsymbol{y}) \ln \frac{q_{t}^{\rightarrow}(\boldsymbol{y})}{\exp(-t \cdot (d - \operatorname{numK}(\boldsymbol{y}))) / \tilde{Z}_{t}} = \ln \tilde{Z}_{t} + \sum_{\boldsymbol{y} \neq [K, \dots, K]} q_{t}^{\rightarrow}(\boldsymbol{y}) \ln \tilde{Z}_{t}$$

$$\leq 2 \ln \tilde{Z}_{t} = 2 \ln \left[1 + (1 + e^{-t})^{d} - 1\right] \leq 2 \cdot (1 + e^{-t})^{d} - 2.$$

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

$$\mathrm{KL}\left(\boldsymbol{q}_{t}^{\rightarrow} \| \tilde{q}_{t}\right) \leq \epsilon \quad \Leftrightarrow \quad (1 + e^{-t})^{d} - 1 \leq \epsilon/2 \quad \Leftrightarrow \quad d\ln(1 + e^{-t}) \leq \ln(1 + \epsilon/2),$$

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

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

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

C EULER DISCRETIZATION ANALYSIS

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

[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 $T = nh + \delta$:

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

C.1 Proof of Theorem 1

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

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

Without loss of generality, assume that $\tilde{R}_t(\boldsymbol{y}',\boldsymbol{y})^{\top}$ satisfies the two sufficient conditions of the transition rate matrix: its off-diagonal entries are non-negative, and each row sums to zero 1 . 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

$$Q_{t,t+h}^{Eu}(\boldsymbol{y},\boldsymbol{y}') = q_{t+h|t}^{Eu}(\boldsymbol{y}'|\boldsymbol{y}) = \delta_{\boldsymbol{y}}(\boldsymbol{y}') + h \cdot \tilde{R}_{t}(\boldsymbol{y}',\boldsymbol{y})$$
(39)

To prove the convergence bound for TV $\left(q_{\delta}^{\leftarrow},q_{T-\delta}^{Eu}\right)$, 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 $\left(q_{T-\delta}^{Eu},q_{T-\delta}^{EI}\right)$ and TV $\left(q_{\delta}^{\leftarrow},q_{T-\delta}^{EI}\right)$ 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

$$\operatorname{TV}\left(q_{(k+1)h}^{Eu}, q_{(k+1)h}^{EI}\right) \leq \operatorname{TV}\left(q_{kh}^{Eu}, q_{kh}^{EI}\right) + \mathbb{E}_{\boldsymbol{y} \sim q_{kh}^{Eu}} \operatorname{TV}\left(q_{(k+1)h|kh}^{Eu}(\cdot \mid \boldsymbol{y}), q_{(k+1)h|kh}^{EI}(\cdot \mid \boldsymbol{y})\right) \tag{40}$$

By the definition of total variation distance, we have

$$\operatorname{TV}\left(q_{(k+1)h|kh}^{Eu}(\cdot\mid\boldsymbol{y}), q_{(k+1)h|kh}^{EI}(\cdot\mid\boldsymbol{y})\right) = \sum_{\boldsymbol{y}'} \left| q_{(k+1)h|kh}^{Eu}(\boldsymbol{y}'\mid\boldsymbol{y}) - q_{(k+1)h|kh}^{EI}(\boldsymbol{y}'\mid\boldsymbol{y}) \right|$$

$$= \sum_{\boldsymbol{y}'} \left| Q_{kh,(k+1)h}^{Eu}(\boldsymbol{y},\boldsymbol{y}') - Q_{kh,(k+1)h}^{EI}(\boldsymbol{y},\boldsymbol{y}') \right|$$
(41)

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

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

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

Therefore,

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

when $h \|\tilde{R}_{kh}\|_{1} \leq 1$. Plugging this into Eq. (40) and (41), we have

$$TV\left(q_{(k+1)h}^{Eu}, q_{(k+1)h}^{EI}\right) \le TV\left(q_{kh}^{Eu}, q_{kh}^{EI}\right) + \left(h \left\|\tilde{R}_{kh}\right\|_{1}\right)^{2},\tag{42}$$

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

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

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

Thus $\|\tilde{R}_{kh}\|_{_1} \leq 2KdM$. By (42) we have

$$\operatorname{TV}\left(q_{nh}^{Eu}, q_{nh}^{EI}\right) \leq \operatorname{TV}\left(q_0^{Eu}, q_0^{EI}\right) + \sum_{k=1}^n \left(h \left\|\tilde{R}_{kh}\right\|_1\right)^2 \\
\leq K^2 d^2 \sum_{k=1}^n h^2 M^2 \leq K^2 d^2 n h^2 M^2 \leq K^2 (T - \delta) h d^2 M^2.$$
(43)

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

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

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

$$\operatorname{KL}\left(q_{(k+1)h}^{\leftarrow} \left\| q_{(k+1)h}^{EI} \right) = \operatorname{KL}\left(q_{kh}^{\leftarrow} \left\| q_{kh}^{EI} \right) + \int_{kh}^{(k+1)h} \frac{\operatorname{dKL}\left(q_t^{\leftarrow} \left\| q_t^{EI} \right)}{\operatorname{d}t} \operatorname{d}t.$$
(44)

By the chain rule of KL divergence (Lemma 15)

$$\frac{\mathrm{d}}{\mathrm{d}t} \mathrm{KL}\left(q_{t}^{\leftarrow} \| q_{t}^{EI}\right) = \lim_{\Delta t \to 0} \frac{\mathrm{KL}\left(q_{t+\Delta t}^{\leftarrow} \| q_{t+\Delta t}^{EI}\right) - \mathrm{KL}\left(q_{t}^{\leftarrow} \| q_{t}^{EI}\right)}{\Delta t}$$

$$\leq \lim_{\Delta t \to 0} \mathbb{E}_{\boldsymbol{y} \sim q_{t}^{\leftarrow}} \frac{\mathrm{KL}\left(q_{t+\Delta t|t}^{\leftarrow} (\cdot \mid \boldsymbol{y}) \| q_{t+\Delta t|t}^{EI} (\cdot \mid \boldsymbol{y})\right)}{\Delta t}$$

$$= \mathbb{E}_{\boldsymbol{y} \sim q_{t}^{\leftarrow}} \underbrace{\lim_{\Delta t \to 0} \frac{\mathrm{KL}\left(q_{t+\Delta t|t}^{\leftarrow} (\cdot \mid \boldsymbol{y}) \| q_{t+\Delta t|t}^{EI} (\cdot \mid \boldsymbol{y})\right)}{\Delta t}}_{\text{Term 1}}$$
(45)

