
Fast Solvers for Discrete Diffusion Models: Theory and Applications of High-Order Algorithms

Yinuo Ren, Haoxuan Chen, Yuchen Zhu, Wei Guo, Yongxin Chen, Grant M. Rotskoff, Molei Tao, and Lexing Ying

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Discrete diffusion models have emerged as a powerful generative modeling frame-
2 work for discrete data with successful applications spanning from text generation
3 to image synthesis. However, their deployment faces challenges due to the high
4 dimensionality of the state space, necessitating the development of efficient in-
5 ference algorithms. Current inference approaches mainly fall into two categories:
6 exact simulation and approximate methods such as τ -leaping. While exact meth-
7 ods suffer from unpredictable inference time and redundant function evaluations,
8 τ -leaping is limited by its first-order accuracy. In this work, we advance the latter
9 category by tailoring the first extension of high-order numerical inference schemes
10 to discrete diffusion models, enabling larger step sizes while reducing error. We
11 rigorously analyze the proposed schemes and establish the second-order accuracy
12 of the θ -trapezoidal method in KL divergence. Empirical evaluations on GPT-2
13 level text and ImageNet-level image generation tasks demonstrate that our method
14 achieves superior sample quality compared to existing approaches under equivalent
15 computational constraints.

16 1 Introduction

17 Diffusion and flow-based models on discrete spaces [1–10] have emerged as a cornerstone of modern
18 generative modeling for categorical data, offering unique advantages in domains where continuity
19 assumptions fail. Unlike their continuous counterparts, discrete diffusion models inherently accom-
20 modate data with discrete structures, *e.g.*, language tokens, molecular sequences, tokenized images,
21 and graphs, enabling principled generation and inference in combinatorially complex spaces. These
22 models have exerted a large impact on numerous applications, from the design of molecules [11],
23 proteins [12], and DNA sequences [13, 14] under biophysical constraints, to the generation of high-
24 fidelity text [15] and images [16] via autoregressive or masked transitions, *etc.*. Beyond standalone
25 tasks, discrete diffusion models also synergize with methodologies, ranging from tensor networks [17]
26 to guidance mechanisms [18–20].

27 Discrete diffusion models, despite their broad applicability, face a critical bottleneck: *inference*
28 *inefficiency*. Current inference methods include: (1) exact simulation methods [21], which ensure
29 unbiased sampling from the pre-trained model but suffer from unpredictable inference time and
30 redundant score evaluations, resulting in poor scaling w.r.t. dimensionality; and (2) approximate
31 methods such as τ -leaping [22], which offer simple and parallelizable implementation but, due to
32 their first-order accuracy, requires small step sizes to control discretization error, forcing a stringent
33 trade-off between speed and sample quality.

34 To address these limitations in possibly computationally constrained environments, we develop
35 high-order numerical schemes tailored for discrete diffusion model inference. Drawing inspirations
36 from acceleration techniques developed for ordinary differential equations (ODEs) [23], stochastic

differential equations (SDEs) [24, 25], chemical reaction simulations [26], and most recently continuous diffusion [27–29], our work represents the *first successful adaptation of high-order numerical schemes to the discrete diffusion domain*. Through careful design, these high-order schemes provide unprecedented efficient and versatile solutions for discrete diffusion model inference.

Our Contributions. The main contributions of this paper are summarized as follows:

- We introduce the *first high-order numerical solvers* for discrete diffusion model inference, namely the θ -Runge-Kutta-2 (θ -RK-2) method and the θ -trapezoidal method;
- We rigorously establish the theoretical properties of both methods, proving *second-order convergence* of the θ -trapezoidal method and conditional second-order convergence of the θ -RK-2 method;
- We empirically validate our theoretical results and demonstrate the *superior performance* of the θ -trapezoidal method through comprehensive evaluations on large-scale text and image generation benchmarks.

1.1 Related Works

Here we briefly review related works and defer a more detailed discussion to App. A.

Discrete Diffusion Models. Since their introduction, discrete diffusion models have undergone significant refinements, including the development of score-entropy loss [30] and flow-matching formulation [31, 32]. These models generally fall into two categories based on their noise distribution: uniform [30, 20] and masked (absorbing state) [33–35, 21], each offering unique advantages in modeling discrete distributions. Recent theoretical advances have emerged through studies [36–38].

High-Order Scheme for Continuous Diffusion Models. The development of high-order numerical schemes for solving ODEs and SDEs represents decades of research, as comprehensively reviewed in [23, 39, 40]. These schemes have recently been adapted to accelerate continuous diffusion model inference, encompassing approaches such as the exponential integrators [41–43], Adams-Bashforth methods [29, 44, 45], Taylor methods [27, 46] and (stochastic) Runge-Kutta methods [47, 28, 48–51].

High-Order Scheme for Chemical Reaction Systems. Regarding approximate methods for simulating compound Poisson processes and chemical reaction systems with state-dependent intensities, efforts have been made on the τ -leaping method [52], and its extensions [53, 54, 26, 55]. For a quick review of the problem setting and these methods, one may refer to [56, 57]. The adaption of these methods to discrete diffusion models presents unique challenges due to the presence of both time and state-inhomogeneous intensities in the underlying Poisson processes.

2 Preliminaries

In this subsection, we review several basic concepts and previous error analysis results of discrete diffusion models.

2.1 Discrete Diffusion Models

In discrete diffusion models, one considers a continuous-time Markov chain (CTMC) $(x_t)_{0 \leq t \leq T}$ on a finite space \mathbb{X} as the *forward process*. We represent the distribution of x_t by a vector $p_t \in \Delta^{|\mathbb{X}|}$, where $\Delta^{|\mathbb{X}|}$ denotes the probability simplex in $\mathbb{R}^{|\mathbb{X}|}$. Given a target distribution p_0 , the CTMC satisfies the following equation:

$$\frac{dp_t}{dt} = Q_t p_t, \quad \text{where } Q_t = (Q_t(y, x))_{x, y \in \mathbb{X}} \quad (2.1)$$

is the rate matrix at time t satisfying

$$(i) \ Q_t(x, x) = - \sum_{y \neq x} Q_t(y, x), \ \forall x \in \mathbb{X}; \ (ii) \ Q_t(x, y) \geq 0, \ \forall x \neq y \in \mathbb{X}.$$

Below, we will use the notation $Q_t^0 = Q_t - \text{diag } Q_t$. It can be shown that the corresponding backward process is of the same form but with a different rate matrix [58]:

$$\frac{d\bar{p}_s}{ds} = \bar{Q}_s \bar{p}_s, \quad \text{where } \bar{Q}_s(y, x) = \begin{cases} \frac{\bar{p}_s(y)}{\bar{p}_s(x)} \bar{Q}_s(x, y), & \forall x \neq y \in \mathbb{X}, \\ - \sum_{y' \neq x} \bar{Q}_s(y', x), & \forall x = y \in \mathbb{X}. \end{cases} \quad (2.2)$$

is the rate matrix and $\tilde{*}_s$ denotes $*_{T-s}$. The rate matrix Q_t is often chosen to possess certain sparse structures such that the forward process converges to a simple distribution that is easy to sample from. Popular choices include the uniform and absorbing state cases [30], where the forward process (2.1) converges to the uniform distribution on \mathbb{X} and a Dirac distribution, respectively.

Common training practice is to define the score function (or the score vector) as $s_t(x) = (s_t(x, y))_{y \in \mathbb{X}} := \frac{p_t}{p_t(x)}$ for any $x \in \mathbb{X}$, $t \in [0, T]$ and estimate it by a neural network $\hat{s}_t^\phi(x)$, where the parameters ϕ are trained by minimizing the score entropy [30, 59] for some weights $\psi_t \geq 0$ as:

$$\min_{\phi} \int_0^T \psi_t \mathbb{E}_{x_t \sim p_t} \left[\sum_{y \neq x_t} Q_t(x_t, y) \left(s_t(x_t, y) \log \frac{s_t(x_t, y)}{\hat{s}_t^\phi(x_t, y)} - s_t(x_t, y) + \hat{s}_t^\phi(x_t, y) \right) \right] dt. \quad (2.3)$$

Similar to the continuous case, the backward process is approximated by another CTMC $\frac{dq_s}{ds} = \hat{Q}_s^\phi q_s$, with $q_0 = p_\infty$ and rate matrix \hat{Q}_s^ϕ , where $\hat{Q}_s^\phi(y, x) = \tilde{s}_s^\phi(x, y) \tilde{Q}_s(x, y)$ for any $x \neq y \in \mathbb{X}$. The inference is done by first sampling from p_∞ and then evolving the CTMC accordingly. For simplicity, we drop the superscript ϕ hereafter.

2.2 Stochastic Integral Formulation of Discrete Diffusion Models

Discrete diffusion models can also be formulated as stochastic integrals, which is especially useful for their theoretical analysis [38]. In this section, we briefly recapitulate relevant results therein and refer to App. B for mathematical details. Below we work on the probability space $(\Omega, \mathcal{B}, \mathbb{P})$ and denote the pairwise difference set of the state space \mathbb{X} by $\mathbb{D} := \{x - y : x \neq y \in \mathbb{X}\}$. In this work, we focus on the case where $\mathbb{X} = [S]^d$ with d data dimensions and S sites along each dimension.

We first introduce the Poisson random measure, a key concept in the formulation.

Definition 2.1 (Informal Definition of Poisson Random Measure). *The random measure $N[\lambda](dt, d\nu)$ on $\mathbb{R}^+ \times \mathbb{D}$ is called a Poisson random measure with evolving intensity λ w.r.t. a measure γ on \mathbb{D} if, roughly speaking, the number of jumps of magnitude ν during the infinitesimal time interval $(t, t + dt]$ is Poisson distributed with mean $\lambda_t(\nu) \gamma(d\nu) dt$.*

The forward process (2.1) can thus be represented by the following stochastic integral:

$$x_t = x_0 + \int_0^t \int_{\mathbb{D}} \nu N[\lambda](ds, d\nu),$$

where the intensity λ is defined as $\lambda_t(\nu, \omega) = Q_t^0(x_{t-}(\omega) + \nu, x_{t-}(\omega))$ if $x_{t-}(\omega) + \nu \in \mathbb{X}$ and 0 otherwise. Here, the outcome $\omega \in \Omega$ and x_{t-} denotes the left limit of the càdlàg process x_t at time t with $x_{0-} = x_0$. We will also omit the variable ω , should it be clear from context. The backward process in discrete diffusion models (2.2) can also be represented similarly as:

$$y_s = y_0 + \int_0^s \int_{\mathbb{D}} \nu N[\mu](ds, d\nu), \quad (2.4)$$

where the intensity μ is defined as $\mu_s(\nu, \omega) = \tilde{s}_s(y_{s-}, y_{s-} + \nu) \tilde{Q}_s^0(y_{s-}, y_{s-} + \nu)$ if $y_{s-} + \nu \in \mathbb{X}$ and 0 otherwise. During inference, $\hat{y}_s = \hat{y}_0 + \int_0^s \int_{\mathbb{D}} \nu N[\hat{\mu}](ds, d\nu)$ is used instead of (2.4), where the estimated intensity $\hat{\mu}$ is defined by replacing the true score s_t with the neural network estimated score \hat{s}_t in $\mu_s(\nu, \omega)$. In the following, we also denote the intensity $\mu_s(\nu, \omega)$ at time s by $\mu_s(\nu, y_{s-})$ with slight abuse of terminology to emphasize its dependency on ω through $y_{s-}(\omega)$.

3 Numerical Schemes for Discrete Diffusion Model Inference

Before introducing the proposed numerical schemes, we first review existing numerical schemes for discrete diffusion models, including exact simulation methods and the τ -leaping method, and discuss their merits and limitations.

3.1 Exact Simulation Methods

Unlike continuous diffusion models, where exact simulation is beyond reach, discrete diffusion models permit inference without discretization error. Notable examples of unbiased samplers include

uniformization [36] for the uniform state case and the First-Hitting Sampler (FHS) [21] for the absorbing state case. The main idea behind these methods is to first sample the next jump time and then the jump itself. Theoretical analysis [38] reveals that such schemes *lack guarantees with finite computation budget*, since the number of required jumps (and thus the inference time) follows a random distribution with expectation $\Omega(d)$. This computational restriction may be less favorable for high-dimensional applications, such as generative modeling of DNA or protein sequences.

Furthermore, *the absence of discretization error does not necessarily translate to superior sample quality*, given the inherent estimation errors in neural network-based score functions. This limitation is further amplified by the *highly skewed distribution* of jumps, with a concentration occurring during the terminal phase of the backward process, when the neural network-based score function exhibits the highest estimation error. This phenomenon stems from the potential singularity of the target distribution p_0 , which induces singularities in the score function, making accurate neural network estimation particularly challenging during that phase (cf. Assump. 4.4 [38]).

Fig. 1 illustrates an application of the uniformization algorithm to discrete diffusion inference for text generation, with detailed experimental parameters presented in Sec. 6.3 and App. D.3. As the process approaches the target distribution ($t \rightarrow T$), the number of jumps (in terms of the number of score function evaluations, NFE) grows unbounded, while perplexity improvements become negligible. This skewness of computational effort results in *redundant function evaluations*. Although early stopping is commonly adopted at $T - \delta$ for some small $\delta \ll 1$ to alleviate this inefficiency, this approach introduces challenges in its selection, particularly under computational constraints or when efficiency-accuracy trade-offs are desired. Moreover, the variable jump schedules across batch samples complicate parallelization efforts in exact methods, highlighting the need for more adaptable and efficient algorithmic solutions.

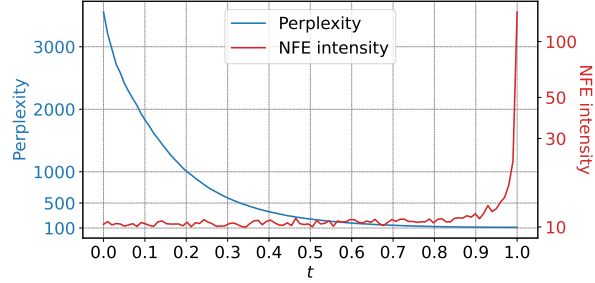


Figure 1: An illustrative application of the uniformization algorithm to discrete diffusion models for text generation. The x -axis denotes the time of the backward process, and the y -axis denotes the frequency of jumps (NFE). Perplexity convergence occurs before the NFE grows unbounded.

3.2 Approximate Method: τ -Leaping Method

The τ -leaping method [52, 22] represents a widely adopted scheme that effectively addresses both dimensionality scaling and inference time control challenges. This Euler-type scheme approximates the backward process with time-dependent intensity $\hat{\mu}_t$ via the following updates:

$$\hat{y}_{t+\Delta} = \hat{y}_t + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}(\hat{\mu}_t(\nu)\Delta), \quad (3.1)$$

where Δ denotes the time step and $\mathcal{P}(\cdot)$ denotes a Poisson random variable. In general, one may design different discretization schemes for τ -leaping, and the summation in (3.1) is parallelizable, underscoring the method's flexibility and efficiency. We refer to Alg. 3 and App. B.2 for a detailed description of the τ -leaping method for discrete diffusion model inference. Regarding convergence properties as the time discretization becomes increasingly refined, theoretical analyses by [22, 38] have established the error bounds of the τ -leaping method, the results of which are summarized in the following theorem. Further discussion can be found in App. B.2.

Theorem 3.1 (Thm. 4.7 in [38]). *Under a certain discretization scheme and technical assumptions, and given an ϵ -accurate score function, the following error bound holds:*

$$D_{\text{KL}}(p_\delta \| \hat{q}_{T-\delta}) \lesssim \exp(-T) + \epsilon + \kappa T, \quad (3.2)$$

where $\delta \ll 1$ is the early stopping time, κ controls the step size, and T is the time horizon. The notation \lesssim indicates the inequality holds up to a constant factor as $\kappa \rightarrow 0$.

The error bound (3.2) decouples three error sources of the τ -leaping scheme: the truncation error $\mathcal{O}(e^{-T})$, the score estimation error ϵ , and the discretization error $\mathcal{O}(\kappa T)$. Similar to the case for

the Euler method for ODEs and the Euler-Maruyama scheme for SDEs, the τ -leaping method is a first-order scheme in terms of the discretization error $\mathcal{O}(\kappa T)$.

4 Algorithms: High-Order Inference Schemes

A natural improvement of τ -leaping is to develop high-order schemes for discrete diffusion models. As a foundational example, consider the second-order Runge-Kutta (RK-2) method with two stages [23] for solving the ODE $dx_t = f_t(x_t)dt$. This method represents one of the simplest high-order numerical schemes:

$$\hat{x}_{t+\theta\Delta}^* = \hat{x}_t + f_t(\hat{x}_t)\theta\Delta, \quad \hat{x}_{t+\Delta} = \hat{x}_t + \left[(1 - \frac{1}{2\theta})f_t(\hat{x}_t) + \frac{1}{2\theta}f_{t+\theta\Delta}(\hat{x}_{t+\theta\Delta}^*)\right]\Delta. \quad (4.1)$$

This scheme reduces to the exact midpoint method when $\theta = \frac{1}{2}$ and Heun's method when $\theta = 1$. The underlying intuition stems from the observation that for $f \in C^2(\mathbb{R})$, $[(1 - \frac{1}{2\theta})f(0) + \frac{1}{2\theta}f(\theta\Delta)]\Delta$ offers a second-order approximation of $\int_0^\Delta f(x)dx$ in contrast to $f(0)\Delta$, which is only first-order. This approach has been successfully adapted for SDE simulation [24] and continuous diffusion model inference [48, 28, 29, 49, 51]. Notably, these methods enhance sample quality and computational efficiency without requiring additional model training, making the development of high-order schemes for discrete diffusion inference both theoretically appealing and practically viable.

In this section, we propose two different high-order solvers for discrete diffusion model inference. We will primarily focus on two-stage algorithms aiming for second-order accuracy. Specifically, we will introduce the θ -RK-2 method and the θ -Trapezoidal method. Throughout this section, we assume a time discretization scheme $(s_i)_{i \in [0:N]}$ with $0 = s_0 < \dots < s_N = T - \delta$, where δ is the early stopping time and use the shorthand notations $*_+ = \max\{0, *\}$. For any $s \in (s_n, s_{n+1}]$ and $n \in [0 : N - 1]$, we define $\lfloor s \rfloor = s_n$, $\rho_s = (1 - \theta)s_n + \theta s_{n+1}$, $\Delta_n = s_{n+1} - s_n$, and θ -section points as $\rho_n = (1 - \theta)s_n + \theta s_{n+1}$. We choose $\gamma(d\nu)$ to be the counting measure on \mathbb{D} .

4.1 θ -RK-2 Method

We first present the θ -RK-2 method, which is simple in design and serves as a natural analog of the second-order RK method for ODEs (4.1) in terms of time and state-dependent Poisson random measures, as a warm-up for the θ -trapezoidal method. We note that similar methods have been proposed for simulating SDEs driven by Brownian motions or Poisson processes, such as the stochastic [24] and the Poisson [54] RK methods. A summary of this method is given in Alg. 1.

Intuitively, the θ -RK-2 method is a two-stage algorithm that:

- (i) Firstly, it runs τ -leaping with step size $\theta\Delta_n$, obtains an intermediate state $\hat{y}_{\rho_n}^*$ at the θ -section point ρ_n , and evaluates the intensity $\hat{\mu}_{\rho_n}^*$ there;
- (ii) Then another step of τ -leaping for a full step Δ_n is run using a weighted sum of the intensities at the current time point s_n and the θ -section point ρ_n .

Algorithm 1: θ -RK-2 Method

Input: $\hat{y}_0 \sim q_0$, $\theta \in (0, 1]$, $(s_n, \rho_n)_{n \in [0:N-1]}$, $\hat{\mu}$, $\hat{\mu}^*$.

Output: A sample $\hat{y}_{s_N} \sim \hat{q}_{t_N}^{\text{RK}}$.

```

1 for  $n = 0$  to  $N - 1$  do
2    $\hat{y}_{\rho_n}^* \leftarrow \hat{y}_{s_n} + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}(\nu)\theta\Delta_n)$ ;
3    $\hat{y}_{s_{n+1}} \leftarrow \hat{y}_{s_n} + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}\left(\mathbf{1}_{\hat{\mu}_{s_n} > 0} \left[(1 - \frac{1}{2\theta})\hat{\mu}_{s_n} + \frac{1}{2\theta}\hat{\mu}_{\rho_n}^*\right]_+(\nu)\Delta_n\right)$ ;
4 end
```

We emphasize that our method is different from the midpoint method proposed in [52] for simulating chemical reactions, where the Poisson random variable in the first step is replaced by its expected magnitude. Such modification is in light of the lack of continuity and orderliness of the state space.

4.2 θ -Trapezoidal Method

As to be shown theoretically and empirically, the conceptually simple θ -RK-2 method may have limitations in terms of both accuracy and efficiency. To this end, we propose the following θ -trapezoidal method, which is developed based on existing methods proposed for simulating SDEs [25] and chemical reactions [26]. Below, we introduce two parameters that will be used extensively later:

$$\alpha_1 = \frac{1}{2\theta(1-\theta)} \text{ and } \alpha_2 = \frac{(1-\theta)^2 + \theta^2}{2\theta(1-\theta)}, \text{ with } \alpha_1 - \alpha_2 = 1.$$

216 The θ -trapezoidal method is sum-
 217 marized in Alg. 2. Intuitively,
 218 this method separates each inter-
 219 val $(s_n, s_{n+1}]$ into two sub-intervals
 220 $(s_n, \rho_n]$ and $(\rho_n, s_{n+1}]$, on which sim-
 221 ulations are detached with different in-
 222 tensities designed in a balanced way.

223 Compared to the θ -RK-2 method,
 224 the θ -trapezoidal method is also two-
 225 stage with an identical first step. The
 226 second step, however, differs in two
 227 major aspects:

Algorithm 2: θ -Trapezoidal Method

Input: $\hat{y}_0 \sim q_0, \theta \in (0, 1], (s_n, \rho_n)_{n \in [0:N-1]}, \hat{\mu}, \hat{\mu}^*$.

Output: A sample $\hat{y}_{s_N} \sim \hat{q}_{t_N}^{\text{trap}}$.

```

1 for  $n = 0$  to  $N - 1$  do
2    $\hat{y}_{\rho_n}^* \leftarrow \hat{y}_{s_n} + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}(\nu) \theta \Delta_n);$ 
3    $\hat{y}_{s_{n+1}} \leftarrow \hat{y}_{\rho_n}^* +$ 
       $\sum_{\nu \in \mathbb{D}} \nu \mathcal{P}\left((\alpha_1 \hat{\mu}_{\rho_n}^* - \alpha_2 \hat{\mu}_{s_n})_+(\nu)(1 - \theta) \Delta_n\right);$ 
4 end
```

- 228 (1) The second step starts from the intermediate state $\hat{y}_{\rho_n}^*$ instead of \hat{y}_{s_n} and only runs for a fractional
 229 step $(1 - \theta)\Delta_n$ rather than a full step Δ_n ;
- 230 (2) The weighted sum is comprised of an altered pair of coefficients $(\alpha_1, -\alpha_2)$, performing an
 231 *extrapolation* instead of interpolation with coefficients $(1 - \frac{1}{2\theta}, \frac{1}{2\theta})$ as in the θ -RK-2 method with
 232 $\theta \in [\frac{1}{2}, 1]$. This feature will be shown to render the algorithm unconditionally second-order.

