

ACCELERATING DISCRETE DIFFUSION MODELS WITH PARALLEL SAMPLING

Anonymous authors

Paper under double-blind review

ABSTRACT

Discrete diffusion models are widely used for learning and generating discrete distributions. As the generation process is inherently sequential, the acceleration of sampling is of significant importance. In this work, we parallelize the main-stream τ -leaping algorithm for absorbing discrete diffusion in a Continuous-Time Markov Chain (CTMC) framework. By leveraging the continuous-time stochastic integral form of the τ -leaping algorithm and the Picard iteration method, we achieve parallel-in-time sampling acceleration. We implement a predictor-corrector structure based on the Markov chain Monte Carlo (MCMC) method to control for additional errors and provide a proof of exponential convergence for our algorithm. We improve the overall time complexity of τ -leaping from $\mathcal{O}(d^2 \log^2 d)$ to at most $\mathcal{O}(d^{3/2} \log^{5/2} d)$. In practice, our accelerated algorithm can achieve at most 5 to 8-fold sampling speedup over the traditional time-sequential τ -leaping with respect to wall-clock time on convex, non-convex, and high-dimensional discrete distributions. Our research broadens the scope of leveraging discrete diffusion models in various challenging areas like molecular structure generation and Large Language Models (LLMs).

1 INTRODUCTION

Diffusion models over discrete spaces (Austin et al., 2021) have become central to generative modeling for categorical data, with increasing impact across applications including molecular design (Vignac et al., 2022), protein engineering (Gruver et al., 2023), DNA sequence generation under biophysical constraints (Sarkar et al., 2024), and high-fidelity generation of text (Zheng et al., 2023), music (Yang et al., 2023), and images (Lezama et al., 2022).

However, due to the sequential nature of the sampling process, to generate a high-quality sample via diffusion models often requires a large number of sequential rounds, and sequential evaluation of the neural network-based score function with expensive time cost (Song et al., 2020). However, the algorithms underlying the above results are highly sequential and fail to fully exploit contemporary parallel computing resources such as multi-core central processing units (CPUs) and many-core graphics processing units (GPUs).

The primary focus of existing parallel sampling algorithms is on parallel-in-tokens (Wu et al., 2025; Xie et al., 2025). While parallel-in-time algorithms for continuous diffusion models have been extensively studied (Shih et al., 2023; Anari et al., 2024; Chen et al., 2024a; Zhou & Sugiyama, 2025), their counterparts for discrete diffusion models remain largely unexplored.. This gap motivates our investigation into the question:

Whether parallelization can fundamentally accelerate sampling for discrete diffusion models?

In this paper, we make significant progress in both theoretical and practical sides. The **key contributions** of this work are summarized as follows:

- **A fast parallel sampler.** We introduce the first *parallel-in-time* τ -leaping algorithm (Cao et al., 2006) for discrete diffusion inference, namely the Picard τ -leaping method. Our method has been shown as a fast sampler that reduces the time complexity from $\tilde{\mathcal{O}}(d^2)$

Table 1: Comparisons of our parallel methods and sequential methods for τ -leaping given δ^2 -accuracy of the score. Here, S denotes the size of vocabulary.

	Time Complexity	Space Complexity
(Campbell et al., 2022, Theorem 1)	$\mathcal{O}(\frac{d^4}{\delta})$	$\mathcal{O}(d)$
(Ren et al., 2025, Theorem 4.7)	$\mathcal{O}(\frac{d^2 \log^2(d/\delta^2)}{\delta^2})$	$\mathcal{O}(d)$
(Liang et al., 2025, Theorem 1)	$\mathcal{O}(\frac{d^2 S}{\delta^2})$	$\mathcal{O}(d)$
Ours, Theorem 3.1	$\mathcal{O}(\frac{d^{3/2} \log^{3/2}(d\delta^{-2})}{\delta})$	$\mathcal{O}(d^{3/2} \log^{1/2}(d))$

to $\tilde{\mathcal{O}}(d^{3/2})$ with slightly larger space complexity¹ (See Theorem 3.1). Here, $\tilde{\mathcal{O}}$ omits the logarithmic factors. We summarize the comparison between existing sequential methods and our results in Table 1.