For each $y \in \mathcal{Y}$, we focus on Term 1 of Eq. (45), and have

Term 1 =
$$\lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \sum_{\mathbf{y}' \in \mathcal{Y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \cdot \ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{q_{t+\Delta t|t}^{EI}(\mathbf{y}'|\mathbf{y})} \right]$$

$$= \lim_{\Delta t \to 0} \left[\sum_{\mathbf{y}' \neq \mathbf{y}} \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{q_{t+\Delta t|t}^{EI}(\mathbf{y}'|\mathbf{y})} \right] + \lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right) \cdot \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})} \right] .$$
(46)

For Term 1.1, we have

Term 1.1 =
$$\sum_{\mathbf{y}' \neq \mathbf{y}} \lim_{\Delta t \to 0} \left[\frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{\Delta t} \right] \cdot \lim_{\Delta t \to 0} \left[\ln \frac{q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y})}{q_{t+\Delta t|t}^{EI}} \right]$$

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

$$= \sum_{\mathbf{y}' \neq \mathbf{y}} R_t^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_t^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_{kh}(\mathbf{y}', \mathbf{y})},$$
(47)

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

Term 1.2 =
$$\lim_{\Delta t \to 0} \left[1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y}) \right] \cdot \lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y})}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{EI}(\boldsymbol{y}'|\boldsymbol{y})} \right]$$

$$= \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \left(\tilde{R}_{kh}(\boldsymbol{y}', \boldsymbol{y}) - R_t^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y}) \right) = \tilde{R}_{kh}(\boldsymbol{y}) - R_t^{\leftarrow}(\boldsymbol{y})$$
(48)

where the first inequality follows from Lemma 9. Plugging Eq. (47), Eq. (48) and Eq. (46), into Eq. (45) we have

1299
1300
$$\frac{\operatorname{dKL}\left(q_{t}^{\leftarrow} \| q_{t}^{EI}\right)}{\operatorname{d}t} \leq \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \left(\sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_{t}^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y})}{\tilde{R}_{kh}(\boldsymbol{y}', \boldsymbol{y})} + \tilde{R}_{kh}(\boldsymbol{y}) - R_{t}^{\leftarrow}(\boldsymbol{y})\right)$$
1302
$$= \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \left(\sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_{t}^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y})}{\tilde{R}_{kh}(\boldsymbol{y}', \boldsymbol{y})} + \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \tilde{R}_{kh}(\boldsymbol{y}', \boldsymbol{y}) - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_{t}^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y})\right)$$
1305
$$= \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') \cdot \left[-\frac{q_{t}^{\leftarrow}(\boldsymbol{y}')}{q_{t}^{\leftarrow}(\boldsymbol{y})} + \tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}') + \frac{q_{t}^{\leftarrow}(\boldsymbol{y}')}{q_{t}^{\leftarrow}(\boldsymbol{y})} \ln \frac{q_{t}^{\leftarrow}(\boldsymbol{y}')}{q_{t}^{\leftarrow}(\boldsymbol{y})\tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}')} \right]$$
1308
$$= \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') \cdot \left[-v_{t,\boldsymbol{y}}(\boldsymbol{y}') + \tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}') + v_{t,\boldsymbol{y}}(\boldsymbol{y}') \ln \frac{v_{t,\boldsymbol{y}}(\boldsymbol{y}')}{\tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}')} \right]$$
1310
$$= \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}') \cdot \left[-v_{t,\boldsymbol{y}}(\boldsymbol{y}') + \tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}') + v_{t,\boldsymbol{y}}(\boldsymbol{y}') \ln \frac{v_{t,\boldsymbol{y}}(\boldsymbol{y}')}{\tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}')} \right]$$
1311
$$= \sum_{\boldsymbol{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\boldsymbol{y}) \cdot \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \sum_{\boldsymbol{y} \in \mathcal{Y}} \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \left[-v_{t,\boldsymbol{y}}(\boldsymbol{y}') + \tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}') + v_{t,\boldsymbol{y}}(\boldsymbol{y}') \ln \frac{v_{t,\boldsymbol{y}}(\boldsymbol{y}')}{\tilde{v}_{kh,\boldsymbol{y}}(\boldsymbol{y}')} \right]$$
1314
1315

For $y'=y[y_i\to k]$, by Eq. (54) we have $v_{t,y}(y')=\frac{q_t^\leftarrow(y[y_i\to k])}{q_t^\leftarrow(y)}\leq \frac{1}{e^{(T-t)}-1}$. By (Liang et al., 2025, Lemma 2), there exist c>0 such that $v_{t,y}(y')\geq \frac{1}{c}e^{-(T-t)}$. Therefore, by (Zhang et al., 2024, Proposition 3), letting $C=\max\{M,ce^T\}$, Term 2 satisfies

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

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

Term
$$2 \lesssim CK$$
 numK $(\boldsymbol{y}) h^2 + C^2 \sum_{\text{Ham}(\boldsymbol{y}, \boldsymbol{y}') = 1 \text{ and } \boldsymbol{y}_{\text{DiffIdx}(\boldsymbol{y}, \boldsymbol{y}')} = K} (D_{\phi}(v_{kh, \boldsymbol{y}}(\boldsymbol{y}') || \tilde{v}_{kh, \boldsymbol{y}}(\boldsymbol{y}')))$

$$\lesssim CKdh^2 + C^2 \sum_{\text{Ham}(\boldsymbol{y}, \boldsymbol{y}') = 1 \text{ and } \boldsymbol{y}_{\text{DiffIdx}(\boldsymbol{y}, \boldsymbol{y}')} = K} (D_{\phi}(v_{kh, \boldsymbol{y}}(\boldsymbol{y}') || \tilde{v}_{kh, \boldsymbol{y}}(\boldsymbol{y}')))$$
(50)