233 Following the common practice in the litera-
 234 ture [22], we reject updates with multiple jumps
 235 along one dimension in both algorithms, ensur-
 236 ing their well-posedness. A simple analysis
 237 shows that rejection only happens with proba-
 238 bility $\mathcal{O}(\kappa)$, and we refer to further details
 239 in Rmk. C.4. We refer to Props. C.2 and C.3 for
 240 the stochastic integral formulations of these two
 241 algorithms. We provide a visual comparison be-
 242 tween the θ -RK-2 method and the θ -trapezoidal
 243 method in Fig. 2.

244 5 Theoretical Analysis

245 In this section, we provide the theoretical results
 246 of the θ -trapezoidal and θ -RK-2 methods. The goal of this section is to show that under certain
 247 conditions, both methods are second-order accurate, improving from the first-order accuracy of the
 248 τ -leaping method (cf. Thm. 3.1). Our theoretical analysis also reveals that the θ -trapezoidal method
 249 is more robust to the choice of θ than θ -RK-2, to be confirmed by our empirical results in Sec. 6.

250 5.1 Assumptions

251 For simplicity, we impose a periodic boundary condition on the state space $\mathbb{X} = [S]^d$, i.e., embed the
 252 state space in the d -dimensional torus \mathbb{T}^d , to streamline the proofs (cf. Rmk. C.4).

253 **Assumption 5.1** (Convergence of Forward Process). *The forward process converges to the stationary*
 254 *distribution exponentially fast, i.e., $D_{\text{KL}}(p_T \| p_\infty) \lesssim \exp(-T)$.*

255 This assumption ensures rapid convergence of the forward process, controlling error when terminated
 256 at a sufficiently large time horizon T , and is automatically satisfied in the masked state case and the
 257 uniform state case, given sufficient connectivity of the graph (cf. [38]). The exponential rate aligns
 258 with continuous diffusion models (cf. [60]).

259 **Assumption 5.2** (Regularity of Intensity). *For the true intensity $\mu_s(\nu, y_{s-})$ and the estimated intensity*
 260 *$\hat{\mu}_s(\nu, y_{s-})$, it holds almost everywhere w.r.t. $\mu_s(\nu, y_{s-})\gamma(d\nu)\bar{p}_{s-}(dy_{s-})$ that: (1) Both intensities*
 261 *belong to $C^2([0, T - \delta])$; (2) Both intensities are upper and lower bounded on $[0, T - \delta]$.*

262 This assumes two key requirements of the scores: (1) the forward process maintains sufficient
 263 smoothness, which is achievable through appropriate time reparametrization; and (2) if and only if a
 264 state $y \in \mathbb{X}$ is achievable by the forward process and ν is a permissible jump therefrom, then both its
 265 true and estimated intensity are bounded, corresponding to Assumps. 4.3(i), 4.4, and 4.5 [38].

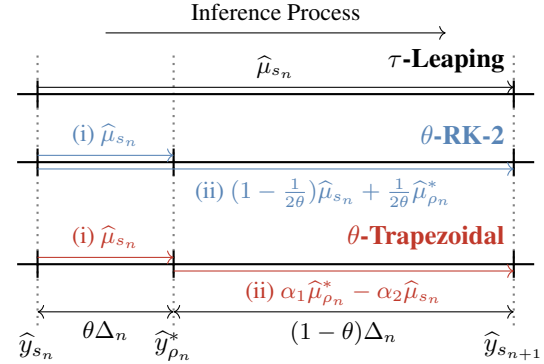


Figure 2: Comparison between τ -leaping method and our proposed second-order schemes.

266 **Assumption 5.3** (Estimation Error). *For all grid points and θ -section points, the estimation error of the neural network-based score is small, i.e., for any $s \in \cup_{n \in [0:N-1]} \{s_n, \rho_n\}$, we have*
 267
 268 (1) $\mathbb{E} \left[\int_{\mathbb{D}} \left(\mu_s(\nu) \left(\log \frac{\mu_s(\nu)}{\hat{\mu}_s(\nu)} - 1 \right) + \hat{\mu}_s(\nu) \right) \gamma(d\nu) \right] \leq \epsilon_I$; (2) $\mathbb{E} \left[\int_{\mathbb{D}} |\mu_s(\nu) - \hat{\mu}_s(\nu)| \gamma(d\nu) \right] \leq \epsilon_{II}$.

269 This assumption quantifies the proximity of the estimated intensity $\hat{\mu}$ to the true intensity μ after
 270 sufficient training. Compared with [38], the additional L^∞ part in (2) is required for technical reasons,
 271 which is similar to [61, 51]. In practice, such additional assumptions may be realized by adding extra
 272 penalty terms to the objective function during training.

273 5.2 Convergence Guarantees

274 The following theorem summarizes our theoretical guarantees for the θ -trapezoidal method:

275 **Theorem 5.4** (Second Order Convergence of θ -Trapezoidal Method). *Suppose $\theta \in (0, 1]$ and*
 276 $\alpha_1 \hat{\mu}_{\rho_s}^* - \alpha_2 \hat{\mu}_{[s]} \geq 0$ *for all $s \in [0, T - \delta]$, then the following error bound holds for Alg. 2*
 277 *under Assumps. 5.1 to 5.3:*

$$D_{\text{KL}}(p_\delta \| \hat{q}_{T-\delta}^{\text{trap}}) \lesssim \exp(-T) + (\epsilon_I + \epsilon_{II})T + \kappa^2 T,$$

278 *where δ is the early stopping time, $\kappa = \max_{n \in [0:N-1]} \Delta_n$, i.e., the largest stepsize, and $\hat{q}_{T-\delta}^{\text{trap}}$ is the*
 279 *distribution obtained by Alg. 1 as defined in Prop. C.2.*

280 The complete proof is presented in App. C.2. The outline is to first bound $D_{\text{KL}}(p_\delta \| \hat{q}_{T-\delta}^{\text{trap}})$ by the
 281 KL divergence between the corresponding path measures, as established in Thm. C.5, and then
 282 decompose the integral in the log-likelihood and bound respectively, where the primary technique
 283 used is *Dynkin's formula* (Thm. C.10). With a term-by-term comparison with Thm. 3.1, we observe
 284 a significant improvement in the discretization error term from $\mathcal{O}(\kappa T)$ to $\mathcal{O}(\kappa^2 T)$. This confirms
 285 that the θ -trapezoidal method achieves second-order accuracy given a sufficient time horizon T and
 286 accurate score estimation, with empirical validation presented in Sec. 6.

287 **Theorem 5.5** (Conditional Second-Order Convergence of θ -RK-2 Method). *Suppose $\theta \in (0, \frac{1}{2}]$ and*
 288 $(1 - \frac{1}{2\theta}) \hat{\mu}_{[s]} + \frac{1}{2\theta} \hat{\mu}_{\rho_s}^* \geq 0$ *for all $s \in [0, T - \delta]$, then the following error bound holds for Alg. 1*
 289 *under Assumps. 5.1 to 5.3:*

$$D_{\text{KL}}(p_\delta \| \hat{q}_{T-\delta}^{\text{RK}}) \lesssim \exp(-T) + (\epsilon_I + \epsilon_{II})T + \kappa^2 T,$$

290 *where δ is the early stopping time, $\kappa = \max_{n \in [0:N-1]} \Delta_n$, i.e., the largest stepsize, and $\hat{q}_{T-\delta}^{\text{RK}}$ is the*
 291 *distribution obtained by Alg. 2 as defined in Prop. C.3.*

292 The proof of the theorem above is provided in App. C.3. The restricted range of θ is caused by one
 293 specific error term (III.4) (C.9) that permits bounding with *Jensen's inequality* only when $\theta \in (0, \frac{1}{2}]$,
 294 similar to its counterpart (II.4) (C.11) in the θ -trapezoidal method. The limitation arises partially
 295 because the weighted sum with coefficients $(1 - \frac{1}{2\theta}, \frac{1}{2\theta})$ becomes an *extrapolation* only if $1 - \frac{1}{2\theta} < 0$,
 296 a feature that naturally holds for all $\theta \in (0, 1]$ in the θ -trapezoidal method. These theoretical findings
 297 are consistent with the empirical observations in Fig. 6 of App. D.3, where the performance of θ -RK-2
 298 method clearly peaks when $\theta \in (0, \frac{1}{2}]$.

299 **Remark 5.6** (Comparison between Trapezoidal and RK-2 Methods). *Trapezoidal methods were*
 300 *originally proposed by [25] as a minimal second-order scheme in the weak sense for simulating*
 301 *SDEs. In simulating chemical reaction contexts, [26] claimed that trapezoidal methods also achieve*
 302 *second-order convergence for covariance error apart from the weak error, a property not shared*
 303 *by midpoint (RK-2) methods. Our empirical results partly reflect these findings, while we defer*
 304 *theoretical investigation of covariance error convergence in discrete diffusion models to future work.*

305 **Remark 5.7** (Remark on the Positivity of Extrapolated Intensity). *Due to the extrapolation nature,*
 306 *both our theorems require an additional assumption on the positivity of the extrapolated intensity,*
 307 *which is classically assumed in [25, 26], and the resolution of this issue is a long-standing open*
 308 *problem. The best result so far is Prop. 5 [26], claiming clamping the intensity above 0 only causes an*
 309 *error of order $\mathcal{O}(\kappa^p)$, for any large integer p . We empirically evaluate this assumption in Tab. 3 with*
 310 *the text generation task (Sec. 6.2) and find that positivity occurs for both methods with high probability*
 311 *over 95%, approaching 100% with increasing NFE. We refer to further discussion in Rmk. C.6.*

6 Experiments

Based on the theoretical analysis, we expect the θ -trapezoidal method to outperform the τ -leaping method and the θ -RK-2 method in terms of sample quality, given the same number of function evaluations. This section empirically validates the anticipated effectiveness of our proposed θ -trapezoidal method (Alg. 2) through comprehensive evaluations across text and image generation tasks. Our comparative analysis includes established discrete diffusion samplers as baselines, *e.g.*, the Euler method [33], τ -leaping [22], Tweedie τ -leaping [30], First-Hitting Sampler (FHS) [21], and Parallel Decoding [62]. We conduct evaluations on both uniform and masked discrete diffusion models, with detailed experimental protocols provided in App. D.

6.1 15-State Toy Model

We first evaluate the performance of the θ -trapezoidal method using a 15-state toy model ($d = 1$, $S = 15$). The target distribution is uniformly generated from Δ^{15} , with rate matrix $Q = \frac{1}{15}E - I$, where E is the all-one and I is the identity matrix. This setup provides analytically available score functions, allowing isolation and quantification of numerical errors introduced by inference algorithms. We apply both the θ -trapezoidal and the θ -RK-2 method to generate 10^6 samples and estimate the KL divergence between the true ground truth p_0 and the generated distribution \hat{q}_T .

For a fair comparison, we choose $\theta = \frac{1}{2}$ for both methods, and the results are presented in Fig. 3. While both methods exhibit super-linear convergence as the total number of steps grows, the θ -trapezoidal method outperforms the θ -RK-2 method in terms of both absolute value and convergence rate, while the θ -RK-2 method takes longer to enter the asymptotic regime. Moreover, the fitted line indicates that the θ -trapezoidal method approximately converges quadratically w.r.t. the step count, confirming our theories.

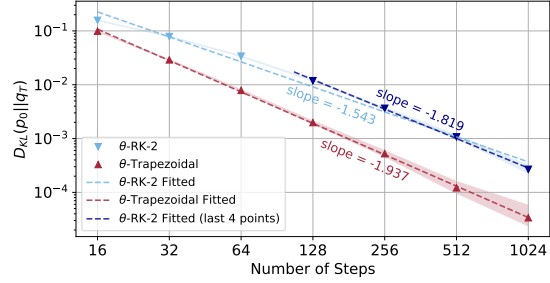


Figure 3: Empirical KL divergence between the true and generated distribution of the toy model vs. number of steps. Data are fitted with linear regression with 95% confidence interval by bootstrapping.

6.2 Text Generation

For the text generation task, we employ the pre-trained score function from RADD [33] as our base model for benchmarking inference algorithms. RADD is a masked discrete diffusion model with GPT-2-level text generation capabilities [63] and is trained on the OpenWebText dataset [64] with $d = 1024$ and $S = 50258$. Our comparative analysis maintains consistent computational resources across methods, quantified through the number of score function evaluations (NFE), and evaluates the sample quality produced by FHS, the Euler method, τ -leaping, Tweedie τ -leaping, and our proposed θ -trapezoidal method. We generate text sequences of 1024 tokens and measure their generative perplexity following the evaluation protocol established in [33].

Table 1: Generative perplexity of texts generated by different sampling algorithms. Lower values are better, with the best in **bold**.

Method	NFE = 128	NFE = 1024
FHS	≤ 122.732	≤ 109.406
Euler	≤ 86.276	≤ 44.686
Tweedie τ -leap.	≤ 85.738	≤ 44.257
τ -leaping	≤ 52.366	≤ 28.797
θ -trapezoidal	$\leq \mathbf{49.051}$	$\leq \mathbf{27.553}$

Tab. 1 presents the results for both low (128) and high (1024) NFE, with comprehensive results across additional NFE values in Tab. 2. The empirical results demonstrate that the θ -trapezoidal method consistently produces better samples under a fixed computation budget compared with existing popular inference algorithms. Notably, it outperforms Euler and Tweedie τ -leaping, two of the best-performing samplers adopted by RADD, by a large margin. It also consistently prevails over FHS, which performs exact simulation at high NFE (1024), supporting again our observations that being free of discretization error does not necessarily imply better sampling quality. These results validate the practical efficiency and accuracy of Alg. 2.

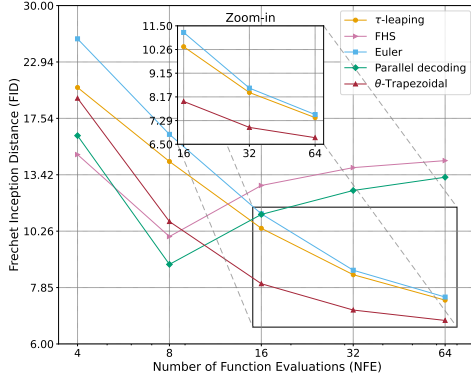


Figure 4: FID of images generated by different sampling algorithms vs. number of function evaluations (NFE). Lower values are better.

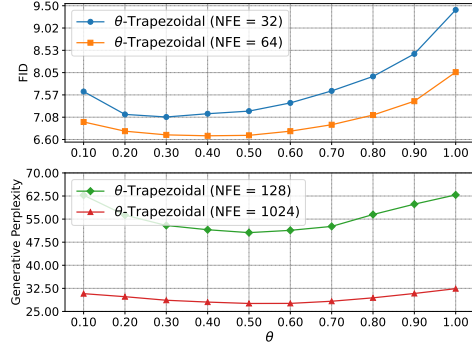


Figure 5: Sampling quality vs. $\theta \in (0, 1]$ in θ -Trapezoidal method. **Upper:** Image generation (FID). **Lower:** Text generation (perplexity). Lower is better.

6.3 Image Generation

Our experiments on image generation utilize the pre-trained score function from MaskGIT [62, 65] as the base model, which can be converted into a masked discrete diffusion model by introducing a noise schedule (see App. D.3). MaskGIT employs a masked image transformer architecture trained on ImageNet [66] of 256×256 resolution, where each image amounts to a sequence of 256 discrete image tokens following VQ-GAN tokenization [67] ($d = 256$, $S = 1025$). We evaluate the θ -trapezoidal method against FHS, the Euler method, τ -leaping, and parallel decoding under equivalent NFE budgets ranging from 4 to 64. Following the setting in [62], we generate 5×10^4 images and compute their Fréchet Inception Distance (FID) against the ImageNet validation split.

Fig. 4 reveals that θ -trapezoidal method (Alg. 2) consistently achieves lower (and thus better) FID values compared to both the Euler method and τ -leaping across all NFE values. While FHS and parallel decoding show advantages at extremely low NFE (≤ 8), their performance saturates with increased computational resources, making them less favorable compared to our rapidly converging method. Additional results, including generated image samples (Fig. 8), are detailed in App. D.

Algorithm Hyperparameters. We evaluate the performance of the θ -trapezoidal method across various θ and NFE values for both text and image generation tasks. As illustrated in Fig. 5, we observe that the θ -trapezoidal method demonstrates robustness to θ , with a flat landscape near the optimal choice. Our empirical analysis suggests that $\theta \in [0.3, 0.5]$ consistently yields competitive performance across different tasks.

7 Conclusion and Future Works

In this work, we introduce the θ -RK-2 and θ -trapezoidal methods as pioneering high-order numerical schemes tailored for discrete diffusion model inference. Through rigorous analysis based on their stochastic integral formulations, we establish second-order convergence of the θ -trapezoidal method and that of the θ -RK-2 method under specified conditions. Our analysis indicates that the θ -trapezoidal method generally provides superior robustness and computational efficiency compared to the θ -RK-2 method. Our empirical evaluations, spanning both a 15-dimensional model with precise score functions and large-scale text and image generation tasks, validate our theoretical findings and demonstrate the superiority performance of our proposed θ -trapezoidal method over existing samplers in terms of sample quality under equivalent computational constraints. Additionally, we provide a comprehensive analysis of the method’s robustness by examining the optimal choice of the parameter θ in our schemes.

Future research directions include comparative analysis of these schemes and development of more sophisticated numerical approaches for discrete diffusion model inference, potentially incorporating adaptive step sizes and parallel sampling methodologies. From the perspective of applications, these methods may also show promise for tasks in computational chemistry and biology, particularly in the design of molecules, proteins, and DNA sequences.

References

- [1] Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202*, 2022.
- [2] 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.
- [3] Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- [4] Griffin Floto, Thorsteinn Jonsson, Mihai Nica, Scott Sanner, and Eric Zhengyu Zhu. Diffusion on the probability simplex. *arXiv preprint arXiv:2309.02530*, 2023.
- [5] 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*, 2022. URL <https://openreview.net/forum?id=Lm8T39vLDTE>.
- [6] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [7] Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35:34532–34545, 2022.
- [8] Pierre H Richemond, Sander Dieleman, and Arnaud Doucet. Categorical sdes with simplex diffusion. *arXiv preprint arXiv:2210.14784*, 2022.
- [9] Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=BYWWwSY2G5s>.
- [10] Javier E Santos, Zachary R Fox, Nicholas Lubbers, and Yen Ting Lin. Blackout diffusion: generative diffusion models in discrete-state spaces. In *International Conference on Machine Learning*, pages 9034–9059. PMLR, 2023.
- [11] Thomas J Kerby and Kevin R Moon. Training-free guidance for discrete diffusion models for molecular generation. *arXiv preprint arXiv:2409.07359*, 2024.
- [12] Nathan C Frey, Daniel Berenberg, Karina Zadorozhny, Joseph Kleinhenz, Julien Lafrance-Vanasse, Isidro Hotzel, Yan Wu, Stephen Ra, Richard Bonneau, Kyunghyun Cho, et al. Protein discovery with discrete walk-jump sampling. *arXiv preprint arXiv:2306.12360*, 2023.
- [13] Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pages 1276–1301. PMLR, 2023.
- [14] Wei Guo, Yuchen Zhu, Molei Tao, and Yongxin Chen. Plug-and-play controllable generation for discrete masked models. *arXiv preprint arXiv:2410.02143*, 2024.
- [15] Do Huu Dat, Do Duc Anh, Anh Tuan Luu, and Wray Buntine. Discrete diffusion language model for long text summarization. *arXiv preprint arXiv:2407.10998*, 2024.
- [16] Minghui Hu, Yujie Wang, Tat-Jen Cham, Jianfei Yang, and Ponnuthurai N Suganthan. Global context with discrete diffusion in vector quantised modelling for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11502–11511, 2022.
- [17] Luke Causer, Grant M Rotskoff, and Juan P Garrahan. Discrete generative diffusion models without stochastic differential equations: a tensor network approach. *arXiv preprint arXiv:2407.11133*, 2024.

- [18] Hunter Nisonoff, Junhao Xiong, Stephan Allenspach, and Jennifer Listgarten. Unlocking guidance for discrete state-space diffusion and flow models. *arXiv preprint arXiv:2406.01572*, 2024.
- [19] Xiner Li, Yulai Zhao, Chenyu Wang, Gabriele Scalia, Gokcen Eraslan, Surag Nair, Tommaso Biancalani, Aviv Regev, Sergey Levine, and Masatoshi Uehara. Derivative-free guidance in continuous and discrete diffusion models with soft value-based decoding. *arXiv preprint arXiv:2408.08252*, 2024.
- [20] Yair Schiff, Subham Sekhar Sahoo, Hao Phung, Guanghan Wang, Sam Boshar, Hugo Dallatorre, Bernardo P de Almeida, Alexander Rush, Thomas Pierrot, and Volodymyr Kuleshov. Simple guidance mechanisms for discrete diffusion models. *arXiv preprint arXiv:2412.10193*, 2024.
- [21] Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. *arXiv preprint arXiv:2409.02908*, 2024.
- [22] 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.
- [23] John Charles Butcher. *The numerical analysis of ordinary differential equations: Runge-Kutta and general linear methods*. Wiley-Interscience, 1987.
- [24] Kevin Burrage and Pamela Marion Burrage. High strong order explicit runge-kutta methods for stochastic ordinary differential equations. *Applied Numerical Mathematics*, 22(1-3):81–101, 1996.
- [25] David F Anderson and Jonathan C Mattingly. A weak trapezoidal method for a class of stochastic differential equations. *Communications in Mathematical Sciences*, 9(1):301–318, 2011.
- [26] Yucheng Hu, Tiejun Li, and Bin Min. A weak second order tau-leaping method for chemical kinetic systems. *The Journal of chemical physics*, 135(2), 2011.
- [27] Hideyuki Tachibana, Mocho Go, Muneyoshi Inahara, Yotaro Katayama, and Yotaro Watanabe. Quasi-taylor samplers for diffusion generative models based on ideal derivatives. *arXiv preprint arXiv:2112.13339*, 2021.
- [28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- [29] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- [30] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*, 2024.
- [31] Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. *arXiv preprint arXiv:2402.04997*, 2024.
- [32] Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky T. Q. Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=GTDKo3Sv9p>.
- [33] 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.