- **Practical parallel sampling via Corrector.** We also propose two Markov Chain Monte Carlo (MCMC)-based correctors designed to address sampling in both simple and complex scenarios for further mitigating the parallel iteration and discretization error with detailed convergence analysis (See Section 4).
- **Experimental evaluation.** Our method achieves substantial improvements in sampling efficiency over the original sequential algorithm, measured by both the Number of Function Evaluations (NFE) and wall-clock runtime in experiments (See Section 5).

In Appendix B, we discuss more related works.

2 PRELIMINARIES ON DISCRETE DIFFUSION MODELS

In discrete diffusion models, the forward process is a continuous-time Markov chain (CTMC) $(\mathbf{x}_t)_{t \in [0, T]}$ on a finite state space \mathbb{X} . Let $\mathbf{p}_t \in \Delta^{|\mathbb{X}|}$ denote the law of \mathbf{x}_t as a column vector. The forward equation is

$$\frac{d\mathbf{p}_t}{dt} = \mathbf{Q}_t \mathbf{p}_t, \quad \mathbf{Q}_t = (\mathbf{Q}_t(y, x))_{x, y \in \mathbb{X}}, \quad (1)$$

where \mathbf{Q}_t is a rate matrix with (i) $\mathbf{Q}_t(x, y) \geq 0$ for $x \neq y$ and (ii) $\mathbf{Q}_t(x, x) = -\sum_{y \neq x} \mathbf{Q}_t(y, x)$. We write $\tilde{\mathbf{Q}}_t := \mathbf{Q}_t - \text{diag}(\mathbf{Q}_t)$ for the off-diagonal part. The time-reversed (backward) process $(\tilde{\mathbf{x}}_t)_{t \in [0, T]}$, with $\tilde{\mathbf{x}}_t := \mathbf{x}_{T-t}$, is again a CTMC with law $\tilde{\mathbf{p}}_s$ and generator $\tilde{\mathbf{Q}}_t$ satisfying (Kelly, 2011)

$$\frac{d\tilde{\mathbf{p}}_t}{dt} = \tilde{\mathbf{Q}}_t \tilde{\mathbf{p}}_t, \quad (2)$$

where for $x \neq y$, $\tilde{\mathbf{Q}}_t(y, x) = \frac{\tilde{\mathbf{p}}_t(y)}{\tilde{\mathbf{p}}_t(x)} \tilde{\mathbf{Q}}_t(x, y)$, and $\tilde{\mathbf{Q}}_t(x, x) = -\sum_{y' \neq x} \tilde{\mathbf{Q}}_t(y', x)$, and $\tilde{\mathbf{Q}}_t := \mathbf{Q}_{T-t}$.

According to Proposition 3.2 in (Ren et al., 2025), discrete diffusion models can also be interpreted as stochastic integrals with Poisson random measure. The forward process in discrete diffusion models (1) can thus be represented by the following stochastic integral:

$$\mathbf{x}_t = \mathbf{x}_0 + \int_0^t \int_{\mathbb{D}} \nu N[\lambda](dt, d\nu), \quad (3)$$

where the intensity λ is defined as $\lambda_t(\nu, \omega) = \tilde{\mathbf{Q}}_t(x_{t-}(\omega) + \nu, x_{t-}(\omega))$ 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) can also be represented similarly as

$$\mathbf{y}_t = \mathbf{y}_0 + \int_0^t \int_{\mathbb{D}} \nu N[\mu](ds, d\nu), \quad (4)$$

¹We note, in this paper, that the space complexity refers to the number of words (Cohen-Addad et al., 2023) instead of the number of bits (Goldreich, 2008) to denote the approximate required storage.