Since by Eq. (44), we have

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

Then, by Eq. (49), Eq. (50), Eq. (6) and Assumption [A1]-, we have

$$\mathrm{KL}\left(q_{T-\delta}^{\leftarrow} \middle\| q_{T-\delta}^{EI}\right) \lesssim (T-\delta)C^2\epsilon_{\mathrm{score}}^2 + C(T-\delta)Kdh^2.$$

By Pinsker's inequality, we have

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

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

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

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

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 $\operatorname{TV}(q_{\delta}^{\leftarrow}, q_{T-\delta}^{Eu}) \lesssim \varepsilon$.

Hence, the proof is completed.

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

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

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

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

which automatically makes

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

Under this condition, we have

$$\ln \frac{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y})}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{EI}(\boldsymbol{y}'|\boldsymbol{y})} = \ln \left[1 + \frac{\sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \left(q_{t+\Delta t|t}^{EI}(\boldsymbol{y}'|\boldsymbol{y}) - q_{t+\Delta t|t}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y}) \right)}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t}^{EI}(\boldsymbol{y}'|\boldsymbol{y})} \right]$$

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

which implies (with the dominated convergence theorem)

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

Only when i = 1, we have

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

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

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

Hence, the proof is completed.

D TRUNCATED UNIFORMIZATION INFERENCE ANALYSIS

D.1 THE PROOF OF LEMMA 3

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

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

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

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

Define the function to transfer the *i*-th element of $y(y_i)$ from k' to k as

$$y[y_i: k' \to k] = [y_1, y_2, \dots, y_{i-1}, k, y_{i+1}, \dots, y_d]$$

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

$$R_t^{\leftarrow}(\boldsymbol{y}) = \sum_{i, \boldsymbol{y}_i = K} \left[\sum_{k=1}^{K-1} R^{\rightarrow}(\boldsymbol{y}, \boldsymbol{y}[\boldsymbol{y}_i : K \rightarrow k]) \cdot \frac{q_t^{\leftarrow}(\boldsymbol{y}[\boldsymbol{y}_i : K \rightarrow k])}{q_t^{\leftarrow}(\boldsymbol{y})} \right]. \tag{51}$$

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

$$\frac{q_t^{\leftarrow}(\boldsymbol{y}[\boldsymbol{y}_i\colon K\to k])}{q_t^{\leftarrow}(\boldsymbol{y})} = \frac{q_{T-t}^{\rightarrow}(\boldsymbol{y}[\boldsymbol{y}_i\colon K\to k])}{q_{T-t}^{\rightarrow}(\boldsymbol{y})} = \frac{\sum_{\boldsymbol{y}_0\in\mathcal{Y}}q_0^{\rightarrow}(\boldsymbol{y}_0)\cdot q_{T-t|0}^{\rightarrow}(\boldsymbol{y}[\boldsymbol{y}_i\colon K\to k]|\boldsymbol{y}_0)}{\sum_{\boldsymbol{y}_0\in\mathcal{Y}}q_0^{\rightarrow}(\boldsymbol{y}_0)\cdot q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_0)}$$

$$= \frac{\sum_{\boldsymbol{y}_{0} \in \mathcal{Y}} q_{0}^{\rightarrow}(\boldsymbol{y}_{0}) \cdot q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0}) \cdot \frac{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0}: K \rightarrow k]|\boldsymbol{y}_{0})}{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}_{0}) \cdot q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0})}} = \mathbb{E}_{\boldsymbol{y}_{0} \sim q_{0|T-t}^{\rightarrow}(\cdot|\boldsymbol{y})} \left[\frac{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0}: K \rightarrow k]|\boldsymbol{y}_{0})}{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0})} \right],$$
(52)

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

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

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

$$q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0}) = \prod_{j \in \{1,...,i-1,i+1,...,d\}} \left[\mathbf{1}_{(K,K)}(\boldsymbol{y}_{j},\boldsymbol{y}_{0,j}) + \left(1 - \mathbf{1}_{(K,K)}(\boldsymbol{y}_{j},\boldsymbol{y}_{0,j})\right) \cdot \mathbf{1}_{0}(\boldsymbol{y}_{j} - \boldsymbol{y}_{0,j}) \cdot e^{-(T-t)} + \left(1 - \mathbf{1}_{(K,K)}(\boldsymbol{y}_{j},\boldsymbol{y}_{0,j})\right) \cdot \mathbf{1}_{K}(\boldsymbol{y}_{j}) \cdot (1 - e^{-(T-t)}) \right] \cdot \left[\mathbf{1}_{(K,K)}(K,\boldsymbol{y}_{0,i}) + \left(1 - \mathbf{1}_{(K,K)}(K,\boldsymbol{y}_{0,i})\right) \cdot (1 - e^{-(T-t)}) \right]$$

and

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

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

$$\frac{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}[\boldsymbol{y}_{i} \to k]|\boldsymbol{y}_{0})}{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_{0})} = \frac{\left(1 - \mathbf{1}_{(K,K)}(k, \boldsymbol{y}_{0,i})\right) \cdot \mathbf{1}_{0}(k - \boldsymbol{y}_{0,i}) \cdot e^{-(T-t)}}{\mathbf{1}_{(K,K)}(K, \boldsymbol{y}_{0,i}) + \left(1 - \mathbf{1}_{(K,K)}(K, \boldsymbol{y}_{0,i})\right) \cdot \left(1 - e^{-(T-t)}\right)}
= \frac{\mathbf{1}_{0}(k - \boldsymbol{y}_{0,i}) \cdot e^{-(T-t)}}{1 - e^{-(T-t)}} \le \frac{e^{-(T-t)}}{1 - e^{-(T-t)}} = \frac{1}{e^{(T-t)} - 1}.$$
(53)