- 497 [34] Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis K Titsias. Simplified and
498 generalized masked diffusion for discrete data. *arXiv preprint arXiv:2406.04329*, 2024.
- 499 [35] Subham Sekhar Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin,
500 Justin T Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked
501 diffusion language models. *arXiv preprint arXiv:2406.07524*, 2024.
- 502 [36] Hongrui Chen and Lexing Ying. Convergence analysis of discrete diffusion model: Exact
503 implementation through uniformization. *arXiv preprint arXiv:2402.08095*, 2024.
- 504 [37] Zikun Zhang, Zixiang Chen, and Quanquan Gu. Convergence of score-based discrete diffusion
505 models: A discrete-time analysis. *arXiv preprint arXiv:2410.02321*, 2024.
- 506 [38] Yinuo Ren, Haoxuan Chen, Grant M Rotskoff, and Lexing Ying. How discrete and continuous
507 diffusion meet: Comprehensive analysis of discrete diffusion models via a stochastic integral
508 framework. *arXiv preprint arXiv:2410.03601*, 2024.
- 509 [39] Peter Eris Kloeden and Eckhard Platen. *Numerical solution of stochastic differential equations*.
510 Stochastic Modelling and Applied Probability, Applications of Mathematics, Springer, 1992.
- 511 [40] Peter Eris Kloeden, Eckhard Platen, and Henri Schurz. *Numerical solution of SDE through*
512 *computer experiments*. Springer Science & Business Media, 2012.
- 513 [41] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential
514 integrator. In *The Eleventh International Conference on Learning Representations*, 2023. URL
515 <https://openreview.net/forum?id=Loek7hfb46P>.
- 516 [42] Qinsheng Zhang, Molei Tao, and Yongxin Chen. gddim: Generalized denoising diffusion
517 implicit models. In *The Eleventh International Conference on Learning Representations*, 2023.
518 URL <https://openreview.net/forum?id=1hKE9qjvz->.
- 519 [43] Martin Gonzalez, Nelson Fernandez Pinto, Thuy Tran, Hatem Hajri, Nader Masmoudi, et al.
520 Seeds: Exponential sde solvers for fast high-quality sampling from diffusion models. *Advances*
521 *in Neural Information Processing Systems*, 36, 2024.
- 522 [44] Shuchen Xue, Mingyang Yi, Weijian Luo, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and
523 Zhi-Ming Ma. Sa-solver: Stochastic adams solver for fast sampling of diffusion models.
524 *Advances in Neural Information Processing Systems*, 36, 2024.
- 525 [45] Qinsheng Zhang, Jiaming Song, and Yongxin Chen. Improved order analysis and design of
526 exponential integrator for diffusion models sampling. *arXiv preprint arXiv:2308.02157*, 2023.
- 527 [46] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion
528 solvers. *Advances in Neural Information Processing Systems*, 35:30150–30166, 2022.
- 529 [47] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion
530 models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- 531 [48] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of
532 diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:
533 26565–26577, 2022.
- 534 [49] Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. Dpm-solver-v3: Improved diffusion ode
535 solver with empirical model statistics. *Advances in Neural Information Processing Systems*,
536 36:55502–55542, 2023.
- 537 [50] Gen Li, Yu Huang, Timofey Efimov, Yuting Wei, Yuejie Chi, and Yuxin Chen. Accelerating
538 convergence of score-based diffusion models, provably. *arXiv preprint arXiv:2403.03852*,
539 2024.
- 540 [51] Yuchen Wu, Yuxin Chen, and Yuting Wei. Stochastic Runge-Kutta methods: Provable
541 acceleration of diffusion models. *arXiv preprint arXiv:2410.04760*, 2024.
- 542 [52] Daniel T Gillespie. Approximate accelerated stochastic simulation of chemically reacting
543 systems. *The Journal of chemical physics*, 115(4):1716–1733, 2001.

- [53] Yang Cao, Linda R Petzold, Muruhan Rathinam, and Daniel T Gillespie. The numerical stability of leaping methods for stochastic simulation of chemically reacting systems. *The Journal of chemical physics*, 121(24):12169–12178, 2004.
- [54] K Burrage and T Tian. Poisson runge-kutta methods for chemical reaction systems in advances in scientific computing and applications, 2004.
- [55] Yucheng Hu and Tiejun Li. Highly accurate tau-leaping methods with random corrections. *The Journal of chemical physics*, 130(12), 2009.
- [56] Desmond J Higham. Modeling and simulating chemical reactions. *SIAM review*, 50(2): 347–368, 2008.
- [57] Weinan E, Tiejun Li, and Eric Vanden-Eijnden. *Applied stochastic analysis*, volume 199. American Mathematical Soc., 2021.
- [58] Frank P Kelly. *Reversibility and stochastic networks*. Cambridge University Press, 2011.
- [59] 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.
- [60] Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly d -linear convergence bounds for diffusion models via stochastic localization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=r5njV3BsuD>.
- [61] Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability flow ode is provably fast. *Advances in Neural Information Processing Systems*, 36, 2024.
- [62] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [63] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [64] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. <http://SkyLion007.github.io/OpenWebTextCorpus>, 2019.
- [65] Victor Besnier and Mickael Chen. A pytorch reproduction of masked generative image transformer. *arXiv preprint arXiv:2310.14400*, 2023.
- [66] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [67] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.
- [68] Ilia Igashov, Arne Schneuing, Marwin Segler, Michael Bronstein, and Bruno Correia. Retrobridge: Modeling retrosynthesis with markov bridges. *arXiv preprint arXiv:2308.16212*, 2023.
- [69] Yang Li, Jinpei Guo, Runzhong Wang, and Junchi Yan. From distribution learning in training to gradient search in testing for combinatorial optimization. *Advances in Neural Information Processing Systems*, 36, 2024.
- [70] Zhiqing Sun and Yiming Yang. Difusco: Graph-based diffusion solvers for combinatorial optimization. *Advances in Neural Information Processing Systems*, 36:3706–3731, 2023.

- [71] Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Bac Nguyen, Stefano Ermon, and Yuki Mitsufuji. G2d2: Gradient-guided discrete diffusion for image inverse problem solving. *arXiv preprint arXiv:2410.14710*, 2024.
- [72] Wenda Chu, Yang Song, and Yisong Yue. Split gibbs discrete diffusion posterior sampling. *arXiv preprint arXiv:2503.01161*, 2025.
- [73] Cheuk Kit Lee, Paul Jeha, Jes Frellsen, Pietro Lio, Michael Samuel Albergo, and Francisco Vargas. Debiasing guidance for discrete diffusion with sequential monte carlo. *arXiv preprint arXiv:2502.06079*, 2025.
- [74] Peter Holderrieth, Michael S Albergo, and Tommi Jaakkola. Leaps: A discrete neural sampler via locally equivariant networks. *arXiv preprint arXiv:2502.10843*, 2025.
- [75] Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex X Lu, Nicolo Fusi, Ava P Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *BioRxiv*, pages 2023–09, 2023.
- [76] Patrick Emami, Aidan Perreault, Jeffrey Law, David Biagioni, and Peter St John. Plug & play directed evolution of proteins with gradient-based discrete mcmc. *Machine Learning: Science and Technology*, 4(2):025014, 2023.
- [77] Dmitry Penzar, Daria Nogina, Elizaveta Noskova, Arsenii Zinkevich, Georgy Meshcheryakov, Andrey Lando, Abdul Muntakim Rafi, Carl De Boer, and Ivan V Kulakovskiy. Legnet: a best-in-class deep learning model for short dna regulatory regions. *Bioinformatics*, 39(8):btad457, 2023.
- [78] Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.
- [79] John J Yang, Jason Yim, Regina Barzilay, and Tommi Jaakkola. Fast non-autoregressive inverse folding with discrete diffusion. *arXiv preprint arXiv:2312.02447*, 2023.
- [80] Hannes Stark, Bowen Jing, Chenyu Wang, Gabriele Corso, Bonnie Berger, Regina Barzilay, and Tommi Jaakkola. Dirichlet flow matching with applications to dna sequence design. *arXiv preprint arXiv:2402.05841*, 2024.
- [81] Kai Yi, Bingxin Zhou, Yiqing Shen, Pietro Liò, and Yuguang Wang. Graph denoising diffusion for inverse protein folding. *Advances in Neural Information Processing Systems*, 36, 2024.
- [82] Yiheng Zhu, Jialu Wu, Qiuyi Li, Jiahuan Yan, Mingze Yin, Wei Wu, Mingyang Li, Jieping Ye, Zheng Wang, and Jian Wu. Bridge-if: Learning inverse protein folding with markov bridges. *arXiv preprint arXiv:2411.02120*, 2024.
- [83] Patrick Esser, Robin Rombach, Andreas Blattmann, and Bjorn Ommer. Imagebart: Bidirectional context with multinomial diffusion for autoregressive image synthesis. *Advances in neural information processing systems*, 34:3518–3532, 2021.
- [84] Jose Lezama, Tim Salimans, Lu Jiang, Huiwen Chang, Jonathan Ho, and Irfan Essa. Discrete predictor-corrector diffusion models for image synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- [85] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022.
- [86] Ari Seff, Wenda Zhou, Farhan Damani, Abigail Doyle, and Ryan P Adams. Discrete object generation with reversible inductive construction. *Advances in neural information processing systems*, 32, 2019.
- [87] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. Permutation invariant graph generation via score-based generative modeling. In *International Conference on Artificial Intelligence and Statistics*, pages 4474–4484. PMLR, 2020.

- [88] Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang. Graphaf: a flow-based autoregressive model for molecular graph generation. *arXiv preprint arXiv:2001.09382*, 2020.
- [89] Yiming Qin, Clement Vignac, and Pascal Frossard. Sparse training of discrete diffusion models for graph generation. *arXiv preprint arXiv:2311.02142*, 2023.
- [90] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734*, 2022.
- [91] Kilian Konstantin Haefeli, Karolis Martinkus, Nathanaël Perraudin, and Roger Wattenhofer. Diffusion models for graphs benefit from discrete state spaces. *arXiv preprint arXiv:2210.01549*, 2022.
- [92] Yiming Qin, Manuel Madeira, Dorina Thanou, and Pascal Frossard. Defog: Discrete flow matching for graph generation. *arXiv preprint arXiv:2410.04263*, 2024.
- [93] Jun Hyeon Kim, Seonghwan Kim, Seokhyun Moon, Hyeonwoo Kim, Jeheon Woo, and Woo Youn Kim. Discrete diffusion schrödinger bridge matching for graph transformation. *arXiv preprint arXiv:2410.01500*, 2024.
- [94] Naoto Inoue, Kotaro Kikuchi, Edgar Simo-Serra, Mayu Otani, and Kota Yamaguchi. Layoutdm: Discrete diffusion model for controllable layout generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10167–10176, 2023.
- [95] Junyi Zhang, Jiaqi Guo, Shizhao Sun, Jian-Guang Lou, and Dongmei Zhang. Layoutdiffusion: Improving graphic layout generation by discrete diffusion probabilistic models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7226–7236, 2023.
- [96] Seunggeun Chi, Hyung-gun Chi, Hengbo Ma, Nakul Agarwal, Faizan Siddiqui, Karthik Ramani, and Kwonjoon Lee. M2d2m: Multi-motion generation from text with discrete diffusion models. *arXiv preprint arXiv:2407.14502*, 2024.
- [97] Yunhong Lou, Linchao Zhu, Yaxiong Wang, Xiaohan Wang, and Yi Yang. Diversemotion: Towards diverse human motion generation via discrete diffusion. *arXiv preprint arXiv:2309.01372*, 2023.
- [98] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:1720–1733, 2023.
- [99] Sam Bond-Taylor, Peter Hessey, Hiroshi Sasaki, Toby P Breckon, and Chris G Willcocks. Unleashing transformers: Parallel token prediction with discrete absorbing diffusion for fast high-resolution image generation from vector-quantized codes. In *European Conference on Computer Vision*, pages 170–188. Springer, 2022.
- [100] Zhicong Tang, Shuyang Gu, Jianmin Bao, Dong Chen, and Fang Wen. Improved vector quantized diffusion models. *arXiv preprint arXiv:2205.16007*, 2022.
- [101] Ye Zhu, Yu Wu, Kyle Olszewski, Jian Ren, Sergey Tulyakov, and Yan Yan. Discrete contrastive diffusion for cross-modal music and image generation. *arXiv preprint arXiv:2206.07771*, 2022.
- [102] Zhichao Wu, Qiulin Li, Sixing Liu, and Qun Yang. Dctts: Discrete diffusion model with contrastive learning for text-to-speech generation. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 11336–11340. IEEE, 2024.
- [103] Jun Han, Zixiang Chen, Yongqian Li, Yiwen Kou, Eran Halperin, Robert E Tillman, and Quanquan Gu. Guided discrete diffusion for electronic health record generation. *arXiv preprint arXiv:2404.12314*, 2024.

- [104] Juntong Shi, Minkai Xu, Harper Hua, Hengrui Zhang, Stefano Ermon, and Jure Leskovec. Tabdiff: a multi-modal diffusion model for tabular data generation. *arXiv preprint arXiv:2410.20626*, 2024.
- [105] Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusionbert: Improving generative masked language models with diffusion models. *arXiv preprint arXiv:2211.15029*, 2022.
- [106] Nikolay Savinov, Junyoung Chung, Mikolaj Binkowski, Erich Elsen, and Aaron van den Oord. Step-unrolled denoising autoencoders for text generation. *arXiv preprint arXiv:2112.06749*, 2021.
- [107] Tong Wu, Zhihao Fan, Xiao Liu, Hai-Tao Zheng, Yeyun Gong, Jian Jiao, Juntao Li, Jian Guo, Nan Duan, Weizhu Chen, et al. Ar-diffusion: Auto-regressive diffusion model for text generation. *Advances in Neural Information Processing Systems*, 36:39957–39974, 2023.
- [108] Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. Diffuseq-v2: Bridging discrete and continuous text spaces for accelerated seq2seq diffusion models. *arXiv preprint arXiv:2310.05793*, 2023.
- [109] Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. A reparameterized discrete diffusion model for text generation. *arXiv preprint arXiv:2302.05737*, 2023.
- [110] Kun Zhou, Yifan Li, Wayne Xin Zhao, and Ji-Rong Wen. Diffusion-nat: Self-prompting discrete diffusion for non-autoregressive text generation. *arXiv preprint arXiv:2305.04044*, 2023.
- [111] Minkai Xu, Tomas Geffner, Karsten Kreis, Weili Nie, Yilun Xu, Jure Leskovec, Stefano Ermon, and Arash Vahdat. Energy-based diffusion language models for text generation. *arXiv preprint arXiv:2410.21357*, 2024.
- [112] Yinuo Ren, Grant M Rotskoff, and Lexing Ying. A unified approach to analysis and design of denoising markov models. *arXiv preprint arXiv:2504.01938*, 2025.
- [113] Yong-Hyun Park, Chieh-Hsin Lai, Satoshi Hayakawa, Yuhta Takida, and Yuki Mitsufuji. Jump your steps: Optimizing sampling schedule of discrete diffusion models. *arXiv preprint arXiv:2410.07761*, 2024.
- [114] Yixiu Zhao, Jiaxin Shi, Lester Mackey, and Scott Linderman. Informed correctors for discrete diffusion models. *arXiv preprint arXiv:2407.21243*, 2024.
- [115] Zixiang Chen, Huizhuo Yuan, Yongqian Li, Yiwen Kou, Junkai Zhang, and Quanquan Gu. Fast sampling via discrete non-markov diffusion models with predetermined transition time. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [116] Lingxiao Zhao, Xueying Ding, Lijun Yu, and Leman Akoglu. Improving and unifying discrete&continuous-time discrete denoising diffusion. *arXiv preprint arXiv:2402.03701*, 2024.
- [117] Machel Reid, Vincent Josua Hellendoorn, and Graham Neubig. Diffuser: Diffusion via edit-based reconstruction. In *The Eleventh International Conference on Learning Representations*, 2023.
- [118] Satoshi Hayakawa, Yuhta Takida, Masaaki Imaizumi, Hiromi Wakaki, and Yuki Mitsufuji. Distillation of discrete diffusion through dimensional correlations. *arXiv preprint arXiv:2410.08709*, 2024.
- [119] Ludwig Winkler, Lorenz Richter, and Manfred Opper. Bridging discrete and continuous state spaces: Exploring the ehrenfest process in time-continuous diffusion models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=8GYclcxQXB>.
- [120] Harshit Varma, Dheeraj Nagaraj, and Karthikeyan Shanmugam. Glauber generative model: Discrete diffusion models via binary classification. *arXiv preprint arXiv:2405.17035*, 2024.

- 733 [121] Nate Gruver, Samuel Stanton, Nathan Frey, Tim GJ Rudner, Isidro Hotzel, Julien Lafrance-
734 Vanasse, Arvind Rajpal, Kyunghyun Cho, and Andrew G Wilson. Protein design with guided
735 discrete diffusion. *Advances in neural information processing systems*, 36, 2024.
- 736 [122] Severi Rissanen, Markus Heinonen, and Arno Solin. Improving discrete diffusion models via
737 structured preferential generation. *arXiv preprint arXiv:2405.17889*, 2024.
- 738 [123] Sulin Liu, Juno Nam, Andrew Campbell, Hannes Stärk, Yilun Xu, Tommi Jaakkola, and Rafael
739 Gómez-Bombarelli. Think while you generate: Discrete diffusion with planned denoising.
740 *arXiv preprint arXiv:2410.06264*, 2024.
- 741 [124] Oscar Davis, Samuel Kessler, Mircea Petrache, Avishek Joey Bose, et al. Fisher flow matching
742 for generative modeling over discrete data. *arXiv preprint arXiv:2405.14664*, 2024.
- 743 [125] GN Mil'shtejn. Approximate integration of stochastic differential equations. *Theory of*
744 *Probability & Its Applications*, 19(3):557–562, 1975.
- 745 [126] Assy Abdulle and Stephane Cirilli. S-rock: Chebyshev methods for stiff stochastic differential
746 equations. *SIAM Journal on Scientific Computing*, 30(2):997–1014, 2008.
- 747 [127] Evelyn Buckwar and Renate Winkler. Multistep methods for sdes and their application to
748 problems with small noise. *SIAM journal on numerical analysis*, 44(2):779–803, 2006.
- 749 [128] Denis Talay and Luciano Tubaro. Expansion of the global error for numerical schemes solving
750 stochastic differential equations. *Stochastic analysis and applications*, 8(4):483–509, 1990.
- 751 [129] Kevin Burrage and Pamela M Burrage. Order conditions of stochastic runge–kutta methods by
752 b-series. *SIAM Journal on Numerical Analysis*, 38(5):1626–1646, 2000.
- 753 [130] Kevin Burrage and Tianhai Tian. Predictor-corrector methods of runge–kutta type for stochastic
754 differential equations. *SIAM Journal on Numerical Analysis*, 40(4):1516–1537, 2002.
- 755 [131] Andreas Rössler. Runge-kutta methods for the numerical solution of stochastic differential
756 equations. *Shaker-Verlag, Aachen*, 2003.
- 757 [132] Andreas Rößler. Runge–kutta methods for the strong approximation of solutions of stochastic
758 differential equations. *SIAM Journal on Numerical Analysis*, 48(3):922–952, 2010.
- 759 [133] James M Foster, Goncalo Dos Reis, and Calum Strange. High order splitting methods for sdes
760 satisfying a commutativity condition. *SIAM Journal on Numerical Analysis*, 62(1):500–532,
761 2024.
- 762 [134] Lei Li, Jianfeng Lu, Jonathan Mattingly, and Lihan Wang. Numerical methods for stochastic
763 differential equations based on gaussian mixtures. *Communications in Mathematical Sciences*,
764 19(6):1549–1577, 2021.
- 765 [135] Ruoqi Shen and Yin Tat Lee. The randomized midpoint method for log-concave sampling.
766 *Advances in Neural Information Processing Systems*, 32, 2019.
- 767 [136] Nima Anari, Sinho Chewi, and Thuy-Duong Vuong. Fast parallel sampling under isoperimetry.
768 *arXiv preprint arXiv:2401.09016*, 2024.
- 769 [137] Lu Yu and Arnak Dalalyana. Parallelized midpoint randomization for langevin monte carlo.
770 *arXiv preprint arXiv:2402.14434*, 2024.
- 771 [138] Changyou Chen, Nan Ding, and Lawrence Carin. On the convergence of stochastic gradient
772 mcmc algorithms with high-order integrators. *Advances in neural information processing*
773 *systems*, 28, 2015.
- 774 [139] Alain Durmus, Umut Simsekli, Eric Moulines, Roland Badeau, and Gaël Richard. Stochastic
775 gradient richardson-romberg markov chain monte carlo. *Advances in neural information*
776 *processing systems*, 29, 2016.

- 777 [140] Xuechen Li, Yi Wu, Lester Mackey, and Murat A Erdogdu. Stochastic runge-kutta accelerates
778 langevin monte carlo and beyond. *Advances in neural information processing systems*, 32,
779 2019.
- 780 [141] Sotirios Sabanis and Ying Zhang. Higher order langevin monte carlo algorithm. *Electron. J.*
781 *Statist*, 13(2):3805–3850, 2019.
- 782 [142] Wenlong Mou, Yi-An Ma, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan. High-
783 order langevin diffusion yields an accelerated mcmc algorithm. *Journal of Machine Learning*
784 *Research*, 22(42):1–41, 2021.
- 785 [143] Pierre Monmarché. High-dimensional mcmc with a standard splitting scheme for the under-
786 damped langevin diffusion. *Electronic Journal of Statistics*, 15(2):4117–4166, 2021.
- 787 [144] James Foster, Terry Lyons, and Harald Oberhauser. The shifted ode method for underdamped
788 langevin mcmc. *arXiv preprint arXiv:2101.03446*, 2021.
- 789 [145] Kevin Burrage, PM Burrage, and Tianhai Tian. Numerical methods for strong solutions of
790 stochastic differential equations: an overview. *Proceedings of the Royal Society of London.*
791 *Series A: Mathematical, Physical and Engineering Sciences*, 460(2041):373–402, 2004.
- 792 [146] Grigori N Milstein and Michael V Tretyakov. *Stochastic numerics for mathematical physics*,
793 volume 39. Springer, 2004.
- 794 [147] Alfred B Bortz, Malvin H Kalos, and Joel L Lebowitz. A new algorithm for monte carlo
795 simulation of ising spin systems. *Journal of Computational physics*, 17(1):10–18, 1975.
- 796 [148] Daniel T Gillespie. A general method for numerically simulating the stochastic time evolution
797 of coupled chemical reactions. *Journal of computational physics*, 22(4):403–434, 1976.
- 798 [149] Daniel T Gillespie. Exact stochastic simulation of coupled chemical reactions. *The journal of*
799 *physical chemistry*, 81(25):2340–2361, 1977.
- 800 [150] Yang Cao, Dan Gillespie, and Linda Petzold. Multiscale stochastic simulation algorithm
801 with stochastic partial equilibrium assumption for chemically reacting systems. *Journal of*
802 *Computational Physics*, 206(2):395–411, 2005.
- 803 [151] Yang Cao, Daniel T Gillespie, and Linda R Petzold. The slow-scale stochastic simulation
804 algorithm. *The Journal of chemical physics*, 122(1), 2005.
- 805 [152] Weinan E, Di Liu, Eric Vanden-Eijnden, et al. Nested stochastic simulation algorithm for
806 chemical kinetic systems with disparate rates. *The Journal of chemical physics*, 123(19), 2005.
- 807 [153] Weinan E, Di Liu, and Eric Vanden-Eijnden. Nested stochastic simulation algorithms for
808 chemical kinetic systems with multiple time scales. *Journal of computational physics*, 221(1):
809 158–180, 2007.
- 810 [154] Michael A Gibson and Jehoshua Bruck. Efficient exact stochastic simulation of chemical
811 systems with many species and many channels. *The journal of physical chemistry A*, 104(9):
812 1876–1889, 2000.
- 813 [155] David F Anderson. A modified next reaction method for simulating chemical systems with
814 time dependent propensities and delays. *The Journal of chemical physics*, 127(21), 2007.
- 815 [156] Casper HL Beentjes and Ruth E Baker. Uniformization techniques for stochastic simulation of
816 chemical reaction networks. *The Journal of Chemical Physics*, 150(15), 2019.
- 817 [157] Muruhan Rathinam, Linda R Petzold, Yang Cao, and Daniel T Gillespie. Stiffness in stochastic
818 chemically reacting systems: The implicit tau-leaping method. *The Journal of Chemical*
819 *Physics*, 119(24):12784–12794, 2003.
- 820 [158] Daniel T Gillespie and Linda R Petzold. Improved leap-size selection for accelerated stochastic
821 simulation. *The journal of chemical physics*, 119(16):8229–8234, 2003.