Plugging this result into Eq. (52), the density ratio of the reverse process will have

$$\frac{q_t^{\leftarrow}(\boldsymbol{y}[\boldsymbol{y}_i \to k])}{q_t^{\leftarrow}(\boldsymbol{y})} = \mathbb{E}_{\boldsymbol{y}_0 \sim q_{0|T-t}^{\rightarrow}(\cdot|\boldsymbol{y})} \left[\frac{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}[\boldsymbol{y}_i : K \to k]|\boldsymbol{y}_0)}{q_{T-t|0}^{\rightarrow}(\boldsymbol{y}|\boldsymbol{y}_0)} \right] \le \frac{1}{e^{(T-t)} - 1}.$$
 (54)

Combining with the fact, i.e.,

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

from Eq. (7), Eq. (51) can be upper bounded as

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

Hence, the proof is completed.

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

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

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

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

The latter upper bound is strictly better compared with the former one, since the number of mask tokens, i.e., $\operatorname{numK}(\boldsymbol{y}) \leq d$. Besides, with the time growth (from 0 to T), $\operatorname{numK}(\boldsymbol{y})$ will be monotonic decrease for $R_t^{\leftarrow}(\boldsymbol{y})$ (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 $0 < t_0 \le t$, we have $q_{t|t_0}^{\leftarrow}(\boldsymbol{y}|\boldsymbol{y}_0) \ne 0$ if and only if

$$y \in \mathcal{Y}^{\leftarrow}(y_0) = \{y' | \forall i, \quad y_{0,i} = K \text{ or } y'_i = y_{0,i} \}.$$

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

$$q_{t,t_0}^{\leftarrow}(\boldsymbol{y},\boldsymbol{y}_0) = q_{t|t_0}^{\leftarrow}(\boldsymbol{y}|\boldsymbol{y}_0) \cdot q_{t_0}^{\leftarrow}(\boldsymbol{y}_0) = q_{T-t,T-t_0}^{\rightarrow}(\boldsymbol{y},\boldsymbol{y}_0)$$

$$= q_{T-t_0,T-t}^{\rightarrow}(\boldsymbol{y}_0,\boldsymbol{y}) = q_{T-t_0|T-t}^{\rightarrow}(\boldsymbol{y}_0|\boldsymbol{y}) \cdot q_{T-t}^{\rightarrow}(\boldsymbol{y}),$$
(55)

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

$$q_{T-t_0|T-t}^{\rightarrow}(\boldsymbol{y}_0|\boldsymbol{y}) \neq 0 \quad \Leftrightarrow \quad \boldsymbol{y}_0 \in \mathcal{Y}^{\rightarrow}(\boldsymbol{y}) = \{\boldsymbol{y}'| \ \forall i, \quad \boldsymbol{y}_i' = \boldsymbol{y}_i \text{ or } \boldsymbol{y}_i' = K\}.$$
 (56)

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

$$q_{t_0}^{\leftarrow}(\boldsymbol{y}_0) = q_{T-t_0}^{\rightarrow}(\boldsymbol{y}_0) > 0 \quad \text{and} \quad q_t^{\rightarrow}(\boldsymbol{y}) > 0.$$
 (57)

Then, we can summarize

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

Hence, the proof is completed.

D.2 THE CONVERGENCE OF ALG. 1

 Suppose, with the infinitesimal reverse transition rate, the particles in Alg. 1 during the reverse process are denotes 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 $\hat{\mathbf{y}}_t = \mathbf{y}$,

- 1519 1. With probability $\Delta t \cdot \beta_{t_w}(y_0)$, allow a state transition.
 - 2. Conditioning on an allowed transition, move from y to y' with probability

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

Here we should note that

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

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

$$\hat{q}_{t+\Delta t|t,t_{w-1}}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0}) = \begin{cases} \Delta t \cdot \beta_{t_{w}}(\boldsymbol{y}_{0}) \cdot \hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y}',\boldsymbol{y}) \cdot \beta_{t_{w}}^{-1}(\boldsymbol{y}_{0}) & \boldsymbol{y}' \neq \boldsymbol{y} \\ 1 - \Delta t \cdot \beta_{t_{w}}(\boldsymbol{y}_{0}) + \Delta t \cdot \beta_{t_{w}}(\boldsymbol{y}_{0}) \cdot (1 - \beta_{t_{w}}(\boldsymbol{y}_{0})^{-1} \cdot \hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y})) & \boldsymbol{y}' = \boldsymbol{y} \end{cases}$$

$$= \begin{cases} \Delta t \cdot \hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y}',\boldsymbol{y}) & \boldsymbol{y}' \neq \boldsymbol{y} \\ 1 - \Delta t \cdot \hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y}) & \boldsymbol{y}' = \boldsymbol{y} \end{cases} . \tag{58}$$

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

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

Proof. Since we have required $\Delta t \to 0$, for any $y' \neq y$, it has

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

where the first row follows from Eq. (58) and the second row follows from Lemma. 1. This automatically makes

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

Under this condition, we have

$$\frac{1560}{1561} \quad \ln \frac{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0})} = \ln \left[1 + \frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0}) - q_{t+\Delta t|t}^{\leftarrow}(\mathbf{y}'|\mathbf{y}) \right)}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0})} \right] \\
\frac{1563}{1564} \quad = \sum_{i=1}^{\infty} \frac{(-1)^{i+1}}{i} \cdot \left[\frac{\sum_{\mathbf{y}' \neq \mathbf{y}} \left(\hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0}) - q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0})}{1 - \sum_{\mathbf{y}' \neq \mathbf{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\mathbf{y}'|\mathbf{y},\mathbf{y}_{0})} \right]^{i},$$

which implies (with the dominated convergence theorem)

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

Only when i = 1, we have

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

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

$$\lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right]$$

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

$$= \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \left(\hat{R}_{t,\boldsymbol{y}_{0}}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y}) - R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y}) \right) = \hat{R}_{t,\boldsymbol{y}_{0}}^{\leftarrow}(\boldsymbol{y}) - R_{t}^{\leftarrow}(\boldsymbol{y}),$$

where the second equation follows from Eq. (58) and the second row follows from Lemma. 1. Hence, the proof is completed.

Theorem 3 (The convergence of Alg. 1). Suppose Assumption [A1] and [A2] hold, if Alg. 1 has $t_0=0, \quad t_W=T-\delta, \quad and \quad \epsilon_{score} \leq T^{-1/2} \cdot \epsilon \quad where \quad T=\ln(4d/\epsilon^2) \quad and \quad \delta \leq d^{-1}\epsilon,$ the TV distance between the target discrete distribution q_* and the underlying distribution of the output particle $\hat{q}_{T-\delta}$ will satisfy $\operatorname{TV}(q_*,\hat{q}_{T-\delta}) \leq 2\epsilon$.

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

$$\operatorname{KL}\left(q_{t_{w}}^{\leftarrow} \| \hat{q}_{t_{w}}\right) \leq \operatorname{KL}\left(q_{t_{w-1}}^{\leftarrow} \| \hat{q}_{t_{w-1}}\right) + \mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[\operatorname{KL}\left(q_{t_{w}|t_{w-1}}^{\leftarrow}(\cdot | \mathbf{y}_{0}) \| \hat{q}_{t_{w}|t_{w-1}}(\cdot | \mathbf{y}_{0})\right) \right]$$

$$= \operatorname{KL}\left(q_{t_{w-1}}^{\leftarrow} \| \hat{q}_{t_{w-1}}\right) + \int_{t_{w-1}}^{t_{w}} d\mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[\operatorname{KL}\left(q_{t|t_{w-1}}^{\leftarrow}(\cdot | \mathbf{y}_{0}) \| \hat{q}_{t|t_{w-1}}(\cdot | \mathbf{y}_{0})\right) \right]$$

$$(59)$$

Then it has

$$d\mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[KL \left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0}^{\leftarrow}) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_{0}^{\leftarrow}) \right) \right] / dt$$

$$= \lim_{\Delta \to 0} (\Delta t)^{-1} \cdot \mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[KL \left(q_{t+\Delta t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0}) \| \hat{q}_{t+\Delta t|t_{w-1}}(\cdot|\mathbf{y}_{0}) \right) - KL \left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0}) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_{0}) \right) \right]$$

$$\leq \lim_{\Delta \to 0} (\Delta t)^{-1} \cdot \mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[\mathbb{E}_{\mathbf{y} \sim q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0})} \left(KL \left(q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y},\mathbf{y}_{0}) \| \hat{q}_{t+\Delta t|t,t_{w-1}}(\cdot|\mathbf{y},\mathbf{y}_{0}) \right) \right) \right]$$

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

$$\frac{d\mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[KL \left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0}) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_{0}) \right) \right] / dt}{\leq \sum_{\mathbf{y} \in \mathcal{Y}, \mathbf{y}_{0} \in \mathcal{Y}^{\rightarrow}(\mathbf{y})} q_{t,t_{w-1}}^{\leftarrow}(\mathbf{y}, \mathbf{y}_{0}) \cdot \underbrace{\lim_{\Delta t \to 0} \left[\frac{KL \left(q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}, \mathbf{y}_{0}) \| \hat{q}_{t+\Delta t|t,t_{w-1}}(\cdot|\mathbf{y}, \mathbf{y}_{0}) \right)}{\Delta t} \right]}_{\text{Term 1}} \tag{60}$$

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. (60), and have

Term 1 =
$$\lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \sum_{\boldsymbol{y}' \in \mathcal{Y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0}) \cdot \ln \frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\hat{q}_{t+\Delta t|t,t_{w-1}}^{-}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right]$$

$$= \lim_{\Delta t \to 0} \left[\sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\Delta t} \cdot \ln \frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\hat{q}_{t+\Delta t|t,t_{w-1}}^{-}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right] + \lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \left(1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0}) \right) \cdot \ln \frac{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right].$$
Term 1.2 (61)

For Term 1.1, we have

Term 1.1 =
$$\sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \lim_{\Delta t \to 0} \left[\frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\Delta t} \right] \cdot \lim_{\Delta t \to 0} \left[\ln \frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right]$$

$$= \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y}) \cdot \ln \left[\lim_{\Delta t \to 0} \left(\frac{q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{\Delta t} \cdot \frac{\Delta t}{\hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right) \right]$$

$$= \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y})}{\hat{R}_{t,\boldsymbol{y}_{0}}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y})},$$
(62)

where the last equation follows from Lemma 1 and Eq. (58). For Term 1.2, we have

Term 1.2 =
$$\lim_{\Delta t \to 0} \left[1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0}) \right]$$

$$\cdot \lim_{\Delta t \to 0} \left[\Delta t^{-1} \cdot \ln \frac{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} q_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})}{1 - \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} \hat{q}_{t+\Delta t|t,t_{w-1}}^{\leftarrow}(\boldsymbol{y}'|\boldsymbol{y},\boldsymbol{y}_{0})} \right] \leq 1 \cdot (\hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y}) - R_{t}^{\leftarrow}(\boldsymbol{y}))$$
(63)

where the first inequality follows from Lemma 11. Plugging Eq. (62), Eq. (63) and Eq. (61), into Eq. (60) we have

$$\frac{\mathrm{d}\mathbb{E}_{\mathbf{y}_{0}\sim q_{t_{w-1}}^{\leftarrow}}\left[\mathrm{KL}\left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0})\|\hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_{0})\right)\right]/\mathrm{d}t}{\leq \sum_{\boldsymbol{y}\in\mathcal{Y},\boldsymbol{y}_{0}\in\mathcal{Y}^{\rightarrow}(\boldsymbol{y})}q_{t,t_{w-1}}^{\leftarrow}(\boldsymbol{y},\boldsymbol{y}_{0})\cdot\left(\sum_{\boldsymbol{y}'\neq\boldsymbol{y}}R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y})\cdot\ln\frac{R_{t}^{\leftarrow}(\boldsymbol{y}',\boldsymbol{y})}{\hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y}',\boldsymbol{y})}+\hat{R}_{t,\boldsymbol{y}_{0}}(\boldsymbol{y})-R_{t}^{\leftarrow}(\boldsymbol{y})\right).}$$
(64)