- [159] Kevin Burrage, Tianhai Tian, and Pamela Burrage. A multi-scaled approach for simulating chemical reaction systems. *Progress in biophysics and molecular biology*, 85(2-3):217–234, 2004.
- [160] Yang Cao, Daniel T Gillespie, and Linda R Petzold. Avoiding negative populations in explicit poisson tau-leaping. *The Journal of chemical physics*, 123(5), 2005.
- [161] Anne Auger, Philippe Chatelain, and Petros Koumoutsakos. R-leaping: Accelerating the stochastic simulation algorithm by reaction leaps. *The Journal of chemical physics*, 125(8), 2006.
- [162] Yang Cao, Daniel T Gillespie, and Linda R Petzold. Adaptive explicit-implicit tau-leaping method with automatic tau selection. *The Journal of chemical physics*, 126(22), 2007.
- [163] Basil Bayati, Philippe Chatelain, and Petros Koumoutsakos. D-leaping: Accelerating stochastic simulation algorithms for reactions with delays. *Journal of Computational Physics*, 228(16): 5908–5916, 2009.
- [164] Yang Cao and Linda Petzold. Slow-scale tau-leaping method. *Computer methods in applied mechanics and engineering*, 197(43-44):3472–3479, 2008.
- [165] Zhouyi Xu and Xiaodong Cai. Unbiased τ -leap methods for stochastic simulation of chemically reacting systems. *The Journal of chemical physics*, 128(15), 2008.
- [166] Mary Sehl, Alexander V Alekseyenko, and Kenneth L Lange. Accurate stochastic simulation via the step anticipation τ -leaping (sal) algorithm. *Journal of Computational Biology*, 16(9): 1195–1208, 2009.
- [167] Krishna A Iyengar, Leonard A Harris, and Paulette Clancy. Accurate implementation of leaping in space: The spatial partitioned-leaping algorithm. *The Journal of chemical physics*, 132(9), 2010.
- [168] David F Anderson and Desmond J Higham. Multilevel monte carlo for continuous time markov chains, with applications in biochemical kinetics. *Multiscale Modeling & Simulation*, 10(1):146–179, 2012.
- [169] Alvaro Moraes, Raúl Tempone, and Pedro Vilanova. Hybrid chernoff tau-leap. *Multiscale Modeling & Simulation*, 12(2):581–615, 2014.
- [170] Jill Padgett and Silvana Ilie. An adaptive tau-leaping method for stochastic simulations of reaction-diffusion systems. *AIP Advances*, 6(3), 2016.
- [171] Jana Lipková, Georgios Arampatzis, Philippe Chatelain, Bjoern Menze, and Petros Koumoutsakos. S-leaping: an adaptive, accelerated stochastic simulation algorithm, bridging τ -leaping and r-leaping. *Bulletin of mathematical biology*, 81(8):3074–3096, 2019.
- [172] Muruhan Rathinam, Linda R Petzold, Yang Cao, and Daniel T Gillespie. Consistency and stability of tau-leaping schemes for chemical reaction systems. *Multiscale Modeling & Simulation*, 4(3):867–895, 2005.
- [173] Tiejun Li. Analysis of explicit tau-leaping schemes for simulating chemically reacting systems. *Multiscale Modeling & Simulation*, 6(2):417–436, 2007.
- [174] Yucheng Hu, Tiejun Li, and Bin Min. The weak convergence analysis of tau-leaping methods: revisited. *Communications in Mathematical Sciences*, 9(4):965–996, 2011.
- [175] David F Anderson, Desmond J Higham, and Yu Sun. Complexity of multilevel monte carlo tau-leaping. *SIAM Journal on Numerical Analysis*, 52(6):3106–3127, 2014.
- [176] Chuchu Chen and Di Liu. Error analysis for d-leaping scheme of chemical reaction system with delay. *Multiscale Modeling & Simulation*, 15(4):1797–1829, 2017.
- [177] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

- [178] Linfeng Zhang, Weinan E, and Lei Wang. Monge-ampère flow for generative modeling. *arXiv preprint arXiv:1809.10188*, 2018.
- [179] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [180] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [181] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [182] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in neural information processing systems*, 34: 1415–1428, 2021.
- [183] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [184] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- [185] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- [186] Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- [187] Haoxuan Chen, Yinuo Ren, Lexing Ying, and Grant M Rotskoff. Accelerating diffusion models with parallel sampling: Inference at sub-linear time complexity. *arXiv preprint arXiv:2405.15986*, 2024.
- [188] Zhenyu Zhou, Defang Chen, Can Wang, and Chun Chen. Fast ode-based sampling for diffusion models in around 5 steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7777–7786, 2024.
- [189] Yilun Xu, Mingyang Deng, Xiang Cheng, Yonglong Tian, Ziming Liu, and Tommi Jaakkola. Restart sampling for improving generative processes. *Advances in Neural Information Processing Systems*, 36:76806–76838, 2023.
- [190] Hanzhong Guo, Cheng Lu, Fan Bao, Tianyu Pang, Shuicheng Yan, Chao Du, and Chongxuan Li. Gaussian mixture solvers for diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [191] Jonathan Heek, Emiel Hoogetboom, and Tim Salimans. Multistep consistency models. *arXiv preprint arXiv:2403.06807*, 2024.
- [192] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- [193] Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *arXiv preprint arXiv:2310.14189*, 2023.
- [194] Cheng Lu and Yang Song. Simplifying, stabilizing and scaling continuous-time consistency models. *arXiv preprint arXiv:2410.11081*, 2024.
- [195] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.
- [196] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306, 2023.

- [197] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- [198] Vinh Tong, Trung-Dung Hoang, Anji Liu, Guy Van den Broeck, and Mathias Niepert. Learning to discretize denoising diffusion odes. *arXiv preprint arXiv:2405.15506*, 2024.
- [199] Eric Frankel, Sitan Chen, Jerry Li, Pang Wei Koh, Lillian J Ratliff, and Sewoong Oh. S4s: Solving for a diffusion model solver. *arXiv preprint arXiv:2502.17423*, 2025.
- [200] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped langevin diffusion. *arXiv preprint arXiv:2112.07068*, 2021.
- [201] Hongkai Zheng, Weili Nie, Arash Vahdat, Kamyar Azizzadenesheli, and Anima Anandkumar. Fast sampling of diffusion models via operator learning. In *International Conference on Machine Learning*, pages 42390–42402. PMLR, 2023.
- [202] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15762–15772, 2024.
- [203] Valentin De Bortoli, Alexandre Galashov, Arthur Gretton, and Arnaud Doucet. Accelerated diffusion models via speculative sampling. *arXiv preprint arXiv:2501.05370*, 2025.
- [204] Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021.
- [205] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [206] Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [207] Hyungjin Chung, Jeongsol Kim, Sehui Kim, and Jong Chul Ye. Parallel diffusion models of operator and image for blind inverse problems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6059–6069, 2023.
- [208] Zhiwei Tang, Jiasheng Tang, Hao Luo, Fan Wang, and Tsung-Hui Chang. Accelerating parallel sampling of diffusion models. In *Forty-first International Conference on Machine Learning*, 2024.
- [209] Jiezhong Cao, Yue Shi, Kai Zhang, Yulun Zhang, Radu Timofte, and Luc Van Gool. Deep equilibrium diffusion restoration with parallel sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2824–2834, 2024.
- [210] Nikil Roashan Selvam, Amil Merchant, and Stefano Ermon. Self-refining diffusion samplers: Enabling parallelization via parareal iterations. *arXiv preprint arXiv:2412.08292*, 2024.
- [211] Saravanan Kandasamy and Dheeraj Nagaraj. The poisson midpoint method for langevin dynamics: Provably efficient discretization for diffusion models. *arXiv preprint arXiv:2405.17068*, 2024.
- [212] Shivam Gupta, Linda Cai, and Sitan Chen. Faster diffusion-based sampling with randomized midpoints: Sequential and parallel. *arXiv preprint arXiv:2406.00924*, 2024.
- [213] Philip Protter. Point process differentials with evolving intensities. In *Nonlinear stochastic problems*, pages 467–472. Springer, 1983.
- [214] Bernt Øksendal and Agnes Sulem. *Applied Stochastic Control of Jump Diffusions*. Springer, 2019.

A Further Discussion on Related Works

In this section, we provide a more detailed literature review of both continuous and discrete diffusion models, as well as several studies on the numerical methods for SDEs and chemical reaction systems, which are highly related to our work.

Discrete Diffusion Models: Methodology, Theory, and Applications. Discrete diffusion and flow-based models [1–4, 6–10, 22] have recently been proposed as generalizations of continuous diffusion models to model discrete distributions.

Such models have been widely used in various areas of science and engineering, including but not limited to modeling retrosynthesis [68], combinatorial optimization [69, 70], solving inverse problems [71, 72] and sampling high-dimensional discrete distributions [73, 74], designing molecules, proteins, and DNA sequences [75, 13, 76, 12, 77–79, 31, 80, 11, 81, 82], image synthesis [83–85], text summarization [15], as well as the generation of graph [86–93], layout [94, 95], motion [96, 97], sound [22, 98], image [16, 99–101], speech [102], electronic health record [103], tabular data [104] and text [105–110, 34, 35, 111, 14]. Inspired by the huge success achieved by discrete diffusion models in practice, researchers have also conducted some studies on the theoretical properties of these models, such as [36–38, 112].

An extensive amount of work has also explored the possibility of making discrete diffusion models more effective from many aspects, such as optimizing the sampling schedule [113], adding correctors [114], developing fast samplers [115], designing correctors based on information learnt by the model [114], simplifying the loss function for training [116], adding editing-based refinements [117], synergizing these models with other techniques and methodologies like distillation [118], Ehrenfest processes [119], Glauber dynamics [120], tensor networks [17], enhanced guidance mechanisms [121, 18–20], structured preferential generation [122], the plan-and-denoise framework [123] and alternative metrics, *e.g.*, the Fisher information metric [124]. However, to the best of our knowledge, existing work on accelerating the inference of discrete diffusion models is relatively sparse compared to the ones we listed above, which makes it a direction worth exploring and serves as one of the main motivations behind this work.

Numerical Methods for SDEs and Chemical Reaction Systems. Below, we review advanced numerical methods proposed for simulating SDEs and chemical reaction systems, which are the main techniques adopted in our work. For the simulation of SDEs driven by Brownian motions, many studies have been performed to design more accurate numerical schemes, which have been widely applied to tackle problems in computational physics, optimization, and Monte Carlo sampling. Examples of such work include the Milstein method [125], explicit methods [126], multistep methods [127], extrapolation-type methods [128, 25], stochastic Runge Kutta methods [24, 129–132], splitting methods [133], methods based on gaussian mixtures [134], randomized midpoint method [135], parallel sampling methods [136, 137] as well as high-order methods for stochastic gradient Markov Chain Monte Carlo [138, 139], underdamped and overdamped Langevin Monte Carlo [140–144]. For a more comprehensive list of related numerical methods, one may refer to [39, 145, 146, 40, 57].

Regarding the simulation of chemical reaction systems, numerical methods can be categorized into two classes. The first class consists of exact simulation methods, which are similar to the Kinetic Monte Carlo (KMC) method [147] developed for simulating spin dynamics and crystal growth in condensed matter physics. Examples of such methods include the Gillespie algorithm (or the Stochastic Simulation Algorithm, a.k.a. SSA) [148, 149] and its variants for multiscale modeling [150–153], the next reaction method and its variants [154, 155], uniformization-based methods [156], etc. The second class of methods are approximate simulation methods, including but not limited to the τ -leaping method [52] and its variants [157, 158, 53, 54, 159–165, 55, 166, 167, 26, 168–171]. For a subset of the methods listed above, numerical analysis has also been performed in many works [172–176] to justify their validity.

Continuous Diffusion Models: Methodology, Theory, and Acceleration. Continuous diffusion and probability flow-based models [177–186] have also been the most popular methods in generative modeling, with a wide range of applications in science and engineering. For a list of related work on the theoretical studies and applications of these models, one may refer to the literature review

conducted in [187, 38]. Here we will only review studies on accelerating the inference of continuous diffusion models, which motivates our work.

An incomplete list of accelerating methods includes approximate mean direction solver [188], restart sampling [189], gaussian mixture solvers [190], self-consistency [191–194], knowledge distillation [195–199], combination with underdamped Langevin dynamics [200], operator learning [201] and more recently ideas from accelerating large language models (LLMs) like caching [202] and speculative decoding [203]. Among all the proposed accelerating methods, one major class of methods are developed based on techniques from numerical analysis like adaptive step sizes [204], exponential integrators [41–43], predictor-corrector solver [205], Adams-Bashforth methods [29, 44, 45], Taylor methods [27, 46], Picard iteration and parallel sampling [206–210, 187], (stochastic) Runge-Kutta methods [47, 28, 48–51] and randomized midpoint method [211, 212]. In contrast, there have been fewer studies on the acceleration of discrete diffusion models via techniques from numerical analysis, which inspires the study undertaken in this paper.

B Mathematical Background

In this section, we provide the mathematical background for the stochastic integral formulation of discrete diffusion models, the error analysis of the τ -leaping method, and useful lemmas for the theoretical analysis of high-order schemes for discrete diffusion models.

B.1 Stochastic Integral Formulation of Discrete Diffusion Models

Throughout this section, we will assume that $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space, \mathbb{X} is a finite-state space, and denote the pairwise difference set of the state space by $\mathbb{D} := \{x - y : x \neq y \in \mathbb{X}\}$. We also assume that the pairwise difference set \mathbb{X} is equipped with a metric $\|\cdot\|$, a finite measure γ , and a σ -algebra \mathcal{B} .

As a warm-up, we introduce the definition of the Poisson random measure for a time-homogeneous counting process.

Definition B.1 (Poisson Random Measure [38, Definition A.1]). *The random measure $N(dt, d\nu)$ on $\mathbb{R}^+ \times \mathbb{D}$ is called a Poisson random measure w.r.t. measure γ if it is a random counting measure satisfying the following properties:*

(i) *For any $B \in \mathcal{B}$ and $0 \leq s < t$,*

$$N((s, t] \times B) \sim \mathcal{P}(\gamma(B)(t - s));$$

(ii) *For any $t \geq 0$ and pairwise disjoint sets $\{B_i\}_{i \in [n]} \subset \mathcal{B}$,*

$$\{N_t(B_i) := N((0, t] \times B_i)\}_{i \in [n]}$$

are independent stochastic processes.

Then we define the Poisson random measure with evolving intensities. The term “evolving” refers to that the intensity is both time and state-dependent.

Definition B.2 (Poisson Random Measure with Evolving Intensity [38, Definition A.3]). *Suppose $\lambda_t(y)$ is a non-negative predictable process on $\mathbb{R}^+ \times \mathbb{D} \times \Omega$ satisfying that for any $0 \leq T < \bar{T}$, $\int_0^T \lambda_t(\nu) dt < \infty$, a.s..*

The random measure $N[\lambda](dt, d\nu)$ on $\mathbb{R}^+ \times \mathbb{D}$ is called a Poisson random measure with evolving intensity $\lambda_t(\nu)$ w.r.t. measure γ if it is a random counting measure satisfying the following properties:

(i) *For any $B \in \mathcal{B}$ and $0 \leq s < t$,*

$$N[\lambda]((s, t] \times B) \sim \mathcal{P}\left(\int_s^t \int_B \lambda_\tau(\nu) \gamma(d\nu) d\tau\right);$$

(ii) *For any $t \geq 0$ and pairwise disjoint sets $\{B_i\}_{i \in [n]} \subset \mathcal{B}$,*

$$\{N_t[\lambda](B_i) := N[\lambda]((0, t] \times B_i)\}_{i \in [n]}$$

are independent stochastic processes.

1047 **Remark B.3** (Construction of Poisson Random Measure with Evolving Intensity). *As discussed*
 1048 *in Thm. A.4 in [38] and originally proposed by [213], the Poisson random measure with evolving*
 1049 *intensity can be constructed in the following way.*

1050 *One first augments the $(\mathbb{X}, \mathcal{B}, \nu)$ measure space to a product space $(\mathbb{D} \times \mathbb{R}, \mathcal{B} \times \mathcal{B}(\mathbb{R}), \gamma \times m)$,*
 1051 *where m is the Lebesgue measure on \mathbb{R} , and $\mathcal{B}(\mathbb{R})$ is the Borel σ -algebra on \mathbb{R} . The Poisson random*
 1052 *measure with evolving intensity $\lambda_t(\nu)$ can be defined in the augmented measure space as*

$$N[\lambda]((s, t] \times B) := \int_s^t \int_B \int_{\mathbb{R}} \mathbf{1}_{0 \leq \xi \leq \lambda_\tau(\nu)} N(d\tau, d\nu, d\xi), \quad (\text{B.1})$$

1053 *where $N(d\tau, d\nu, d\xi)$ is the Poisson random measure on $\mathbb{R}^+ \times \mathbb{D} \times \mathbb{R}$ w.r.t. measure $\nu(dy)d\xi$.*

1054 The following theorem provides the change of measure theorem for Poisson random measure with
 1055 evolving intensity, which is crucial for the theoretical analysis of numerical schemes for discrete
 1056 diffusion models.

1057 **Theorem B.4** (Change of Measure for Poisson Random Measure with Evolving Density [38,
 1058 Thm. 3.3]). *Let $N[\lambda](dt, d\nu)$ be a Poisson random measure with evolving intensity $\lambda_t(\nu)$, and*
 1059 *$h_t(\nu)$ a positive predictable process on $\mathbb{R}^+ \times \mathbb{D} \times \Omega$. Suppose the following exponential process is a*
 1060 *local \mathcal{F}_t -martingale:*

$$Z_t[h] := \exp \left(\int_0^t \int_{\mathbb{D}} \log h_t(\nu) N[\lambda](dt \times d\nu) - \int_0^t \int_{\mathbb{D}} (h_t(\nu) - 1) \lambda_t(\nu) \gamma(d\nu) \right), \quad (\text{B.2})$$

1061 *and \mathbb{Q} is another probability measure on (Ω, \mathcal{F}) such that $\mathbb{Q} \ll \mathbb{P}$ with Radon-Nikodym derivative*
 1062 *$d\mathbb{Q}/d\mathbb{P}|_{\mathcal{F}_t} = Z_t[h]$.*

1063 *Then the Poisson random measure $N[\lambda](dt, d\nu)$ under the measure \mathbb{Q} is a Poisson random measure*
 1064 *with evolving intensity $\lambda_t(\nu)h_t(\nu)$.*

1065 B.2 Error Analysis of τ -leaping

1066 The τ -leaping method was originally proposed by [52] and adopted for the inference of discrete
 1067 diffusion models by [22]. A summary of the algorithm is given in Alg. 3. In this subsection, we
 1068 provide a sketch for the error analysis of the τ -leaping method when applied to discrete diffusion
 1069 models, which will be compared with that of high-order schemes later on.

Algorithm 3: τ -Leaping Method for Discrete Diffusion Model Inference

Input: $\hat{y}_0 \sim q_0$, $\theta \in [0, 1]$, time discretization $(s_n, \rho_n)_{n \in [0:N-1]}$, $\hat{\mu}$, $\hat{\mu}^*$ as defined in Prop. C.2.

Output: A sample $\hat{y}_{s_N} \sim \hat{q}_{t_N}^{\text{RK}}$.

```

1 for  $n = 0$  to  $N - 1$  do
2    $\hat{y}_{s_{n+1}} \leftarrow \hat{y}_{s_n} + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}(\nu) \Delta_n);$ 
3 end
```

1070 *Proof of Thm. 3.1.* As we are considering the case where $\mathbb{X} = [S]^d$, i.e. the state space is a d -
 1071 dimensional grid with S states along each dimension, we have $\log |\mathbb{X}| = d \log S$. Then we consider a
 1072 simple time-homogeneous transition matrix $\mathbf{Q}_t \equiv \mathbf{Q}$ that allows jumps between neighboring states
 1073 with equal probability. Specifically, we have

$$Q(y, x) = \begin{cases} 1, & \|x - y\|_1 = 1, \\ -2d, & x = y, \end{cases}$$

1074 which can be verified to satisfy Assumption 4.3(i) in [38] with $C = 1$ and $\underline{D} = \overline{D} = 2d$. Assump-
 1075 tion 4.3(ii) is also satisfied, as shown in Example B.10 of [38].

1076 Then we may apply Thm. 4.7 in [38] by using the required time discretization scheme according to
 1077 the properties of the target distribution and plugging in the corresponding values of $C, \underline{D}, \overline{D}$. The
 1078 result follows by scaling the transition matrix \mathbf{Q} by $\frac{1}{d}$, equivalent to scaling the time by d . \square

1079 C Proofs

1080 In this section, we provide the missing proofs in the main text. We will first provide the proofs of
 1081 the stochastic integral formulations of high-order schemes for discrete diffusion models in App. C.1.
 1082 Then we will provide the proofs of the main results for the θ -trapezoidal method in App. C.2 and the
 1083 θ -RK-2 method in App. C.3. We remark that the proof for the θ -trapezoidal method requires more
 1084 techniques and is more involved, to which the proof for the θ -RK-2 method is analogous. In App. C.4,
 1085 we provide the detailed lemmas and computations omitted in the proofs of Thm. 5.4 and Thm. 5.5.

1086 C.1 Stochastic Integral Formulations of High-Order Schemes

1087 In order to rigorously analyze the θ -RK-2 method, we need the following definition:

1088 **Definition C.1** (Intermediate Process). *We define the intermediate process \hat{y}_s^* piecewisely on*
 1089 *$(s_n, s_{n+1}]$ as follows:*

$$\hat{y}_s^* = \hat{y}_{s_n} + \int_{s_n}^s \int_{\mathbb{D}} \nu N[\hat{\mu}_{s_n}](ds, d\nu), \quad (\text{C.1})$$

1090 *where the intensity $\hat{\mu}_{s_n}$ is given by $\hat{\mu}_{s_n}(\nu, \hat{y}_{s_n}) = \tilde{s}_{s_n}(\hat{y}_{s_n}, \hat{y}_{s_n} + \nu) \tilde{Q}_{s_n}^0(\hat{y}_{s_n}, \hat{y}_{s_n} + \nu)$, i.e., \hat{y}_s^* is the*
 1091 *process obtained by performing τ -leaping from time s_n to s with intensity $\hat{\mu}$.*

1092 The following proposition provides the stochastic integral formulation of this method.

1093 **Proposition C.2** (Stochastic Integral Formulation of θ -RK-2 Method). *The θ -RK-2 method (Alg. 1)*
 1094 *is equivalent to solving the following stochastic integral:*

$$\hat{y}_s^{\text{RK}} = \hat{y}_0^{\text{RK}} + \int_0^s \int_{\mathbb{D}} \nu N[\hat{\mu}^{\text{RK}}](ds, d\nu), \quad (\text{C.2})$$

1095 *in which the intensity $\hat{\mu}^{\text{RK}}$ is defined as a weighted sum*

$$\hat{\mu}_s^{\text{RK}}(\nu) = (1 - \frac{1}{2\theta}) \hat{\mu}_{\lfloor s \rfloor}(\nu, \hat{y}_{\lfloor s \rfloor}^{\text{RK}}) + \frac{1}{2\theta} \hat{\mu}_{\rho_s}^*(\nu, \hat{y}_{\rho_s}^*), \quad (\text{C.3})$$

1096 *and the intermediate intensity $\hat{\mu}^*$ is defined piecewisely as*

$$\hat{\mu}_s^*(\nu, \hat{y}_s^*) = \tilde{s}_s(\hat{y}_s^*, \hat{y}_s^* + \nu) \tilde{Q}_s^0(\hat{y}_s^*, \hat{y}_s^* + \nu), \quad (\text{C.4})$$

1097 *with the intermediate process \hat{y}_s^* defined in (C.1) for the corresponding interval. We will call \hat{y}_s^{RK} the*
 1098 *interpolating process of the θ -RK-2 method and denote the distribution of \hat{y}_s^{RK} by \hat{q}_s^{RK} .*

1099 The following proposition establishes the stochastic integral formulation of the θ -trapezoidal method,
 1100 whose proof can be found in App. C.1.

1101 **Proposition C.3** (Stochastic Integral Formulation of θ -Trapezoidal Method). *The θ -trapezoidal*
 1102 *method (Alg. 2) is equivalent to solving the following stochastic integral:*

$$\hat{y}_s^{\text{trap}} = \hat{y}_0^{\text{trap}} + \int_0^s \int_{\mathbb{D}} N[\hat{\mu}^{\text{trap}}](ds, d\nu) \quad (\text{C.5})$$

1103 *where the intensity $\hat{\mu}^{\text{trap}}$ is defined piecewisely as*

$$\hat{\mu}_s^{\text{trap}}(\nu) = \mathbf{1}_{s < \rho_s} \hat{\mu}_{\lfloor s \rfloor}(\nu, \hat{y}_{\lfloor s \rfloor}^{\text{trap}}) + \mathbf{1}_{s \geq \rho_s} \left(\alpha_1 \hat{\mu}_{\rho_s}^*(\nu, \hat{y}_{\rho_s}^*) - \alpha_2 \hat{\mu}_{\lfloor s \rfloor}(\nu, \hat{y}_{\lfloor s \rfloor}^{\text{trap}}) \right)_+. \quad (\text{C.6})$$

1104 *Above, $\mathbf{1}_{(\cdot)}$ denotes the indicator function and the intermediate process \hat{y}_s^* is defined in (C.1) for the*
 1105 *corresponding interval. We will call the process \hat{y}_s^{trap} the interpolating process of the θ -trapezoidal*
 1106 *method and denote the distribution of \hat{y}_s^{trap} by \hat{q}_s^{trap} .*

1107 *Proof of Prop. C.2 and Prop. C.3.* Without loss of generality, we give the proof on the interval
 1108 $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, and the generalization to the whole interval $[0, T]$ is straightforward.

1109 Notice that once we condition on the filtration \mathcal{F}_{s_n} and construct the intermediate process \hat{y}_s^* as
 1110 specified in (C.1) along the interval $(s_n, s_{n+1}]$, the intermediate intensity $\hat{\mu}^*$ and the piecewise
 1111 intensity $\hat{\mu}_{\lfloor s \rfloor}$ do not evolve with time s or the interpolating processes \hat{y}_s^{RK} (or \hat{y}_s^{trap} , respectively)

1112 since it only depends on the state, the intensity at the beginning of the interval s_n and other randomness
 1113 that is independent of the interpolating process.

1114 Therefore, the stochastic integral on this interval can be rewritten as for the θ -RK-2 scheme that

$$\begin{aligned}\hat{y}_{s_{n+1}}^{\text{RK}} &= \hat{y}_{s_n}^{\text{RK}} + \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \nu N[\hat{\mu}^{\text{trap}}](ds, d\nu) \\ &= \hat{y}_{s_n}^{\text{RK}} + \int_{\mathbb{D}} \nu N[\hat{\mu}^{\text{RK}}]((s_n, s_{n+1}], d\nu) \\ &= \hat{y}_{s_n}^{\text{RK}} + \int_{\mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}^{\text{RK}}(\nu)(s_{n+1} - s_n)) \gamma(d\nu),\end{aligned}$$

1115 and for the θ -Trapezoidal scheme that

$$\begin{aligned}\hat{y}_{s_{n+1}}^{\text{trap}} &= \hat{y}_{s_n}^{\text{trap}} + \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \nu N[\hat{\mu}^{\text{trap}}](ds, d\nu) \\ &= \hat{y}_{s_n}^{\text{trap}} + \int_{\mathbb{D}} \nu N[\hat{\mu}^{\text{trap}}]((s_n, s_{n+1}], d\nu) \\ &= \hat{y}_{s_n}^{\text{trap}} + \int_{\mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}^{\text{trap}}(\nu)(s_{n+1} - s_n)) \gamma(d\nu),\end{aligned}$$

1116 and the statement follows by taking $\gamma(d\nu)$ as the counting measure. \square

1117 **Remark C.4** (Remark on Rejection Sampling and Periodicity Assumption). *The rejection sampling*
 1118 *procedure in both algorithms (Algs. 1 and 2) guarantees well-posedness in the rare scenarios where*
 1119 *a large drawn value of Poisson random variables or multiple simultaneous jumps in one coordinate*
 1120 *would result in an update out of the state space $\mathbb{X} = [S]^d$. To enforce this, we simply allow at most*
 1121 *one jump per update across the summation, for example, in the update*

$$\hat{y}_{\rho_n}^* \leftarrow \hat{y}_{s_n} + \sum_{\nu \in \mathbb{D}} \nu \mathcal{P}(\hat{\mu}_{s_n}(\nu) \theta \Delta_n),$$

1122 *as the standard practice in the literature [22, 38]. The indicator function $\mathbf{1}_{\hat{\mu}_{s_n} > 0}$ in Alg. 1 is also*
 1123 *used to ensure that only valid jumps from the current state \hat{y}_{s_n} are considered, while in Alg. 2, this*
 1124 *is implicitly guaranteed by taking the positive part of $\alpha_1 \hat{\mu}_{\rho_n}^* - \alpha_2 \hat{\mu}_{s_n}$, which implies the positivity*
 1125 *of $\alpha_1 \hat{\mu}_{\rho_n}^*$ and thus the validity of the jumps $\hat{y}_{\rho_n}^*$. We point out that the single-jump rule is only*
 1126 *a convenient sufficient condition, one should notice that this condition is not necessary for the*
 1127 *well-posedness of our algorithms, since our setting of the state space \mathbb{X} carries both orderliness*
 1128 *and algebraic structure, and thus one could in principle admit multiple simultaneous jumps without*
 1129 *ambiguity.*

1130 *Over the full inference process, the total probability of rejection is at most $\mathcal{O}(\kappa)$. Below, we give a*
 1131 *brief justification and we refer to Proposition A.14 in [38] for a complete proof of this claim. During*
 1132 *the update aforementioned, the probability of at least two jumps occurring is bounded by*

$$\begin{aligned}\mathbb{P}\left(\sum_{\nu \in \mathbb{D}} \mathcal{P}(\hat{\mu}_{s_n}(\nu) \theta \Delta_n) > 1\right) &= 1 - \mathbb{P}\left(\mathcal{P}\left(\sum_{\nu \in \mathbb{D}} \hat{\mu}_{s_n}(\nu) \theta \Delta_n\right) \leq 1\right) \\ &= 1 - \exp\left(-\sum_{\nu \in \mathbb{D}} \hat{\mu}_{s_n}(\nu) \theta \Delta_n\right) \left(1 + \sum_{\nu \in \mathbb{D}} \hat{\mu}_{s_n}(\nu) \theta \Delta_n\right) \\ &\lesssim \left(\sum_{\nu \in \mathbb{D}} \hat{\mu}_{s_n}(\nu) \theta \Delta_n\right)^2 \lesssim \Delta_n^2.\end{aligned}$$

1133 *Summing $\mathcal{O}(\Delta_n^2)$ over N steps gives $\sum_{n=0}^{N-1} \Delta_n^2 \lesssim \kappa T$, and an identical argument applies to the*
 1134 *second update in each iteration. Hence, the overall rejection rate is at most $\mathcal{O}(\kappa)$.*

1135 *When we impose periodic boundary conditions, $\mathbb{X} = [S]^d$ is equipped with a convenient algebraic*
 1136 *structure: addition and scalar multiplication are globally well-defined. In that case, Algs. 1 and 2*
 1137 *match exactly the stochastic integral formulations in Props. C.2 and C.3. This alignment removes*
 1138 *the need for per-step rejection, streamlines the application of the change-of-measure argument, and*
 1139 *greatly simplifies the convergence proofs of Thms. 5.4 and 5.5. Even without periodicity, those*
 1140 *theorems hold with probability at least $1 - \mathcal{O}(\kappa)$, as shown above.*

1141 C.2 Convergence Analysis of the θ -Trapezoidal Method

1142 **Theorem C.5.** Let $\tilde{p}_{0:T-\delta}$ and $\hat{q}_{0:T-\delta}^{\text{trap}}$ be the path measures of the backward process with the
 1143 stochastic integral formulation (2.4) and the interpolating process (C.5) of the θ -trapezoidal method
 1144 (Alg. 2), then it holds that

$$\begin{aligned} D_{\text{KL}}(\tilde{p}_{T-\delta} \|\hat{q}_{T-\delta}^{\text{trap}}) &\leq D_{\text{KL}}(\tilde{p}_{0:T-\delta} \|\hat{q}_{0:T-\delta}^{\text{trap}}) \\ &\leq D_{\text{KL}}(\tilde{p}_0 \|\hat{q}_0) + \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \right], \end{aligned} \quad (\text{C.7})$$

1145 where the intensity $\hat{\mu}^{\text{trap}}$ is defined in (C.5), and the expectation is taken w.r.t. both paths generated
 1146 by the backward process (2.4) and the randomness of the Poisson random measure used in the first
 1147 step of each iteration of the algorithm, i.e., the construction of the intermediate process (C.1), which
 1148 is assumed to be independent of that of the backward process.

1149 *Proof.* First, we will handle the randomness introduced by the Poisson random measure in the first
 1150 step of each iteration of the θ -trapezoidal method. For the ease of presentation, we encode the
 1151 aforementioned randomness as a random variable ζ and suppose it is still supported on the probability
 1152 space $(\Omega, \mathcal{F}, \mathbb{P})$ while being independent of the backward process. Then for each realization of ζ ,
 1153 the intermediate process \hat{y}_s^* is constructed as in (C.1) and the corresponding intensity $\hat{\mu}_s^*$ is defined
 1154 in (C.4).

1155 Given the stochastic integral formulation of the backward process (2.4) and the interpolating process
 1156 of the θ -trapezoidal method (C.5), we have by Thm. B.4 that this particular realization of the path
 1157 measure $\hat{q}_{0:T-\delta}^{\text{trap}}$ can be obtained by changing the path measure $\tilde{p}_{0:T-\delta}$ with the Radon-Nikodym
 1158 derivative

$$Z_t \left[\frac{\hat{\mu}^{\text{trap}}}{\mu} \right] = \exp \left(- \int_0^t \int_{\mathbb{D}} \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} N[\mu](ds, d\nu) + \int_0^t \int_{\mathbb{D}} (\mu_s(\nu) - \hat{\mu}_s^{\text{trap}}(\nu)) \gamma(d\nu) ds \right),$$

1159 i.e.,

$$\begin{aligned} D_{\text{KL}}(\tilde{p}_{0:T-\delta} \|\hat{q}_{0:T-\delta}^{\text{trap}} | \zeta) &= \mathbb{E} \left[\log Z_{T-\delta}^{-1} \left[\frac{\hat{\mu}^{\text{trap}}}{\mu} \right] \right] \\ &= \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \right]. \end{aligned}$$

1160 Then it is easy to see by the data processing inequality and the chain rule of KL divergence that

$$\begin{aligned} D_{\text{KL}}(\tilde{p}_{T-\delta} \|\hat{q}_{T-\delta}^{\text{trap}}) &\leq D_{\text{KL}}(\tilde{p}_{0:T-\delta} \|\hat{q}_{0:T-\delta}^{\text{trap}}) \leq \mathbb{E} [D_{\text{KL}}(\tilde{p}_{T-\delta} \|\hat{q}_{T-\delta}^{\text{trap}} | \zeta)] \\ &= D_{\text{KL}}(\tilde{p}_0 \|\hat{q}_0) + \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \right], \end{aligned}$$

1161 and the proof is complete. \square

1162 In the following, we will provide the outline of the proof of Thm. 5.4, where we leave the proof of
 1163 several lemmas and detailed calculations to App. C.4 for the clarity of presentation.

1164 *Proof of Thm. 5.4.* Throughout this proof, including the subsequent lemmas and propositions that
 1165 will be detailed in App. C.4, we will assume that $(y_s)_{s \in [0, T]}$ is a process generated by the path
 1166 measure $\tilde{p}_{0:T}$ of the backward process with the stochastic integral formulation (2.4) and set it as the
 1167 underlying paths of the expectation in (C.7) as required by Thm. C.5. Especially, $y_s \sim \tilde{p}_s$ holds for
 1168 any $s \in [0, T]$. For simplicity, we will assume that the process y_s is left-continuous at each grid point
 1169 s_i for $i \in [0 : N]$, which happens with probability one.

1170 We first consider the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, and thus we have $\lfloor s \rfloor = s_n$ and $\rho_s = \rho_n$.
 1171 Within this interval, we will denote its intermediate process as appeared in (C.1) as y_s^* , and the
 1172 corresponding intermediate intensity as appeared in (C.4) as $\hat{\mu}_s^*$. In the following discussion, we will
 1173 assume implicitly that the processes are conditioned on the filtration \mathcal{F}_{s_n} .

1174 By the definition of the intensity $\widehat{\mu}^{\text{trap}}(\nu)$ as specified in (C.6)

$$\widehat{\mu}_s^{\text{trap}} = \mathbf{1}_{s < \rho_s} \widehat{\mu}_{[s]} + \mathbf{1}_{s \geq \rho_s} (\alpha_1 \widehat{\mu}_{\rho_s}^* - \alpha_2 \widehat{\mu}_{[s]})_+,$$

1175 we can rewrite the corresponding part of the integral in (C.7) as

$$\begin{aligned} & \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\widehat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \widehat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \\ &= \left(\int_{s_n}^{\rho_n} + \int_{\rho_n}^{s_{n+1}} \right) \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\widehat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \widehat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \\ &= \underbrace{\int_{s_n}^{\rho_n} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\widehat{\mu}_{s_n}(\nu)} - \mu_s(\nu) + \widehat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds}_{\text{(I)}} \\ &+ \underbrace{\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)} - \mu_s(\nu) + \alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds}_{\text{(II)}}, \end{aligned}$$

1176 where the assumption that $\alpha_1 \widehat{\mu}_{\rho_s}^* - \alpha_2 \widehat{\mu}_{[s]} \geq 0$ for all $s \in [0, T - \delta]$ is applied here for the second
1177 term (II) above.

1178 **Decomposition of the Integral.** Next, we decompose the integral (I) and (II) into several terms,
1179 the magnitudes of which or combinations of which are to be bounded.

1180 (i) The first term is decomposed as

$$(I) = (I.1) + (I.2) + (I.3) + (I.4),$$

1181 where each term is defined as

$$\begin{aligned} (I.1) &= \int_{s_n}^{\rho_n} \int_{\mathbb{D}} \left(\mu_{s_n}(\nu) \log \frac{\mu_{s_n}(\nu)}{\widehat{\mu}_{s_n}(\nu)} - \mu_{s_n}(\nu) + \widehat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds, \\ (I.2) &= \int_{s_n}^{\rho_n} \int_{\mathbb{D}} (\mu_s(\nu) \log \mu_s(\nu) - \mu_s(\nu) - \mu_{s_n}(\nu) \log \mu_{s_n}(\nu) + \mu_{s_n}(\nu)) \gamma(d\nu) ds, \\ (I.3) &= \int_{s_n}^{\rho_n} \int_{\mathbb{D}} (\mu_s(\nu) - \mu_{s_n}(\nu)) (\log (\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) - \log \widehat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds, \\ (I.4) &= \int_{s_n}^{\rho_n} \int_{\mathbb{D}} \mu_{s_n}(\nu) \log (\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \\ &\quad - \int_{s_n}^{\rho_n} \int_{\mathbb{D}} \mu_s(\nu) \log (\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds. \end{aligned}$$

1182 (ii) The second term is decomposed as

$$(II) = (II.1) + (II.2) + (II.3) + (II.4) + (II.5) + (II.6),$$

1183 where each term is defined as

$$\begin{aligned} (II.1) &= \alpha_1 \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_{\rho_n}(\nu) \log \frac{\mu_{\rho_n}(\nu)}{\widehat{\mu}_{\rho_n}(\nu)} - \mu_{\rho_n}(\nu) + \widehat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \\ &\quad - \alpha_2 \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_{s_n}(\nu) \log \frac{\mu_{s_n}(\nu)}{\widehat{\mu}_{s_n}(\nu)} - \mu_{s_n}(\nu) + \widehat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds, \\ (II.2) &= \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\mu_s(\nu) \log \mu_s(\nu) - \mu_s(\nu)) \gamma(d\nu) ds \\ &\quad - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 (\mu_{\rho_n}(\nu) \log \mu_{\rho_n}(\nu) - \mu_{\rho_n}(\nu)) - \alpha_2 (\mu_{s_n}(\nu) \log \mu_{s_n}(\nu) - \mu_{s_n}(\nu))) \gamma(d\nu) ds, \end{aligned}$$

$$(II.3) = \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\hat{\mu}_{\rho_n}^*(\nu) - \hat{\mu}_{\rho_n}(\nu)) \gamma(d\nu) ds,$$

$$(II.4) = \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \\ - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds,$$

$$(II.5) = \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \\ - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds,$$

$$(II.6) = \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \\ - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \mu_s(\nu) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds.$$

1184 **Bounding the Error Terms.** Then we briefly summarize the intuitions and related techniques
1185 used in the bounds of the terms above, and the detailed calculations and proofs of the lemmas and
1186 propositions are deferred to App. C.4.

1187 (i) *Error due to estimation error associated with the intensity:* The terms (I.1) and (II.1) are
1188 bounded by the assumption on the estimation error of the intensity $\hat{\mu}_s$ (Assump. 5.3), as

$$\mathbb{E}[(I.1) + (II.1)] \leq \theta \Delta_n \epsilon_I + \alpha_1 (1 - \theta) \Delta_n \epsilon_I = \theta \Delta_n \epsilon_I + \frac{1}{2\theta} \Delta_n \epsilon_I \lesssim \Delta_n \epsilon_I,$$

1189 for any $\theta \in (0, 1]$.

1190 The term (II.4) is bounded by Prop. C.9, as

$$\mathbb{E}[(II.4)] \lesssim \Delta_n \epsilon_{II},$$

1191 where Jensen's inequality is applied here based on the convexity of the loss.

1192 (ii) *Error related to the smoothness of intensity:* By Cor. C.13, the terms (I.2) and (II.2) are
1193 bounded by

$$\mathbb{E}[(I.2) + (II.2)] \leq \Delta_n^3.$$

1194 By Cor. C.14, the terms (I.4) and (II.6) are bounded by

$$\mathbb{E}[(I.4) + (II.6)] \leq \Delta_n^3.$$

1195 Intuitively, the bounds on these terms closely relate to the properties of the jump process and
1196 quantify the smoothness assumption on the intensity μ_s (Assump. 5.2), especially when the
1197 intensity does not vary significantly within the interval $(s_n, s_{n+1}]$. The main technique used
1198 for bounding these terms is Dynkin's Formula (Thm. C.10). The third-order accuracy here
1199 directly follows from the intuition provided in Sec. 4 based on numerical quadrature.

1200 (iii) *Error involving the intermediate process:* The terms (II.3) and (II.5) are bounded by Prop. C.18
1201 and Cor. C.19 respectively as follows

$$\mathbb{E}[(II.3)] \lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{II}, \quad \text{and} \quad \mathbb{E}[(II.5)] \lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{II},$$

1202 The term (I.3) is bounded by Prop. C.20 as below

$$\mathbb{E}[(I.3)] \lesssim \Delta_n^3.$$

1203 The three terms above all involve the intermediate process y_s^* and the corresponding intermedi-
1204 ate density $\hat{\mu}_s^*$.