Then, for any $y \in \mathcal{Y}$ and $y_0 \in \mathcal{Y}^{\rightarrow}(y)$, we have

$$\sum_{\mathbf{y}'\neq\mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\hat{R}_{t,\mathbf{y}_{0}}(\mathbf{y}', \mathbf{y})} + \hat{R}_{t,\mathbf{y}_{0}}(\mathbf{y}) - R_{t}^{\leftarrow}(\mathbf{y})$$

$$= \sum_{\mathbf{y}'\neq\mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_{t}(\mathbf{y}', \mathbf{y})} + \tilde{R}_{t}(\mathbf{y}) - R_{t}^{\leftarrow}(\mathbf{y})$$

$$+ \sum_{\mathbf{y}'\neq\mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \ln \frac{\tilde{R}_{t}(\mathbf{y}', \mathbf{y})}{\hat{R}_{t,\mathbf{y}_{0}}(\mathbf{y}', \mathbf{y})} + \hat{R}_{t,\mathbf{y}_{0}}(\mathbf{y}) - \tilde{R}_{t}(\mathbf{y}).$$
(65)

When $\tilde{R}_t(y) \leq \beta_{t,...}(y_0)$, due to Eq. (18), we have

$$\hat{R}_{t, oldsymbol{y}_0}(oldsymbol{y}', oldsymbol{y}) = \tilde{R}_t(oldsymbol{y}', oldsymbol{y}) \quad ext{and} \quad \hat{R}_{t, oldsymbol{y}_0}(oldsymbol{y}) = \sum_{oldsymbol{y}'
eq oldsymbol{y}} \hat{R}_{t, oldsymbol{y}_0}(oldsymbol{y}', oldsymbol{y}) = \sum_{oldsymbol{y}'
eq oldsymbol{y}} \hat{R}_t(oldsymbol{y}', oldsymbol{y}) = \hat{R}_t(oldsymbol{y})$$

which implies Term 2 = 0 in Eq. (65). Otherwise, we have

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

which implies

$$\operatorname{Term} 2 = \sum_{\boldsymbol{y}' \neq \boldsymbol{y}} R_t^{\leftarrow}(\boldsymbol{y}', \boldsymbol{y}) \cdot \ln \frac{\tilde{R}_t(\boldsymbol{y})}{\beta_{t_w}(\boldsymbol{y}_0)} + \beta_{t_w}(\boldsymbol{y}_0) - \tilde{R}_t(\boldsymbol{y})$$

$$= R_t^{\leftarrow}(\boldsymbol{y}) \cdot \ln \left[1 + \frac{\tilde{R}_t(\boldsymbol{y}) - \beta_{t_w}(\boldsymbol{y}_0)}{\beta_{t_w}(\boldsymbol{y}_0)} \right] + \beta_{t_w}(\boldsymbol{y}_0) - \tilde{R}_t(\boldsymbol{y})$$

$$\leq \beta_{t_w}(\boldsymbol{y}_0) \cdot \left[\frac{\tilde{R}_t(\boldsymbol{y}) - \beta_{t_w}(\boldsymbol{y}_0)}{\beta_{t_w}(\boldsymbol{y}_0)} \right] + \beta_{t_w}(\boldsymbol{y}_0) - \tilde{R}_t(\boldsymbol{y}) = 0,$$

where the last inequality follows from

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

Combining with Eq. (65) and Eq. (64), we have

$$d\mathbb{E}_{\mathbf{y}_{0} \sim q_{t_{w-1}}^{\leftarrow}} \left[\operatorname{KL} \left(q_{t|t_{w-1}}^{\leftarrow}(\cdot|\mathbf{y}_{0}) \| \hat{q}_{t|t_{w-1}}(\cdot|\mathbf{y}_{0}) \right) \right] / dt$$

$$\leq \sum_{\mathbf{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_{t}(\mathbf{y}', \mathbf{y})} + \tilde{R}_{t}(\mathbf{y}) - R_{t}^{\leftarrow}(\mathbf{y}) \right)$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\mathbf{y}) \cdot \left(\sum_{\mathbf{y}' \neq \mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \cdot \ln \frac{R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y})}{\tilde{R}_{t}(\mathbf{y}', \mathbf{y})} + \sum_{\mathbf{y}' \neq \mathbf{y}} \tilde{R}_{t}(\mathbf{y}', \mathbf{y}) - \sum_{\mathbf{y}' \neq \mathbf{y}} R_{t}^{\leftarrow}(\mathbf{y}', \mathbf{y}) \right)$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') \cdot \left[-\frac{q_{t}^{\leftarrow}(\mathbf{y}')}{q_{t}^{\leftarrow}(\mathbf{y})} + \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') + \frac{q_{t}^{\leftarrow}(\mathbf{y}')}{q_{t}^{\leftarrow}(\mathbf{y})} \ln \frac{q_{t}^{\leftarrow}(\mathbf{y}')}{q_{t}^{\leftarrow}(\mathbf{y})\tilde{v}_{t,\mathbf{y}}(\mathbf{y}')} \right]$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}} q_{t}^{\leftarrow}(\mathbf{y}) \cdot \sum_{\mathbf{y}' \neq \mathbf{y}} R^{\rightarrow}(\mathbf{y}, \mathbf{y}') D_{\phi} \left(\frac{q_{t}^{\leftarrow}(\mathbf{y}')}{q_{t}^{\leftarrow}(\mathbf{y})} \| \tilde{v}_{t,\mathbf{y}}(\mathbf{y}') \right),$$

$$(66)$$

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. (59) can be rewritten as

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

With a recursive manner, we have

$$\operatorname{KL}\left(q_{T-\delta}^{\leftarrow} \left\| \hat{q}_{T-\delta} \right\} \leq \operatorname{KL}\left(q_{0}^{\leftarrow} \left\| \hat{q}_{0} \right\} + L_{\operatorname{SE}}(\tilde{v}) = \operatorname{KL}\left(q_{T}^{\rightarrow} \left\| \hat{q}_{0} \right\} + L_{\operatorname{SE}}(\tilde{v}) \leq (1 + e^{-T})^{d} - 1 + T\epsilon_{\operatorname{score}}^{2},$$