1205 In conclusion, by summing up all these terms, we have

$$\begin{aligned} & \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \\ & \lesssim \Delta_n(\epsilon_I + \epsilon_{II}) + \Delta_n^3 + \Delta_n^2 \epsilon_{II} \lesssim \Delta_n(\epsilon_I + \epsilon_{II}) + \Delta_n^3. \end{aligned}$$

1206 Therefore, the overall error is bounded by first applying Thm. C.5 and then the upper bound derived
1207 above to each interval $(s_n, s_{n+1}]$, which yields

$$\begin{aligned} & D_{\text{KL}}(\tilde{p}_{T-\delta} \| \hat{q}_{T-\delta}^{\text{trap}}) \\ & \leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{trap}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{trap}}(\nu) \right) \gamma(d\nu) ds \right] \\ & \lesssim D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \sum_{n=0}^{N-1} (\Delta_n(\epsilon_I + \epsilon_{II}) + \Delta_n^3) \\ & \lesssim \exp(-T) + T(\epsilon_I + \epsilon_{II}) + \kappa^2 T, \end{aligned}$$

1208 as desired. \square

Remark C.6 (Discussion on the Positivity Assumption). *In the following, we will take the positivity assumption in Thm. 5.4 as an example, and the case of the θ -RK-2 method is similar. In the statement of Thm. 5.4, we have assumed that*

$$\alpha_1 \hat{\mu}_{\rho_s}^*(\nu) - \alpha_2 \hat{\mu}_{[s]}(\nu) \geq 0$$

1209 in (C.6) for all $s \in [0, T - \delta]$, which allows us to replace $(\alpha_1 \hat{\mu}_{\rho_s}^*(\nu) - \alpha_2 \hat{\mu}_{[s]}(\nu))_+$ by the difference
1210 itself. [25] showed that this approximation is at most of $\mathcal{O}(\Delta_n^3)$ within the corresponding interval,
1211 and [26] further proved that for any order $p \geq 1$, there exists a sufficiently small step size Δ such
1212 that this approximation is at least p -th order, i.e., of order $\mathcal{O}(\Delta^p)$ for that step.

1213 We give a brief justification of this assumption here. We consider the expectation of the difference
1214 itself, which is given by

$$\begin{aligned} & \mathbb{E} [\alpha_1 \hat{\mu}_{\rho_s}^*(\nu) - \alpha_2 \hat{\mu}_{[s]}(\nu)] = \mathbb{E} [\hat{\mu}_{[s]}(\nu) + \alpha_1 (\hat{\mu}_{\rho_s}^*(\nu) - \hat{\mu}_{\rho_s}(\nu)) + \alpha_1 (\hat{\mu}_{\rho_s}(\nu) - \hat{\mu}_{[s]}(\nu))] \\ & \gtrsim 1 - \alpha_1(\kappa \epsilon_{II} + \kappa) = 1 - \mathcal{O}(\kappa), \end{aligned}$$

1215 where we used $\mathbb{E} [|\hat{\mu}_{\rho_s}^*(\nu) - \hat{\mu}_{\rho_s}(\nu)|] \lesssim \kappa \epsilon_{II}$, as established in Eq. (C.17) and
1216 $\mathbb{E} [|\hat{\mu}_{\rho_s}(\nu) - \hat{\mu}_{[s]}(\nu)|] \lesssim \kappa$, as shown in Eq. (C.18). Therefore, as long as the step sizes Δ_n
1217 are sufficiently small, the positivity assumption is valid in the sense that the expectation of the
1218 difference is at least $1 - \mathcal{O}(\kappa)$.

1219 C.3 Convergence Analysis of the θ -RK-2 Method

1220 Here we may again apply the data processing inequality and the chain rule of KL divergence to upper
1221 bound the error associated with the θ -RK-2 method. A statement of the upper bound is provided
1222 in Thm. C.7 below, whose proof is omitted here since it is similar to that of Thm. C.5 above.

1223 **Theorem C.7.** *Let $\tilde{p}_{0:T-\delta}$ and $\hat{q}_{0:T-\delta}^{\text{RK}}$ be the path measures of the backward process with the
1224 stochastic integral formulation (2.4) and the interpolating process (C.2) of the θ -RK-2 method (
1225 Alg. 1), then it holds that*

$$\begin{aligned} & D_{\text{KL}}(\tilde{p}_{T-\delta} \| \hat{q}_{T-\delta}^{\text{RK}}) \leq D_{\text{KL}}(\tilde{p}_{0:T-\delta} \| \hat{q}_{0:T-\delta}^{\text{RK}}) \\ & \leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{RK}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{RK}}(\nu) \right) \gamma(d\nu) ds \right], \quad (\text{C.8}) \end{aligned}$$

1226 where the intensity $\hat{\mu}^{\text{RK}}$ is defined in (C.2), and the expectation is taken w.r.t. both paths generated
1227 by the backward process (2.4) and the randomness of the Poisson random measure used in the first
1228 step of each iteration of the algorithm, i.e., the construction of the intermediate process (C.1), which
1229 is assumed to be independent of that of the backward process.

Following the same flow as in the proof of Thm. 5.4, we will first provide an outline of the proof of Thm. 5.5, and defer the proof of several key lemmas and detailed calculations to App. C.4 for the clarity of presentation. We will also comment on the differences that may lead to the less desirable numerical properties of the θ -RK-2 method.

Proof of Thm. 5.5. In the following proof sketch, we will be using the same notation as in the proof of Thm. 5.4, and we will assume that the process y_s is left-continuous at each grid point s_i for $i \in [0 : N]$. We also start by taking a closer look at the integral within each interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, and denote the intermediate process as appeared in (C.1) as y_s^* and the corresponding intermediate intensity as appeared in (C.4) as $\hat{\mu}_s^*$.

As defined in (C.3), the intensity $\hat{\mu}^{\text{RK}}(\nu)$ is given by

$$\hat{\mu}_s^{\text{RK}}(\nu) = \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{\lfloor s \rfloor}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_s}^*(\nu),$$

which helps us rewrite the corresponding part of the integral in (C.8) as

$$\begin{aligned} & \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{RK}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{RK}}(\nu) \right) \gamma(d\nu) ds \\ &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \underbrace{\left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu)} - \mu_s(\nu) + \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right)}_{\text{(III)}} \gamma(d\nu) ds. \end{aligned}$$

Above we again use the positivity assumption that $\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{\lfloor s \rfloor} + \frac{1}{2\theta} \hat{\mu}_{\rho_s}^* \geq 0$ for the term (III) above, just as what we have done in the proof and discussion of Thm. 5.4 above.

Decomposition of the Integral. Then we perform a similar decomposition of the integral as in the proof of Thm. 5.4 as follows:

$$\text{(III)} = \text{(III.1)} + \text{(III.2)} + \text{(III.3)} + \text{(III.4)} + \text{(III.5)} + \text{(III.6)},$$

where each term is defined as

$$\begin{aligned} \text{(III.1)} &= \left(1 - \frac{1}{2\theta}\right) \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_{s_n}(\nu) \log \left(\frac{\mu_{s_n}(\nu)}{\hat{\mu}_{s_n}(\nu)} \right) - \mu_{s_n}(\nu) + \hat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds \\ &\quad + \frac{1}{2\theta} \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_{\rho_n}(\nu) \log \left(\frac{\mu_{\rho_n}(\nu)}{\hat{\mu}_{\rho_n}(\nu)} \right) - \mu_{\rho_n}(\nu) + \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds, \\ \text{(III.2)} &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} (\mu_s(\nu) \log \mu_s(\nu) - \mu_s(\nu)) \gamma(d\nu) ds \\ &\quad - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) (\mu_{s_n}(\nu) \log \mu_{s_n}(\nu) - \mu_{s_n}(\nu)) + \frac{1}{2\theta} (\mu_{\rho_n}(\nu) \log \mu_{\rho_n}(\nu) - \mu_{\rho_n}(\nu)) \right) \gamma(d\nu) ds, \\ \text{(III.3)} &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \frac{1}{2\theta} (\hat{\mu}_{\rho_n}^*(\nu) - \hat{\mu}_{\rho_n}(\nu)) \gamma(d\nu) ds, \\ \text{(III.4)} &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \\ &\quad - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds, \\ \text{(III.5)} &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \\ &\quad - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \gamma(d\nu) ds, \\ \text{(III.6)} &= \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \gamma(d\nu) ds \end{aligned}$$

$$- \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \mu_s(\nu) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \gamma(d\nu) ds.$$

1246 **Bounding the Error Terms.** Then we briefly summarize the intuitions and related techniques used
 1247 in the bound of the terms above,. Detailed calculations and proofs of the lemmas and propositions
 1248 used here are deferred to App. C.4.

1249 (i) *Error due to the intensity estimation:* The terms in (III.1) are bounded by the assumption on
 1250 the estimation error of the intensity $\hat{\mu}_s$ (Assump. 5.3) as follows

$$\mathbb{E}[(\text{III.1})] \leq \left(1 - \frac{1}{2\theta}\right) \Delta_n \epsilon_I + \frac{1}{2\theta} \Delta_n \epsilon_I = \Delta_n \epsilon_I,$$

1251 for any $\theta \in (0, 1]$.

1252 (ii) *Error related to the smoothness of intensity:* By Cor. C.16 and Cor. C.17, the terms (III.2) and
 1253 (III.6) are bounded by

$$\mathbb{E}[(\text{III.2})] \leq \Delta_n^3, \quad \text{and} \quad \mathbb{E}[(\text{III.6})] \leq \Delta_n^3,$$

1254 respectively.

1255 (iii) *Error involving the intermediate process:* The term (III.3) and (III.5) are bounded in almost
 1256 the same way as that of Prop. C.18 and Cor. C.19. By simply altering the integral upper limits,
 1257 we obtain that

$$\mathbb{E}[(\text{III.3})] \lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{\text{II}}, \quad \mathbb{E}[(\text{III.5})] \lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{\text{II}}.$$

1258 The only term that cannot be directly bounded based on results in App. C.4 is (III.4), which is given
 1259 by

$$\begin{aligned} \mathbb{E}[(\text{III.4})] &= \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \right. \\ &\quad \left. - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \right] \end{aligned} \quad (\text{C.9})$$

1260 Recall that in the proof of its counterpart (Prop. C.9), we utilized the convexity of the loss function and
 1261 the extrapolation nature of the second step in the θ -trapezoidal method (C.6) to bound the error term.
 1262 However, the same technique cannot be directly applied to the θ -RK-2 method for any $\theta \in [0, 1]$, as
 1263 the intensity $\hat{\mu}_s^{\text{RK}}$ is an interpolation of the intensity $\hat{\mu}_s$ when $\theta \in (\frac{1}{2}, 1]$. Therefore, below we will
 1264 first focus on the case when $\theta \in (0, \frac{1}{2}]$.

1265 To be specific, by the assumption on the estimation error (Assump. 5.3), we can reduce (C.9) to

$$\begin{aligned} \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) \right) \right. \\ \left. - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \right], \end{aligned} \quad (\text{C.10})$$

1266 which can then be upper bounded based on Jensen's inequality and the convexity of the loss function
 1267 for $\theta \in (0, \frac{1}{2}]$.

1268 Summing up the bounds of the terms above, we have

$$\begin{aligned} &\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{RK}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{RK}}(\nu) \right) \gamma(d\nu) ds \\ &\lesssim \Delta_n (\epsilon_I + \epsilon_{\text{II}}) + \Delta_n^3 + \Delta_n^2 \epsilon_{\text{II}} \lesssim \Delta_n (\epsilon_I + \epsilon_{\text{II}}) + \Delta_n^3, \end{aligned}$$

1269 Consequentially, the overall error of the θ -RK-2 method is bounded by

$$\begin{aligned}
& D_{\text{KL}}(\tilde{p}_{T-\delta} \| \hat{q}_{T-\delta}^{\text{RK}}) \\
& \leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \mathbb{E} \left[\int_0^{T-\delta} \int_{\mathbb{D}} \left(\mu_s(\nu) \log \frac{\mu_s(\nu)}{\hat{\mu}_s^{\text{RK}}(\nu)} - \mu_s(\nu) + \hat{\mu}_s^{\text{RK}}(\nu) \right) \gamma(d\nu) ds \right] \\
& \lesssim D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \sum_{n=0}^{N-1} (\Delta_n(\epsilon_{\text{I}} + \epsilon_{\text{II}}) + \Delta_n^3) \\
& \lesssim \exp(-T) + T(\epsilon_{\text{I}} + \epsilon_{\text{II}}) + \kappa^2 T,
\end{aligned}$$

1270 which suggests that the θ -RK-2 is also of second order when $\theta \in (0, \frac{1}{2}]$. For the other case when
1271 $\theta \in (\frac{1}{2}, 1]$, we will provide a brief discussion in the remark below. \square

1272 **Remark C.8** (Discussions on the case when $\theta \in (\frac{1}{2}, 1]$). For $\theta \in (\frac{1}{2}, 1]$, the term (C.10) is positive
1273 and thus not necessarily bounded. One may wonder if, despite being positive, this term is still of
1274 at least second order. However, the answer seems negative. By applying the Dynkin's formula
1275 (Thm. C.10 and Cor. C.11) to $\mu_s \log \hat{\mu}_s$ in the term (III.4), we have that the first integral in (C.9) can
1276 be expanded as follows

$$\begin{aligned}
& \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) \right) \gamma(d\nu) ds \right] \\
& = \frac{1}{2\theta} \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} (\mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) + \theta \Delta_n \mathcal{L}(\mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu))) \gamma(d\nu) ds \\
& + \left(1 - \frac{1}{2\theta}\right) \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) \gamma(d\nu) ds + \mathcal{O}(\Delta_n^2) \\
& = \Delta_n \int_{\mathbb{D}} \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) \gamma(d\nu) + \frac{1}{2} \Delta_n^2 \int_{\mathbb{D}} \mathcal{L}(\mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) + \mathcal{O}(\Delta_n^3).
\end{aligned}$$

1277 Similarly, applying Dynkin's formula to the following function

$$G_s(\nu, y_{s-}) = \left(\frac{1}{2\theta} \mu_s(\nu, y_{s-}) + \left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu, y_{s-}) \right) \log \left(\frac{1}{2\theta} \hat{\mu}_s(\nu, y_{s-}) + \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu, y_{s-}) \right),$$

1278 with $G_0(\nu, y_{s_n}) = \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n})$ allows us to expand the second integral in (C.9) as
1279 below

$$\begin{aligned}
& \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\frac{1}{2\theta} \mu_{\rho_n}(\nu) + \left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) \right) \log \left(\frac{1}{2\theta} \hat{\mu}_{\rho_n}(\nu) + \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) \right) \gamma(d\nu) ds \right] \\
& = \Delta_n \int_{\mathbb{D}} G_{s_n}(y_{s_n}) \gamma(d\nu) + \theta \Delta_n^2 \int_{\mathbb{D}} \mathcal{L} G_{s_n}(y_{s_n}) \gamma(d\nu) + \mathcal{O}(\Delta_n^3),
\end{aligned}$$

1280 where

$$\begin{aligned}
& \mathcal{L}G_{s_n}(\nu, y_{s_n}) \\
&= \frac{1}{2\theta} \partial_s \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) + \frac{1}{2\theta} \mu_{s_n}(\nu, y_{s_n}) \frac{1}{2\theta} \frac{\partial_s \hat{\mu}_{s_n}(\nu, y_{s_n})}{\hat{\mu}_{s_n}(\nu, y_{s_n})} \\
&+ \frac{1}{2\theta} \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n} + \nu') \log \left(\frac{1}{2\theta} \hat{\mu}_s(\nu, y_{s_n} + \nu') + \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu, y_{s_n} + \nu') \right) \gamma(d\nu') \\
&- \frac{1}{2\theta} \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) \gamma(d\nu') \\
&+ \left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu, y_{s_n}) \frac{1}{2\theta} \frac{\partial_s \hat{\mu}_{s_n}(\nu, y_{s_n})}{\hat{\mu}_{s_n}(\nu, y_{s_n})} \\
&+ \left(1 - \frac{1}{2\theta}\right) \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n} + \nu') \log \left(\frac{1}{2\theta} \hat{\mu}_s(\nu, y_{s_n} + \nu') + \left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu, y_{s_n} + \nu') \right) \gamma(d\nu') \\
&- \left(1 - \frac{1}{2\theta}\right) \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) \gamma(d\nu') \\
&= \frac{1}{2\theta} \partial_s \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) + \frac{1}{2\theta} \mu_{s_n}(\nu, y_{s_n}) \frac{\partial_s \hat{\mu}_{s_n}(\nu, y_{s_n})}{\hat{\mu}_{s_n}(\nu, y_{s_n})} \\
&+ \frac{1}{2\theta} \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n} + \nu') \log \hat{\mu}_s(\nu, y_{s_n} + \nu') \gamma(d\nu') \\
&+ \left(1 - \frac{1}{2\theta}\right) \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n} + \nu') \log \hat{\mu}_s(\nu, y_{s_n} + \nu') \gamma(d\nu') \\
&- \frac{1}{2\theta} \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) \gamma(d\nu') - \left(1 - \frac{1}{2\theta}\right) \int_{\mathbb{D}} \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n}) \gamma(d\nu').
\end{aligned}$$

1281 This further implies that

$$\begin{aligned}
\theta \mathcal{L}G_{s_n}(y_{s_n}) &= \frac{1}{2} \mathcal{L}(\mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \\
&+ \frac{1}{2\theta} \int_{\mathbb{D}} (\mu_{s_n}(\nu, y_{s_n} + \nu') \log \hat{\mu}_s(\nu, y_{s_n} + \nu') - \mu_{s_n}(\nu, y_{s_n}) \log \hat{\mu}_{s_n}(\nu, y_{s_n})) \gamma(d\nu').
\end{aligned}$$

1282 Comparing the first and second order terms in the two expansions of the two integrals in (C.9) above
1283 then implies that the term (III.4) is of at most second order.

1284 C.4 Lemmas and Propositions

1285 In this section, we provide the detailed proofs of the lemmas and propositions omitted in the proof
1286 of Thm. 5.4 and Thm. 5.5.

1287 **Error due to the Intensity Estimation.** Apart from the terms (I.1) and (II.1) in the proof
1288 of Thm. 5.4 and the term (III.1) in the proof of Thm. 5.5, we also need to bound the error terms
1289 (II.4) in terms of the intensity estimation error, which is given by the following proposition. Notably,
1290 the following bound also utilizes the convexity of the loss function and the extrapolation nature of the
1291 second step in the θ -trapezoidal method (C.6).

1292 **Proposition C.9.** For the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, we have the following error bound:

$$\begin{aligned}
\mathbb{E}[(\text{II.4})] &= \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right. \\
&\quad \left. - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right] \quad (\text{C.11}) \\
&\lesssim \Delta_n \epsilon_{\text{II}}.
\end{aligned}$$

1293

1294 *Proof.* We first define and bound three error terms (II.4.1), (II.4.2), and (II.4.3) with score estimation error (Assump. 5.3) as follows:

$$\begin{aligned} \mathbb{E}[(\text{II.4.1})] &= \mathbb{E} \left[\left| \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \alpha_1 (\mu_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) - \hat{\mu}_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu)) \gamma(d\nu) ds \right| \right] \\ &\leq \alpha_1 \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} |\mu_{\rho_n}(\nu) - \hat{\mu}_{\rho_n}(\nu)| |\log \hat{\mu}_{\rho_n}(\nu)| \gamma(d\nu) ds \right] \\ &\lesssim \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} |\mu_{\rho_n}(\nu) - \hat{\mu}_{\rho_n}(\nu)| \gamma(d\nu) ds \right] \lesssim \Delta_n \epsilon_{\text{II}}, \end{aligned}$$

1296 Similarly, we also have

$$\mathbb{E}[(\text{II.4.2})] = \mathbb{E} \left[\left| \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} \alpha_2 (\mu_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu) - \hat{\mu}_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right| \right] \lesssim \Delta_n \epsilon_{\text{II}},$$

1297 and

$$\begin{aligned} \mathbb{E}[(\text{II.4.3})] &= \mathbb{E} \left[\left| \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right. \right. \\ &\quad \left. \left. - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right| \right] \\ &\lesssim \Delta_n \epsilon_{\text{II}}. \end{aligned}$$

1298 The remaining term (II.4.4) = (II.4) - (II.4.1) - (II.4.2) - (II.4.3) is then given by

$$\begin{aligned} (\text{II.4.4}) &= \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \hat{\mu}_{\rho_n}(\nu) \log \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu) \log \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \\ &\quad - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \leq 0, \end{aligned}$$

1299 where the last inequality follows from Jensen's inequality, *i.e.*,

$$\alpha_1 x \log x - \alpha_2 y \log y \leq (\alpha_1 x - \alpha_2 y) \log(\alpha_1 x - \alpha_2 y),$$

1300 for $\alpha_1, \alpha_2 \geq 0$ and $\alpha_1 - \alpha_2 = 1$. Therefore, by summing up the terms above, we have

$$\mathbb{E}[(\text{II.4})] \leq \mathbb{E}[(\text{II.4.1}) + (\text{II.4.2}) + (\text{II.4.3}) + (\text{II.4.4})] \lesssim \Delta_n \epsilon_{\text{II}},$$

1301 and the proof is complete. \square

1302 **Error Related to the Smoothness of Intensity.** Below we first present the Dynkin's formula,
1303 which is the most essential tool for the proof of the error related to the smoothness of the intensity.

1304 **Theorem C.10** (Dynkin's Formula). *Let $(y_t)_{t \in [0, \tau]}$ be the following process:*

$$y_t = y_0 + \int_0^t \int_{\mathbb{D}} \nu N[\mu](ds, d\nu),$$

1305 where $N[\mu](ds, d\nu)$ is a Poisson random measure with intensity μ of the form $\mu_s(\nu, y_{s-})$. For any
1306 $f \in C^1([0, \tau] \times \mathbb{X})$, we define the generator of the process $(y_t)_{t \in [0, \tau]}$ as below

$$\mathcal{L}f_t(y) = \lim_{\tau \rightarrow 0^+} \left[\frac{f_{t+\tau}(y_{t+\tau}) - f_t(y_t)}{\tau} \Big|_{y_t = y} \right] = \partial_t f_t(y) + \int_{\mathbb{D}} (f_t(y + \nu) - f_t(y)) \mu_t(\nu, y) \gamma(d\nu). \quad (\text{C.12})$$

1307 Then we have that

$$\mathbb{E}[f_t(y_t)] = f_0(y_0) + \mathbb{E} \left[\int_0^t \mathcal{L}f_s(y_s) ds \right].$$

1308

1309 *Proof.* The definition and the form of the generator \mathcal{L} , as well as the Dynkin's formula are all
 1310 well-known in the literature of jump processes. We refer readers to detailed discussions on these
 1311 topics in [214].

1312 Here we take $X(t) = (t, y_t)$, $z = (\nu, \xi)$, $\alpha(t, X(t)) = 0$, $\sigma(t, X(t)) = 0$, $\gamma(t, X(t^-), z) =$
 1313 $\nu \mathbf{1}_{0 \leq \xi \leq \mu_t(\nu, y_{t-})}$ in the statement of Thm. 1.19 in [214] and replace the compensated Poisson random
 1314 measure $\tilde{N}(dt, dz)$ with the Poisson random measure $N(ds, d\nu, d\xi)$ defined as Rmk. B.3. Then we
 1315 are allowed to use the ordinary Poisson random measure instead of the compensated one since we are
 1316 working with a finite measure $\gamma(d\nu)$.

1317 From Thm. 1.22 in [214], we have that

$$\begin{aligned}\mathcal{L}f_t(y) &= \partial_t f_t(y) + \int_{\mathbb{D}} \int_{\mathbb{R}} (f_t(y + \nu \mathbf{1}_{0 \leq \xi \leq \mu_t(\nu, y)}) - f_t(y)) \gamma(d\nu) d\xi \\ &= \partial_t f_t(y) + \int_{\mathbb{D}} (f_t(y + \nu) - f_t(y)) \mu_t(\nu, y) \gamma(d\nu),\end{aligned}$$

1318 and the proof is complete. \square

1319 In many cases below, we will need the following first-order expansion of the expectation of the
 1320 function $f_t(y_t)$ by assuming the second-order smoothness of the function f .

1321 **Corollary C.11.** *Suppose that the process $(y_t)_{t \in [0, \tau]}$ and the generator \mathcal{L} are defined as in Thm. C.10.*
 1322 *If we further assume that $f \in C^2([0, \tau] \times \mathbb{X})$, then it holds that*

$$\mathbb{E}[f_t(y_t)] = f_0(y_0) + t\mathcal{L}f_0(y_0) + \mathcal{O}(t^2).$$

1323

1324 *Proof.* We expand the function $f_s(y_s)$ from $t = 0$ as follows

$$\begin{aligned}\mathbb{E}[f_t(y_t)] &= f_0(y_0) + \mathbb{E}\left[\int_0^t \mathcal{L}f_s(y_s) ds\right] \\ &= f_0(y_0) + \mathbb{E}\left[\int_0^t \mathcal{L}\left(f_0(y_0) + \int_0^s \mathcal{L}f_\sigma(y_\sigma) d\sigma\right) ds\right] \\ &= f_0(y_0) + \mathcal{L}f_0(y_0)t + \mathbb{E}\left[\int_0^t \int_0^s \mathcal{L}^2 f_\sigma(y_\sigma) d\sigma ds\right],\end{aligned}$$

1325 where \mathcal{L}^2 is the second-order generator of the process $(y_t)_{t \in [0, \tau]}$ defined as follows

$$\begin{aligned}\mathcal{L}^2 f_\sigma(y) &= \mathcal{L}\left(\partial_\sigma f_\sigma(y) + \int_{\mathbb{D}} (f_\sigma(y + \nu) - f_\sigma(y)) \mu_\sigma(\nu) \gamma(d\nu)\right) \\ &= \partial_\sigma^2 f_\sigma(y) + 2 \int_{\mathbb{D}} (\partial_\sigma f_\sigma(y + \nu) - \partial_\sigma f_\sigma(y)) \mu_\sigma(\nu) \gamma(d\nu) \\ &\quad + \int_{\mathbb{D}} (f_\sigma(y + \nu) - f_\sigma(y)) \partial_\sigma \mu_\sigma(\nu) \gamma(d\nu) \\ &\quad + \int_{\mathbb{D}} \int_{\mathbb{D}} (f_\sigma(y + \nu + \nu') - f_\sigma(y + \nu') - f_\sigma(y + \nu) + f_\sigma(y)) \mu_\sigma(\nu) \mu_\sigma(\nu') \gamma(d\nu) \gamma(d\nu'),\end{aligned}$$

1326 which is bounded uniformly by a constant based on the assumption on the smoothness of the function
 1327 f up to the second order and the boundedness of the measure $\gamma(d\nu)$. Therefore, the second-order
 1328 term above is of magnitude $\mathcal{O}(t^2)$, and the proof is complete. \square

1329 The following lemma provides a general recipe for bounding a combination of errors, which resembles
 1330 standard analysis performed for numerical quadratures. In fact, the following lemma can be easily
 1331 proved by Taylor expansion when the process $(y_t)_{t \in [0, \tau]}$ is constant, *i.e.*, $y_t \equiv y$. Cor. C.11 offers an
 1332 analogous approach to perform the expansion when the process $(y_t)_{t \in [0, \tau]}$ is not constant.

1333 **Lemma C.12.** For any function $f \in C^2([0, \tau] \times \mathbb{X})$ and the true backward process $(y_t)_{t \in [0, \tau]}$ defined
 1334 in (2.4), it holds that

$$\left| \mathbb{E} \left[\int_0^{\theta\tau} f_0(y_0) ds + \int_{\theta\tau}^{\tau} (\alpha_1 f_{\theta\tau}(y_{\theta\tau}) - \alpha_2 f_0(y_0)) ds - \int_0^{\tau} f_s(y_s) ds \right] \right| \lesssim \tau^3.$$

1336 *Proof.* Let \mathcal{L} be the generator defined in Thm. C.10. By applying the Dynkin's formula (Thm. C.10
 1337 and Cor. C.11) to the function $f_t(y_t)$ and plugging in the expression of the generator \mathcal{L} , we have that

$$\begin{aligned} & \mathbb{E} \left[\int_0^{\theta\tau} f_0(y_0) ds - \alpha_2 \int_{\theta\tau}^{\tau} f_0(y_0) ds + \alpha_1 \int_{\theta\tau}^{\tau} f_{\theta\tau}(y_{\theta\tau}) ds - \int_0^{\tau} f_s(y_s) ds \right] \\ &= \theta\tau f_0(y_0) - \alpha_2(1 - \theta)\tau f_0(y_0) + \alpha_1(1 - \theta)\tau (f_0(y_0) + \theta\tau \mathcal{L}f_0(y_0)) \\ & \quad - \int_0^{\tau} (f_0(y_0) + s\mathcal{L}f_0(y_0)) ds + \mathcal{O}(\tau^3) \\ &= (\theta - \alpha_2(1 - \theta) + \alpha_1(1 - \theta) - 1)\tau f_0(y_0) + \alpha_1(1 - \theta)\theta\tau^2 \mathcal{L}f_0(y_0) - \frac{\tau^2}{2} \mathcal{L}f_0(y_0) + \mathcal{O}(\tau^3), \end{aligned}$$

1338 which is of the order $\mathcal{O}(\tau^3)$ by noticing that

$$\begin{aligned} \theta - \alpha_2(1 - \theta) + \alpha_1(1 - \theta) - 1 &= \left(\frac{1}{2\theta(1-\theta)} - \frac{\theta^2 + (1-\theta)^2}{2\theta(1-\theta)} \right) (1 - \theta) - (1 - \theta) = 0 \\ \alpha_1(1 - \theta)\theta - \frac{1}{2} &= \frac{1}{2\theta(1-\theta)}(1 - \theta)\theta - \frac{1}{2} = 0, \end{aligned}$$

1339 and the proof is complete. \square

1340 We remark that in Thm. C.10, Cor. C.11, and Lem. C.12, the smoothness of the function f implies
 1341 that its derivatives up to the relevant order are bounded by constants independent of the time step τ .
 1342 This condition is verified in the subsequent proofs.

1343 Then we are ready to bound some of the error terms in the proof of Thm. 5.4 with Lem. C.12.

1344 **Corollary C.13.** For the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, we have the following error bound:

$$\begin{aligned} & |\mathbb{E}[(\text{I.2}) + (\text{II.2})]| \\ &= \left| \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} (\mu_s(\nu) \log \mu_s(\nu) - \mu_s(\nu)) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. - \int_{s_n}^{\rho_n} \int_{\mathbb{D}} (\mu_{s_n}(\nu) \log \mu_{s_n}(\nu) + \mu_{s_n}(\nu)) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1(\mu_{\rho_n}(\nu) \log \mu_{\rho_n}(\nu) - \mu_{\rho_n}(\nu)) - \alpha_2(\mu_{s_n}(\nu) \log \mu_{s_n}(\nu) - \mu_{s_n}(\nu))) \gamma(d\nu) ds \right] \right| \\ & \lesssim \Delta_n^3. \end{aligned}$$

1345

Proof. The bound is obtained by applying Lem. C.12 with f being the function

$$f_s(y_s) = \int_{\mathbb{D}} \mu_s(\nu) \log \mu_s(\nu) \gamma(d\nu),$$

1346 Strictly speaking, $f_s(y_s)$ is actually in the form of $f_s(y_{s-})$, but the argument can be easily extended
 1347 to this case by assuming time continuity of the function f . \square

1348 **Corollary C.14.** For the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, we have the following error bound:

$$\begin{aligned} & |\mathbb{E}[(\text{I.4}) + (\text{II.6})]| \\ &= \left| \mathbb{E} \left[\int_{s_n}^{\rho_n} \int_{\mathbb{D}} \mu_{s_n}(\nu) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. + \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \mu_s(\nu) \log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right] \right| \lesssim \Delta_n^3. \end{aligned}$$

1349

1350 *Proof.* Note that the intermediate process y_s^* defined in (C.1) is driven by a Poisson random mea-
 1351 sure that is independent of the Poisson random measure driving the process y_s within the interval
 1352 $(s_n, s_{n+1}]$. Therefore, the error bound is obtained by

1353 (1) Taking the expectation w.r.t. the intermediate process y_s^* and thus the intermediate intensity
 1354 $\hat{\mu}_s^*$, and

(2) Then applying Lem. C.12 with f being the following function

$$f_s(y_s) = \int_{\mathbb{D}} \mu_s(\nu) \mathbb{E} [\log (\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu))] \gamma(d\nu).$$

1355 The result follows directly. \square

1356 Now we turn to the error term (III.6) in Thm. 5.5, for which we need the following variant
 1357 of Lem. C.12.

1358 **Lemma C.15.** *For any function $f \in C^2([0, \tau] \times \mathbb{X})$ and the true backward process $(y_t)_{t \in [0, \tau]}$ defined*
 1359 *in (2.4), it holds that*

$$\left| \mathbb{E} \left[\int_0^\tau \left(\left(1 - \frac{1}{2\theta}\right) f_0(y_0) + \frac{1}{2\theta} f_{\theta\tau}(y_{\theta\tau}) \right) ds - \int_0^\tau f_s(y_s) ds \right] \right| \lesssim \tau^3.$$

1360

1361 *Proof.* The proof is similar to that of Lem. C.12. Specifically, we let \mathcal{L} be the generator defined
 1362 in Thm. C.10, apply the Dynkin's formula (Thm. C.10 and Cor. C.11) to the function $f_t(y_t)$ and plug
 1363 in the expression of the generator \mathcal{L} , which yields

$$\begin{aligned} & \mathbb{E} \left[\int_0^\tau \left(\left(1 - \frac{1}{2\theta}\right) f_0(y_0) + \frac{1}{2\theta} f_{\theta\tau}(y_{\theta\tau}) \right) ds - \int_0^\tau f_s(y_s) ds \right] \\ &= \left(1 - \frac{1}{2\theta}\right) \tau f_0(y_0) + \frac{1}{2\theta} \int_0^\tau (f_0(y_0) + \theta \tau \mathcal{L} f_0(y_0)) ds - \int_0^\tau (f_0(y_0) + s \mathcal{L} f_0(y_0)) ds + \mathcal{O}(\tau^3) \\ &= \mathcal{O}(\tau^3), \end{aligned}$$

1364 as desired. \square

1365 **Corollary C.16.** *For the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, we have the following error bound:*

$$\begin{aligned} & |\mathbb{E}[(\text{III.2})]| \\ &= \left| \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} (\mu_s(\nu) \log \mu_s(\nu) - \mu_s(\nu)) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) (\mu_{s_n}(\nu) \log \mu_{s_n}(\nu) - \mu_{s_n}(\nu)) + \frac{1}{2\theta} (\mu_{\rho_n}(\nu) \log \mu_{\rho_n}(\nu) - \mu_{\rho_n}(\nu)) \right) \gamma(d\nu) ds \right] \right| \\ &\lesssim \Delta_n^3. \end{aligned}$$

1366

1367 *Proof.* By applying Lem. C.15 with f being the function

$$f_s(y_s) = \int_{\mathbb{D}} \mu_s(\nu) \log \mu_s(\nu) \gamma(d\nu),$$

1368 we have that the result follows directly. \square

1369 **Corollary C.17.** *For any $n \in [0 : N - 1]$ and the corresponding interval $(s_n, s_{n+1}]$, we have the*
 1370 *following error bound:*

$$\begin{aligned} & |\mathbb{E}[(\text{III.6})]| \\ &= \left| \mathbb{E} \left[\int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \left(\left(1 - \frac{1}{2\theta}\right) \mu_{s_n}(\nu) + \frac{1}{2\theta} \mu_{\rho_n}(\nu) \right) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \gamma(d\nu) ds \right. \right. \\ & \quad \left. \left. - \int_{s_n}^{s_{n+1}} \int_{\mathbb{D}} \mu_s(\nu) \log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \gamma(d\nu) ds \right] \right| \lesssim \Delta_n^3. \end{aligned}$$

1371

1372 *Proof.* Following the arguments in the proof of Cor. C.14, the error bound is obtained by first taking
 1373 the expectation w.r.t. the intermediate process y_s^* and thus the intermediate intensity $\hat{\mu}_s^*$, and then
 1374 applying Lem. C.15 with f being the function

$$f_s(y_s) = \int_{\mathbb{D}} \mu_s(\nu) \mathbb{E} \left[\log \left(\left(1 - \frac{1}{2\theta}\right) \hat{\mu}_{s_n}(\nu) + \frac{1}{2\theta} \hat{\mu}_{\rho_n}^*(\nu) \right) \right] \gamma(d\nu),$$

1375 as desired. \square

1376 **Error involving the Intermediate Process.**

1377 **Proposition C.18.** *For the interval $(s_n, s_{n+1}]$ with $n \in [0 : N - 1]$, we have the following error*
 1378 *bound:*

$$\mathbb{E}[(\text{II.3})] = \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\hat{\mu}_{\rho_n}^*(\nu) - \hat{\mu}_{\rho_n}(\nu)) \gamma(d\nu) ds \right] \lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{\text{II}}.$$

1379

1380 *Proof.* First, we rewrite the error term (II.3) as

$$\begin{aligned} \mathbb{E}[(\text{II.3})] &= \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\hat{\mu}_{\rho_n}^*(\nu) - \hat{\mu}_{\rho_n}(\nu)) \gamma(d\nu) ds \right] \\ &\lesssim \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\mathbb{E}[\hat{\mu}_{\rho_n}^*(\nu)] - \mathbb{E}[\hat{\mu}_{\rho_n}(\nu)]) \gamma(d\nu) ds. \end{aligned} \quad (\text{C.13})$$

1381 Then we expand the integrand by applying the Dynkin's formula (Thm. C.10 and Cor. C.11) to the
 1382 function $\hat{\mu}_s(\nu)$ w.r.t. the intermediate process $(y_s^*)_{s \in [s_n, \rho_n]}$ and the process $(y_s)_{s \in [s_n, \rho_n]}$ respectively
 1383 as follows

$$\begin{aligned} &\mathbb{E}[\hat{\mu}_{\rho_n}^*(\nu)] - \mathbb{E}[\hat{\mu}_{\rho_n}(\nu)] \\ &= \mathbb{E}[\hat{\mu}_{s_n}(\nu) + \mathcal{L}^* \hat{\mu}_{s_n}(\nu) \Delta_n + \mathcal{O}(\Delta_n^2)] - \mathbb{E}[\hat{\mu}_{s_n}(\nu) + \mathcal{L} \hat{\mu}_{s_n}(\nu) \Delta_n + \mathcal{O}(\Delta_n^2)] \\ &= \mathbb{E}[(\mathcal{L}^* - \mathcal{L}) \hat{\mu}_{s_n}(\nu) \Delta_n] + \mathcal{O}(\Delta_n^2), \end{aligned}$$

1384 where the generators \mathcal{L}^* and \mathcal{L} are defined as in (C.12) w.r.t. the processes $(y_s^*)_{s \in [s_n, \rho_n]}$ and
 1385 $(y_s)_{s \in [s_n, \rho_n]}$, respectively, i.e., for any function $f \in C^1([s_n, \rho_n] \times \mathbb{X})$, we have

$$\begin{aligned} \mathcal{L}^* f_s(y) &= \partial_s f_s(y) + \int_{\mathbb{D}} (f_s(y + \nu) - f_s(y)) \hat{\mu}_{s_n}(\nu) \gamma(d\nu), \\ \mathcal{L} f_s(y) &= \partial_s f_s(y) + \int_{\mathbb{D}} (f_s(y + \nu) - f_s(y)) \mu_s(\nu) \gamma(d\nu). \end{aligned} \quad (\text{C.14})$$

1386 Therefore, for the term $\mathbb{E}[(\mathcal{L}^* - \mathcal{L}) \hat{\mu}_{s_n}(\nu)]$ evaluated at $s = s_n$, we have

$$\begin{aligned} \mathbb{E}[(\mathcal{L}^* - \mathcal{L}) \hat{\mu}_{s_n}(\nu)] &= \mathbb{E} \left[\int_{\mathbb{D}} (\hat{\mu}_{s_n}(y + \nu) - \hat{\mu}_{s_n}(y)) (\hat{\mu}_{s_n}(\nu) - \mu_{s_n}(\nu)) \gamma(d\nu) \right] \\ &\lesssim \mathbb{E} \left[\int_{\mathbb{D}} |\hat{\mu}_{s_n}(\nu) - \mu_{s_n}(\nu)| \gamma(d\nu) \right] \lesssim \epsilon_{\text{II}}, \end{aligned} \quad (\text{C.15})$$

1387 where we used the assumption on the estimation error (Assump. 5.3) in the last inequality. Then we
 1388 can further reduce (C.13) to

$$\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\mathbb{E}[\hat{\mu}_{\rho_n}^*(\nu)] - \mathbb{E}[\hat{\mu}_{\rho_n}(\nu)]) \gamma(d\nu) ds \lesssim \int_{\rho_n}^{s_{n+1}} (\epsilon_{\text{II}} \Delta_n + \mathcal{O}(\Delta_n^2)) ds \lesssim \epsilon_{\text{II}} \Delta_n^2 + \Delta_n^3,$$

1389 and the proof is complete. \square

1390 **Corollary C.19.** *For the interval $(s_n, s_{n+1}]$ for $n \in [0 : N - 1]$, we have the following error bound:*

$$\begin{aligned} \mathbb{E}[(\text{II.5})] &= \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log(\alpha_1 \hat{\mu}_{\rho_n}(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right. \\ &\quad \left. - \int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} (\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)) \log(\alpha_1 \hat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \hat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right] \\ &\lesssim \Delta_n^3 + \Delta_n^2 \epsilon_{\text{II}}. \end{aligned}$$

1391

1392 *Proof.* Since the two integrands in (II.5) only differ by replacing $\widehat{\mu}_{\rho_n}^*(\nu)$ with $\widehat{\mu}_{\rho_n}(\nu)$, we have the
 1393 following upper bound by using the assumption on the boundedness of the intensities (Assump. 5.2
 1394 (II))

$$\begin{aligned} & \mathbb{E}[(\text{II.5})] \\ & \lesssim \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} |\alpha_1 \mu_{\rho_n}(\nu) - \alpha_2 \mu_{s_n}(\nu)| \frac{1}{\alpha_1 \widehat{\mu}_{\rho_n}(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)} \alpha_1 |\widehat{\mu}_{\rho_n}(\nu) - \widehat{\mu}_{\rho_n}^*(\nu)| \gamma(d\nu) ds \right] \\ & \lesssim \mathbb{E} \left[\int_{\rho_n}^{s_{n+1}} \int_{\mathbb{D}} |\widehat{\mu}_{\rho_n}(\nu) - \widehat{\mu}_{\rho_n}^*(\nu)| \gamma(d\nu) ds \right] \\ & \lesssim \Delta_n \mathbb{E} \left[\int_{\mathbb{D}} |\widehat{\mu}_{\rho_n}(\nu) - \widehat{\mu}_{\rho_n}^*(\nu)| \gamma(d\nu) \right] \end{aligned} \quad (\text{C.16})$$

1395 Applying the same arguments as in Prop. C.18, which uses the generators \mathcal{L} and \mathcal{L}^* defined in (C.14),
 1396 we can bound the RHS above as follows

$$\begin{aligned} & \mathbb{E} [|\widehat{\mu}_{\rho_n}^*(\nu) - \widehat{\mu}_{\rho_n}(\nu)|] \\ & = \mathbb{E} [|(\widehat{\mu}_{s_n}(\nu) + \mathcal{L}^* \widehat{\mu}_{s_n}(\nu) \Delta_n + \mathcal{O}(\Delta_n^2)) - (\widehat{\mu}_{s_n}(\nu) + \mathcal{L} \widehat{\mu}_{s_n}(\nu) \Delta_n + \mathcal{O}(\Delta_n^2))|] \quad (\text{C.17}) \\ & \lesssim \Delta_n \mathbb{E} [|(\mathcal{L}^* - \mathcal{L}) \widehat{\mu}_{s_n}(\nu)|] + \mathcal{O}(\Delta_n^2) \lesssim \Delta_n \epsilon_{\text{II}} + \mathcal{O}(\Delta_n^2) \end{aligned}$$

1397 where the last inequality follows from (C.15). Substituting (C.17) into (C.16) then yields the desired
 1398 upper bound. \square

1399 **Proposition C.20.** *For the interval $(s_n, s_{n+1}]$ with $n \in [0 : N - 1]$, we have the following error*
 1400 *bound:*

$$\begin{aligned} \mathbb{E}[(\text{I.3})] & = \mathbb{E} \left[\int_{s_n}^{\rho_n} \int_{\mathbb{D}} (\mu_s(\nu) - \mu_{s_n}(\nu)) (\log(\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) - \log \widehat{\mu}_{s_n}(\nu)) \gamma(d\nu) ds \right] \\ & \lesssim \Delta_n^3. \end{aligned}$$

1401

1402 *Proof.* First, we observe by Dynkin's formula (Thm. C.10) that

$$\mathbb{E} [|\mu_s(\nu) - \mu_{s_n}(\nu)|] = \mathbb{E} \left[\left| \int_{s_n}^s \mathcal{L} \mu_{s_n} ds + \mathcal{O}(\Delta_n^2) \right| \right] \lesssim \Delta_n,$$

1403 and also

$$\mathbb{E} [|\widehat{\mu}_s(\nu) - \widehat{\mu}_{s_n}(\nu)|] = \mathbb{E} \left[\left| \int_{s_n}^s \mathcal{L}^* \widehat{\mu}_{s_n} ds + \mathcal{O}(\Delta_n^2) \right| \right] \lesssim \Delta_n. \quad (\text{C.18})$$

1404 Secondly, applying the given assumption (Assump. 5.2 (II)) on the boundedness of the intensities
 1405 yields

$$\begin{aligned} & \mathbb{E} [|\log(\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) - \log \widehat{\mu}_{s_n}(\nu)|] \\ & \lesssim \frac{1}{\widehat{\mu}_{s_n}(\nu)} \mathbb{E} [|\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu) - \widehat{\mu}_{s_n}(\nu)|] \\ & \lesssim \mathbb{E} [|\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu) - \widehat{\mu}_{s_n}(\nu)|] \quad (\text{C.19}) \\ & \lesssim \mathbb{E} [\alpha_1 |\widehat{\mu}_{\rho_n}^*(\nu) - \widehat{\mu}_{s_n}(\nu)|] \\ & \lesssim \mathbb{E} [|\widehat{\mu}_{\rho_n}^*(\nu) - \widehat{\mu}_{\rho_n}(\nu)|] + \mathbb{E} [|\widehat{\mu}_{\rho_n}(\nu) - \widehat{\mu}_{s_n}(\nu)|] \\ & \lesssim \Delta_n + \Delta_n \epsilon_{\text{II}} + \mathcal{O}(\Delta_n^2) \lesssim \Delta_n \end{aligned}$$

1406 where the last inequality follows from (C.17) proved above. Therefore, we may further deduce that

$$\begin{aligned} & \mathbb{E}[(\text{I.3})] \\ & \leq \int_{s_n}^{\rho_n} \int_{\mathbb{D}} \mathbb{E} [|\mu_s(\nu) - \mu_{s_n}(\nu)|] \\ & \quad \mathbb{E} [|\log(\alpha_1 \widehat{\mu}_{\rho_n}^*(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu)) - \log(\alpha_1 \widehat{\mu}_{\rho_n}(\nu) - \alpha_2 \widehat{\mu}_{s_n}(\nu))|] \gamma(d\nu) ds \\ & \lesssim \Delta_n^3, \end{aligned}$$

1407 where the first inequality is due to the independency of y_s and y_s^* for $s \in [s_n, \rho_n]$, and the proof is
 1408 complete. \square

D Details of Numerical Experiments

In Apps. D.1 to D.3, we present additional numerical results for the 15-dimensional toy model, text generation, and image generation, respectively.