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

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

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

Bounding $\mathrm{TV}(q_*,q_\delta^{\rightarrow})$ 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

$$\operatorname{TV}(q_*, q_{\overrightarrow{\delta}}) \coloneqq \inf_{\gamma \in \Gamma(q_*, q_{\overrightarrow{\delta}})} \mathbb{P}_{(u, v) \sim \gamma}[u \neq v] \le \mathbb{P}(\mathbf{y} \neq \mathbf{y}'),$$

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

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

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

$$\delta \le d^{-1}\epsilon \le d^{-1} \cdot \ln\left(\frac{1}{1-\epsilon}\right) \quad \Rightarrow \quad e^{\delta d} \le \frac{1}{1-\epsilon} \quad \Rightarrow \quad \text{TV}\left(q_*, q_{\delta}^{\rightarrow}\right) \le 1 - e^{-\delta d} \le \epsilon. \tag{67}$$

Finally, we have

$$\operatorname{TV}\left(q_{0}^{\rightarrow},\hat{q}_{T-\delta}^{\leftarrow}\right) \leq \operatorname{TV}\left(q_{0}^{\rightarrow},q_{\delta}^{\rightarrow}\right) + \operatorname{TV}\left(q_{T-\delta}^{\leftarrow},\hat{q}_{T-\delta}\right) \leq \epsilon + \sqrt{\frac{1}{2}\operatorname{KL}\left(q_{T-\delta}^{\leftarrow}\left\|\hat{q}_{T-\delta}\right\right)} \leq 2\epsilon.$$

Hence the proof is completed.

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 $w \in \{1, 2, ..., W\}$, $W = (T - \delta)/\eta$, $\eta = \epsilon/2d$, 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 $\mathrm{TV}\left(p_*,\hat{p}\right) \leq 2\epsilon$ 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

$$\mathbb{E}\left[\sum_{w=1}^{W} \beta_{t_w}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_w - t_{w-1})\right] = \sum_{w=1}^{W} \mathbb{E}[\text{numK}(\hat{\mathbf{y}}_{t_{w-1}})] \cdot \frac{K}{e^{(T - t_w)} - 1} \cdot (t_w - t_{w-1}). \quad (68)$$

Although $\mathbb{E}[\operatorname{numK}(\hat{\mathbf{y}}_{t_{w-1}})]$ 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.,

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

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

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

$$+\underbrace{\left(1-\mathbf{1}_{(K,K)}(\boldsymbol{y}_i,\boldsymbol{y}_i')\right)\cdot\mathbf{1}_K(\boldsymbol{y}_i)\cdot(1-e^{-t})}_{\text{turn into mask token}}\right].$$

due to Lemma 8. Let $\mathbb{P}[\operatorname{numK}(\mathbf{y}_t^{\rightarrow}) = k]$ 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

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

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

$$\operatorname{TV}\left(q_{T-t_{w-1}}^{\leftarrow}, \hat{q}_{T-t_{w-1}}\right) \leq \operatorname{TV}\left(q_{T-t_{W}}^{\leftarrow}, \hat{q}_{T-t_{W}}\right) \leq 2\epsilon,$$

which implies

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

Then, we have

$$\mathbb{E}[\operatorname{numK}(\hat{\mathbf{y}}_{t_{w-1}})] \leq \mathbb{E}[\operatorname{numK}(\mathbf{y}_{t_{w-1}}^{\leftarrow})] + 2d\epsilon = \mathbb{E}[\operatorname{numK}(\mathbf{y}_{T-t_{w-1}}^{\rightarrow})] + 2d\epsilon$$

$$= d \cdot (1 - e^{-(T-t_{w-1})}) + 2d\epsilon.$$
(69)

Plugging Eq. (69) into Eq. (68), we have

$$\mathbb{E}\left[\sum_{w=1}^{W} \beta_{t_{w}}(\hat{\mathbf{y}}_{t_{w-1}}) \cdot (t_{w} - t_{w-1})\right]$$

$$\leq \sum_{w=1}^{W} d \cdot (1 - e^{-(T - t_{w-1})}) \frac{K}{e^{(T - t_{w})} - 1} \cdot (t_{w} - t_{w-1})$$

$$+ \sum_{w=1}^{W} 2d\epsilon \cdot \frac{K}{e^{(T - t_{w})} - 1} \cdot (t_{w} - t_{w-1})$$

$$= Kd \cdot \sum_{w=1}^{W} e^{-(T - t_{w})} \cdot (t_{w} - t_{w-1}) \cdot \frac{1 - e^{-(T - t_{w-1})}}{1 - e^{-(T - t_{w})}}$$

$$+ 2Kd\epsilon \cdot \sum_{w=1}^{W} \left(e^{T - t_{w}} - 1\right)^{-1} \cdot (t_{w} - t_{w-1})$$
(70)

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

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

Under these conditions, we have

$$\eta \leq \frac{\delta}{2} \leq \ln(\frac{1}{2} + \frac{e^{\delta}}{2}) \quad \Rightarrow \quad e^{\eta} \leq \frac{e^{\delta}}{2} + \frac{1}{2} \leq \frac{e^{(T - t_{w-1})}}{2} + \frac{1}{2} \quad \forall \ w \in \{1, \dots, W\}
\Rightarrow \quad e^{\eta} \leq \frac{1 + e^{-(T - t_{w-1})}}{2e^{-(T - t_{w-1})}} \quad \Rightarrow \quad 2 \cdot e^{-(T - t_{w-1} - \eta)} \leq 1 + e^{-(T - t_{w-1})}
\Rightarrow \quad 1 - e^{-(T - t_{w-1})} \leq 2 - 2e^{-(T - t_{w-1} - \eta)} \quad \Rightarrow \quad \frac{1 - e^{-(T - t_{w-1})}}{1 - e^{-(T - t_{w})}} \leq 2.$$
(71)

Plugging these results into Term 1 of Eq. (70), we have

Term
$$1 = 2Kd \cdot \sum_{w=1}^{W} e^{-(T-t_w)} \cdot \eta = 2Kd \cdot \sum_{w=1}^{W} e^{-(T-w\eta)} \cdot \eta$$

 $= 2Kd \cdot \eta \cdot e^{-T} \cdot \frac{e^{(W+1)\eta} - e^{\eta}}{e^{\eta} - 1} \le 2Kd \cdot e^{\eta} \cdot (e^{-\delta} - e^{-T}) \le 2Kd \cdot (1 - e^{-T})$ (72)
 $= 2Kd \cdot \left(1 - \frac{\epsilon^2}{4d}\right)$.

Moreover, we have

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

where the first inequality follows from Eq. (71) and the last inequality is established when $\eta \leq 1$. Then, Term 2 of Eq. (70) can be upper bounded as

$$\operatorname{Term} 2 = 2Kd\epsilon \cdot \sum_{w=1}^{W} \frac{\eta}{e^{T - t_w} - 1} \le 4e \cdot Kd\epsilon \cdot \sum_{w=1}^{W} \frac{\eta}{e^{T - t_{w-1}} - 1} \le 4e \cdot Kd\epsilon \cdot \sum_{w=1}^{W} \frac{\eta}{T - t_{w-1}}$$

$$\le 4e \cdot dK\epsilon \cdot \int_{0}^{T - \delta} \frac{1}{T - t} dt = 4e \cdot dK\epsilon \cdot \ln \frac{T}{\delta} \le 4e \cdot dK\epsilon \cdot \ln \frac{4d^2}{\epsilon^3} \le 12e \cdot Kd \ln d \cdot \epsilon \ln \frac{1}{\epsilon}$$

$$(73)$$

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

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

which follows from the monotonicity of the function $x \ln x$. Under this condition, the RHS of Eq. (73) has the following bound

Term
$$2 \le 12 \cdot Kd \ln d$$
. (74)

Finally, plugging Eq. (72) and Eq. (74) into Eq. (70), the expected calls of discrete scores will be

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

Hence, the proof is completed.

E TECHNICAL LEMMAS

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

$$\overline{A} \coloneqq A_1 \otimes A_2 \otimes \ldots \otimes A_n$$

then we have

$$\overline{\boldsymbol{A}}_{[a_{1,i},a_{2,i},...,a_{n,i}],[a_{1,j},a_{2,j},...,a_{n,j}]} \coloneqq \overline{\boldsymbol{A}}_{\sum_{k=1}^{n} a_{k,i} \cdot d^{n-k}, \sum_{k=1}^{n} a_{k,j} \cdot d^{n-k}} = \prod_{k=1}^{n} [\boldsymbol{A}_{k}]_{a_{k,i},a_{k,j}}.$$

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

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 starts 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,

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

Similarly,

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

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

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

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

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

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. \Box

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 $p(\boldsymbol{y},t|\boldsymbol{x},s)$, then it solves the backward Kolmogorov equation

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

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 $(\mathbf{x}, \tilde{\mathbf{z}})$ and $(\tilde{\mathbf{x}}, \tilde{\mathbf{z}})$, which we write in terms of the conditionals and marginals as

$$p_{x,z}(\boldsymbol{x},\boldsymbol{z}) = p_{x|z}(\boldsymbol{x}|\boldsymbol{z}) \cdot p_z(\boldsymbol{z}) = p_{z|x}(\boldsymbol{z}|\boldsymbol{x}) \cdot p_x(\boldsymbol{x})$$
$$q_{x,z}(\boldsymbol{x},\boldsymbol{z}) = q_{x|z}(\boldsymbol{x}|\boldsymbol{z}) \cdot q_z(\boldsymbol{z}) = q_{z|x}(\boldsymbol{z}|\boldsymbol{x}) \cdot q_x(\boldsymbol{x}).$$

then we have

$$KL\left(p_{x,z} \| q_{x,z}\right) = KL\left(p_z \| q_z\right) + \mathbb{E}_{\mathbf{z} \sim p_z} \left[KL\left(p_{x|z}(\cdot | \mathbf{z}) \| q_{x|z}(\cdot | \mathbf{z})\right)\right]$$
$$= KL\left(p_x \| q_x\right) + \mathbb{E}_{\mathbf{x} \sim p_x} \left[KL\left(p_{z|x}(\cdot | \mathbf{x}) \| q_{z|x}(\cdot | \mathbf{x})\right)\right]$$

where the latter equation implies

$$\mathrm{KL}\left(p_{x} \| q_{x}\right) \leq \mathrm{KL}\left(p_{x,z} \| q_{x,z}\right).$$

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

1944 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 $p_{x,z}(\boldsymbol{x}, \boldsymbol{z}) = p_{x|z}(\boldsymbol{x}|\boldsymbol{z}) \cdot p_z(\boldsymbol{z}) = p_{z|x}(\boldsymbol{z}|\boldsymbol{x}) \cdot p_x(\boldsymbol{x})$ 1947 $q_{x,z}(\boldsymbol{x}, \boldsymbol{z}) = q_{x|z}(\boldsymbol{x}|\boldsymbol{z}) \cdot q_z(\boldsymbol{z}) = q_{z|x}(\boldsymbol{z}|\boldsymbol{x}) \cdot q_x(\boldsymbol{x}).$

then we have

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

Besides, we have

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

F THE USE OF LARGE LANGUAGE MODELS (LLMS)

In writing, we used LLMs for grammar checking and sentence polishing.