D.1 15-Dimensional Toy Model

We first derive the closed-form formula of the marginal distributions \mathbf{p}_t in this model. Recall that the state space $\mathbb{X} = \{1, 2, \dots, d\}$ with $d = 15$, and the initial distribution is $\mathbf{p}_0 \in \Delta^d$. The rate matrix at any time is $\mathbf{Q} = \frac{1}{d}\mathbf{E} - \mathbf{I}$. By solving (2.1), we see that

$$\mathbf{p}_t = e^{t\mathbf{Q}}\mathbf{p}_0 = \left(\frac{1 - e^{-t}}{d}\mathbf{E} + e^{-t}\mathbf{I} \right) \mathbf{p}_0,$$

and therefore \mathbf{p}_t converges to the uniform distribution $\mathbf{p}_\infty = \frac{1}{d}\mathbf{1}$ as $t \rightarrow \infty$. The formula of \mathbf{p}_t directly yields the scores $s_t(x) = \frac{p_t(x)}{p_t(x)}$.

During inference, we initialize at the uniform distribution $\mathbf{q}_0 = \mathbf{p}_\infty$ and run from time 0 to $T = 12$. The truncation error of this choice of time horizon is of the magnitude of 10^{-12} reflected by $D_{\text{KL}}(\mathbf{p}_T \parallel \mathbf{p}_\infty)$, and therefore negligible. The discrete time points form an arithmetic sequence.

We generate 10^6 samples for each algorithm and use `np.bincount` to obtain the empirical distribution $\hat{\mathbf{q}}_T$ as the output distribution. Finally, the KL divergence is computed by

$$D_{\text{KL}}(\mathbf{p}_0 \parallel \hat{\mathbf{q}}_T) = \sum_{i=1}^d p_0(i) \log \frac{p_0(i)}{\hat{q}_T(i)}.$$

We also perform bootstrapping for 1000 times to obtain the 95% confidence interval of the KL divergence, the results are shown by the shaded area in Fig. 3. The fitted lines are obtained by standard linear regression on the log-log scale with the slopes marked beside each line in Fig. 3.

D.2 Text Generation

For text generation, we use the small version of RADD [33] checkpoint¹ trained with λ -DCE loss. We choose an early stopping time $\delta = 10^{-3}$ for a stable numerical simulation. Since RADD is a masked discrete diffusion model, we can freely choose the noise schedule $\sigma(t)$ used in the inference process. We consider the following log-linear noise schedule used in the model training,

$$\sigma(t) = \frac{1 - \epsilon}{1 - (1 - \epsilon)t}, \quad \bar{\sigma}(t) = \int_0^t \sigma(s)ds = -\log(1 - (1 - \epsilon)t) \quad (\text{D.1})$$

where we choose $\epsilon = 10^{-3}$.

The score function $s_\theta(\mathbf{x}_t, t)$ used for computing the transition rate matrix can be computed from the RADD score model \mathbf{p}_θ using the following formula from [33],

$$s_t^\theta(\mathbf{x}_t) = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \mathbf{p}_\theta(\mathbf{x}_t), \quad (\text{D.2})$$

where the model \mathbf{p}_θ is trained to approximate the conditional distribution of the masked positions given all unmasked positions. More specifically, let d be the length of the sequence and $\{1, 2, \dots, S\}$ be the vocabulary set (not including the mask token). Then given a partially masked sequence $\mathbf{x} = (x^1, \dots, x^d)$, the model $\mathbf{p}_\theta(\mathbf{x})$ outputs a $d \times S$ matrix whose (ℓ, s) element approximates $\mathbb{P}_{\mathbf{X} \sim \mathbf{p}_{\text{data}}}(\mathbf{x}^\ell = s \mid \mathbf{X}^{\text{UM}} = \mathbf{x}^{\text{UM}})$ when x^ℓ is mask, and is $\mathbf{1}_{X^\ell, s}$ if otherwise. Here, \mathbf{x}^{UM} represents the unmasked portion of the sequence \mathbf{x} .

We adopt a uniform discretization of the time interval $(\delta, 1]$. For θ -RK-2 and θ -Trapezoidal, we pick $\theta = \frac{1}{2}$. We compare our proposed θ -RK-2 and θ -Trapezoidal with the Euler method, Tweedie τ -leaping, τ -leaping, and we present full results across all NFEs ranging from 16 to 1024 in Tab. 2. For each method, we generate 1024 samples with it and compute the averaged perplexities. All the experiments are run on a single NVIDIA A100 GPU.

¹<https://huggingface.co/JingyangOu/radd-lambda-dce>

Table 2: Generative perplexity of texts generated by different sampling algorithms. Lower values are better, with the best in **bold**.

Method	NFE = 16	NFE = 32	NFE = 64	NFE = 128
FHS	≤ 307.425	≤ 186.594	≤ 141.625	≤ 122.732
Euler	≤ 277.962	≤ 160.586	≤ 111.597	≤ 86.276
Tweedie τ -leaping	≤ 277.133	≤ 160.248	≤ 110.848	≤ 85.738
τ -leaping	≤ 126.835	≤ 96.321	≤ 69.226	≤ 52.366
θ -RK-2	≤ 127.363	≤ 109.351	≤ 86.102	≤ 64.317
θ -Trapezoidal	$\leq \mathbf{123.585}$	$\leq \mathbf{89.912}$	$\leq \mathbf{66.549}$	$\leq \mathbf{49.051}$

Method	NFE = 256	NFE = 512	NFE = 1024
FHS	≤ 113.310	≤ 113.026	≤ 109.406
Euler	≤ 68.092	≤ 55.622	≤ 44.686
Tweedie τ -leaping	≤ 70.102	≤ 55.194	≤ 44.257
τ -leaping	≤ 41.694	≤ 33.789	≤ 28.797
θ -RK-2	≤ 49.816	≤ 40.375	≤ 33.971
θ -Trapezoidal	$\leq \mathbf{39.959}$	$\leq \mathbf{32.456}$	$\leq \mathbf{27.553}$

From the table, we observe that θ -Trapezoidal consistently outperforms all other approaches and generates samplers with better perplexities across all NFEs. We also noticed that both the Euler method and Tweedie τ -leaping share a similar performance, which is beaten by a large margin by θ -RK-2 and τ -leaping.

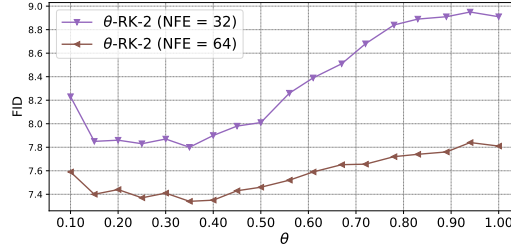


Figure 6: Sampling quality vs. $\theta \in (0, 1]$ in θ -RK-2 algorithm. Sampling quality is quantified through FID.

In Fig. 6, we present the performance of θ -RK-2 with respect to different choices of θ at NFE 32 and 64. We observe that the performance of θ -RK-2 has a flat landscape around the optimal θ choices, which fall in the range $[0.15, 0.4]$. In general, as is evident from the curve, the method performs better when using extrapolation to compute the transition rate matrix, which once again certifies the correctness of our theoretical results (Thm. 5.5) and discussions therebelow.

Table 3: Percentage of positive extrapolated intensities for different algorithms across NFE values.

Method	NFE = 32	NFE = 64	NFE = 128	NFE = 256	NFE = 512	NFE = 1024
θ -RK-2	97.21 ± 3.1	98.31 ± 2.0	98.01 ± 1.3	99.27 ± 0.9	99.44 ± 0.7	99.52 ± 0.6
θ -Trapezoidal	95.67 ± 4.8	97.06 ± 3.6	98.22 ± 2.4	98.87 ± 1.6	99.24 ± 1.1	99.43 ± 0.9

In Tab. 3, we present the percentage of positive extrapolated intensities for different algorithms across NFE values. This partially validates the assumption in our theoretical analysis (Thms. 5.4 and 5.5) that the intensity remains positive throughout the sampling process.

D.3 Image Generation

For the image generation, we use the checkpoint of MaskGIT [62, 65] reproduced in Pytorch². Recall that the MaskGIT is a masked image model which, given a partially masked sequence, outputs the conditional distributions of the masked positions given the unmasked portion, just like the model $p_\theta(\cdot)$ in the aforementioned masked text model, RADD. Therefore, by similarly introducing a time noise schedule $\sigma(t)$ (for which we adopt the same log-linear schedule (D.1) in our experiment), we obtain a masked discrete diffusion model akin to the RADD. The score function can be computed accordingly using the model output as in (D.2).

We choose an early stopping time $\delta = 10^{-3}$, and adopt a uniform discretization of the time interval $(\delta, 1]$ for θ -RK-2, θ -Trapezoidal, τ -leaping and the Euler method. For parallel decoding, we use a linear randomization strategy in the re-masking step and an arccos masking scheduler, the same as the recommended practice in [62]. For each method, we generate 50k samples in a class-conditioned way and compute its FID against the validation split of ImageNet. We use classifier-free guidance to enhance the generation quality and choose the guidance strength to be $w = 3$.

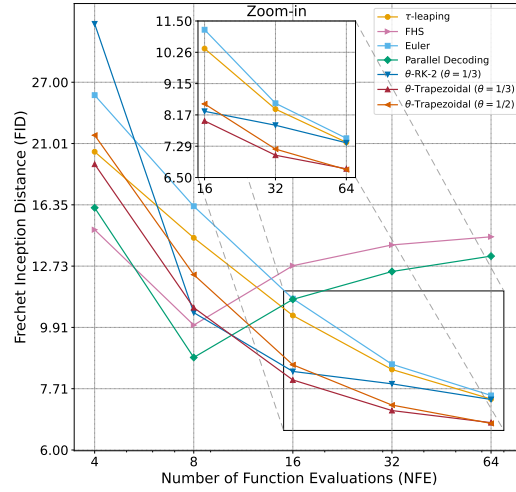


Figure 7: FID of images generated by sampling algorithms vs. number of function evaluations (NFE) with different parameter choices. Lower values are better.

We present the full results for NFE ranging from 4 to 64 in Fig. 7. All the experiments are run on 1 NVIDIA A100. Notably, θ -Trapezoidal with $\theta = \frac{1}{3}$ is the best-performing method except for extremely low NFE budgets. While θ -Trapezoidal with $\theta = \frac{1}{2}$ in general demonstrates a less competitive performance, it converges to the same generation quality as $\theta = \frac{1}{3}$ in the high NFE regime. We also noticed that when using extrapolation with $\theta = \frac{1}{3}$, θ -RK-2 beats τ -leaping for NFE larger than 8, which again accords with our theoretical prediction of its competitive performance in $\theta \in (0, \frac{1}{2}]$ regime.

To investigate the robustness of θ -RK-2 with respect to the choice of θ , we also benchmark its performance across multiple choices at NFE 32 and 64, and we present the results in Fig. 6. Again, similar to the behavior of θ -Trapezoidal, the performance of θ -RK-2 has a flat landscape around the optimal θ choices, which typically falls in the range $[0.3, 0.5]$. In general, as is evident from the curve, the method performs better when using extrapolation to compute the transition rate matrix, which once again certifies the correctness of our theoretical results.

Finally, we visualize some images generated with θ -Trapezoidal on 6 different classes in Fig. 8. θ -Trapezoidal consistently generates high-fidelity images that are visually similar to the ground truth ones and well aligned with the concept.

²<https://github.com/valeoai/Maskgit-pytorch>



Figure 8: Visualization of samples generated by θ -Trapezoidal. **Upper Left:** Aircraft carrier (ImageNet-1k class: **933**); **Upper Middle:** Pirate (ImageNet-1k class: **724**); **Upper Right:** Volcano (ImageNet-1k class: **980**); **Lower Left:** Ostrich (ImageNet-1k class: **009**); **Lower Middle:** Cheeseburger (ImageNet-1k class: **933**); **Lower Right:** Beer bottle (ImageNet-1k class: **440**).

1482 NeurIPS Paper Checklist

1483 The checklist is designed to encourage best practices for responsible machine learning research,
1484 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
1485 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
1486 follow the references and follow the (optional) supplemental material. The checklist does NOT count
1487 towards the page limit.

1488 Please read the checklist guidelines carefully for information on how to answer these questions. For
1489 each question in the checklist:

- 1490 • You should answer [Yes], [No], or [NA].
- 1491 • [NA] means either that the question is Not Applicable for that particular paper or the
1492 relevant information is Not Available.
- 1493 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

1494 **The checklist answers are an integral part of your paper submission.** They are visible to the
1495 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
1496 (after eventual revisions) with the final version of your paper, and its final version will be published
1497 with the paper.

1498 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
1499 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
1500 proper justification is given (e.g., "error bars are not reported because it would be too computationally
1501 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
1502 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
1503 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
1504 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
1505 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
1506 please point to the section(s) where related material for the question can be found.

1507 IMPORTANT, please:

- 1508 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 1509 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 1510 • **Do not modify the questions and only use the provided macros for your answers.**

1511 1. Claims

1512 Question: Do the main claims made in the abstract and introduction accurately reflect the
1513 paper’s contributions and scope?

1514 Answer: [Yes]

1515 Justification: The abstract and introduction accurately reflect the paper’s contributions and
1516 scope.

1517 Guidelines:

- 1518 • The answer NA means that the abstract and introduction do not include the claims
1519 made in the paper.
- 1520 • The abstract and/or introduction should clearly state the claims made, including the
1521 contributions made in the paper and important assumptions and limitations. A No or
1522 NA answer to this question will not be perceived well by the reviewers.
- 1523 • The claims made should match theoretical and experimental results, and reflect how
1524 much the results can be expected to generalize to other settings.
- 1525 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1526 are not attained by the paper.

1527 2. Limitations

1528 Question: Does the paper discuss the limitations of the work performed by the authors?

1529 Answer: [Yes]

1530 Justification: The paper discusses the limitations of the work performed by the authors.

1531 Guidelines:

- 1532 • The answer NA means that the paper has no limitation while the answer No means that
- 1533 the paper has limitations, but those are not discussed in the paper.
- 1534 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1535 • The paper should point out any strong assumptions and how robust the results are to
- 1536 violations of these assumptions (e.g., independence assumptions, noiseless settings,
- 1537 model well-specification, asymptotic approximations only holding locally). The authors
- 1538 should reflect on how these assumptions might be violated in practice and what the
- 1539 implications would be.
- 1540 • The authors should reflect on the scope of the claims made, e.g., if the approach was
- 1541 only tested on a few datasets or with a few runs. In general, empirical results often
- 1542 depend on implicit assumptions, which should be articulated.
- 1543 • The authors should reflect on the factors that influence the performance of the approach.
- 1544 For example, a facial recognition algorithm may perform poorly when image resolution
- 1545 is low or images are taken in low lighting. Or a speech-to-text system might not be
- 1546 used reliably to provide closed captions for online lectures because it fails to handle
- 1547 technical jargon.
- 1548 • The authors should discuss the computational efficiency of the proposed algorithms
- 1549 and how they scale with dataset size.
- 1550 • If applicable, the authors should discuss possible limitations of their approach to
- 1551 address problems of privacy and fairness.
- 1552 • While the authors might fear that complete honesty about limitations might be used by
- 1553 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
- 1554 limitations that aren't acknowledged in the paper. The authors should use their best
- 1555 judgment and recognize that individual actions in favor of transparency play an impor-
- 1556 tant role in developing norms that preserve the integrity of the community. Reviewers
- 1557 will be specifically instructed to not penalize honesty concerning limitations.

1558 3. Theory assumptions and proofs

1559 Question: For each theoretical result, does the paper provide the full set of assumptions and

1560 a complete (and correct) proof?

1561 Answer: [\[Yes\]](#)

1562 Justification: The paper provides the full set of assumptions and a complete (and correct)

1563 proof.

1564 Guidelines:

- 1565 • The answer NA means that the paper does not include theoretical results.
- 1566 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 1567 referenced.
- 1568 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1569 • The proofs can either appear in the main paper or the supplemental material, but if
- 1570 they appear in the supplemental material, the authors are encouraged to provide a short
- 1571 proof sketch to provide intuition.
- 1572 • Inversely, any informal proof provided in the core of the paper should be complemented
- 1573 by formal proofs provided in appendix or supplemental material.
- 1574 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1575 4. Experimental result reproducibility

1576 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

1577 perimental results of the paper to the extent that it affects the main claims and/or conclusions

1578 of the paper (regardless of whether the code and data are provided or not)?

1579 Answer: [\[Yes\]](#)

1580 Justification: The paper fully discloses all the information needed to reproduce the main

1581 experimental results of the paper.

1582 Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The paper will open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper will specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The paper will report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper provides sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: There is no societal impact of the work performed.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

1742 Answer: [NA]

1743 Justification: There is no such risk.

1744 Guidelines:

- 1745 • The answer NA means that the paper poses no such risks.
- 1746 • Released models that have a high risk for misuse or dual-use should be released with
- 1747 necessary safeguards to allow for controlled use of the model, for example by requiring
- 1748 that users adhere to usage guidelines or restrictions to access the model or implementing
- 1749 safety filters.
- 1750 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 1751 should describe how they avoided releasing unsafe images.
- 1752 • We recognize that providing effective safeguards is challenging, and many papers do
- 1753 not require this, but we encourage authors to take this into account and make a best
- 1754 faith effort.

1755 **12. Licenses for existing assets**

1756 Question: Are the creators or original owners of assets (e.g., code, data, models), used in

1757 the paper, properly credited and are the license and terms of use explicitly mentioned and

1758 properly respected?

1759 Answer: [Yes]

1760 Justification: The creators or original owners of assets, used in the paper, are properly

1761 credited and the license and terms of use are explicitly mentioned and properly respected.

1762 Guidelines:

- 1763 • The answer NA means that the paper does not use existing assets.
- 1764 • The authors should cite the original paper that produced the code package or dataset.
- 1765 • The authors should state which version of the asset is used and, if possible, include a
- 1766 URL.
- 1767 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1768 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 1769 service of that source should be provided.
- 1770 • If assets are released, the license, copyright information, and terms of use in the
- 1771 package should be provided. For popular datasets, `paperswithcode.com/datasets`
- 1772 has curated licenses for some datasets. Their licensing guide can help determine the
- 1773 license of a dataset.
- 1774 • For existing datasets that are re-packaged, both the original license and the license of
- 1775 the derived asset (if it has changed) should be provided.
- 1776 • If this information is not available online, the authors are encouraged to reach out to
- 1777 the asset's creators.

1778 **13. New assets**

1779 Question: Are new assets introduced in the paper well documented and is the documentation

1780 provided alongside the assets?

1781 Answer: [Yes]

1782 Justification: The new assets introduced in the paper are well documented and the documen-

1783 tation is provided alongside the assets.

1784 Guidelines:

- 1785 • The answer NA means that the paper does not release new assets.
- 1786 • Researchers should communicate the details of the dataset/code/model as part of their
- 1787 submissions via structured templates. This includes details about training, license,
- 1788 limitations, etc.
- 1789 • The paper should discuss whether and how consent was obtained from people whose
- 1790 asset is used.
- 1791 • At submission time, remember to anonymize your assets (if applicable). You can either
- 1792 create an anonymized URL or include an anonymized zip file.

1793	14. Crowdsourcing and research with human subjects
1794	Question: For crowdsourcing experiments and research with human subjects, does the paper
1795	include the full text of instructions given to participants and screenshots, if applicable, as
1796	well as details about compensation (if any)?
1797	Answer: [NA]
1798	Justification: The paper does not involve crowdsourcing nor research with human subjects.
1799	Guidelines:
1800	• The answer NA means that the paper does not involve crowdsourcing nor research with
1801	human subjects.
1802	• Including this information in the supplemental material is fine, but if the main contribu-
1803	tion of the paper involves human subjects, then as much detail as possible should be
1804	included in the main paper.
1805	• According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1806	or other labor should be paid at least the minimum wage in the country of the data
1807	collector.
1808	15. Institutional review board (IRB) approvals or equivalent for research with human
1809	subjects
1810	Question: Does the paper describe potential risks incurred by study participants, whether
1811	such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1812	approvals (or an equivalent approval/review based on the requirements of your country or
1813	institution) were obtained?
1814	Answer: [NA]
1815	Justification: The paper does not involve research with human subjects.
1816	Guidelines:
1817	• The answer NA means that the paper does not involve crowdsourcing nor research with
1818	human subjects.
1819	• Depending on the country in which research is conducted, IRB approval (or equivalent)
1820	may be required for any human subjects research. If you obtained IRB approval, you
1821	should clearly state this in the paper.
1822	• We recognize that the procedures for this may vary significantly between institutions
1823	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1824	guidelines for their institution.
1825	• For initial submissions, do not include any information that would break anonymity (if
1826	applicable), such as the institution conducting the review.
1827	16. Declaration of LLM usage
1828	Question: Does the paper describe the usage of LLMs if it is an important, original, or
1829	non-standard component of the core methods in this research? Note that if the LLM is used
1830	only for writing, editing, or formatting purposes and does not impact the core methodology,
1831	scientific rigor, or originality of the research, declaration is not required.
1832	Answer: [NA]
1833	Justification: The paper does not involve LLMs as any important, original, or non-standard
1834	components.
1835	Guidelines:
1836	• The answer NA means that the core method development in this research does not
1837	involve LLMs as any important, original, or non-standard components.
1838	• Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM)
1839	for what should or should not be described.