where the intensity μ is defined as

$$\mu_t(\nu, \omega) = \bar{s}_t(\mathbf{y}_{t-}, \mathbf{y}_{t-} + \nu) \bar{Q}_s(\mathbf{y}_{s-}, \mathbf{y}_{s-} + \nu) \quad (5)$$

if $\mathbf{y}_{s-} + \nu \in \mathbb{X}$ and 0 otherwise. During inference,

$$\hat{\mathbf{y}}_t = \hat{\mathbf{y}}_0 + \int_0^t \int_{\mathbb{D}} \nu N[\hat{\mu}](ds, d\nu) \quad (6)$$

is used instead of Equation (4), where the estimated intensity $\hat{\mu}$ is defined by replacing the true score \mathbf{s}_t with the neural network (NN) estimated score $\hat{\mathbf{s}}_t$ in Equation (5).

3 PARALLEL SAMPLING FOR ABSORBING DISCRETE DIFFUSION MODELS

In this section, we introduce the parallel-in-time τ -leaping algorithm for discrete diffusion models in a CTMC framework. The key idea is to extend the application of Picard iteration (See Appendix C.2) from the domain of diffusion models in continuous state spaces to the stochastic integral form of τ -leaping (Cao et al., 2006) in discrete state spaces.

This section is organized as follows. Section 3.1 introduces and details our parallel algorithm based on Picard iteration. Section 3.2 then states the assumptions and their justifications that are necessary for the subsequent theoretical analysis. Finally, Section 3.3 presents the main theorem regarding the algorithm’s error bound and computational complexity, followed by its proof.

3.1 OUR ALGORITHM

Discretization Scheme. In this work, the total time horizon T is firstly segmented into N blocks with the large time grid $(t_n)_{n \in [0, N]}$, $t_0 = 0, t_N = T$. For each block, the time interval $[t_n, t_{n+1}]$ is divided into M steps with the small time grid $(\tau_{n,m})_{n \in [0, N-1], m \in [0, M]}$, $\tau_{n,0} = t_n$, and $\tau_{n,M} = t_{n+1}$. We also define the small step size $\epsilon = \tau_{n,m} - \tau_{n,m-1}$. It is also possible to set $M = M_n$ and $\epsilon = \epsilon_n$ for flexible grid when using early stopping.

When describing the standard τ -leaping algorithm with the form in (equation 6), the main update step in Picard iteration can be described as

$$\hat{\mathbf{y}}_{t_n+m\epsilon}^{(k+1)} = \hat{\mathbf{y}}_{t_n} + \sum_{j=0}^{m-1} \left(\sum_{y' \in \mathbb{X}} (y' - \hat{\mathbf{y}}_{t_n+j\epsilon}^{(k)}) \cdot \mathcal{P}(\hat{\mu}_{t_n+j\epsilon}^\theta(y' | \hat{\mathbf{y}}_{t_n+j\epsilon}^{(k)}) \cdot \epsilon) \right), \quad (7)$$

where $\hat{\mathbf{y}}_{t_n+j\epsilon}^{(k)}$ is the sample state at $t = t_n + j\epsilon$ in the k -th Picard iteration. $\mathcal{P}(\hat{\mu}_{t_n+j\epsilon}^\theta(y' | \hat{\mathbf{y}}_{t_n+j\epsilon}^{(k)}))$ is the numerical form of $N[\hat{\mu}_{t_n+j\epsilon}^\theta]$, which denotes the number of jumps from $\hat{\mathbf{y}}_{t_n+j\epsilon}^{(k)}$ to y' during the small time interval ϵ . In the block-wise parallel algorithm, the computation for each block starts from its initial state $\hat{\mathbf{y}}_{t_n}$ which is also the terminal value of the last block.

One may notice that such update includes a summation among all possible adjacent states, i.e. $\sum_{y' \in \mathbb{X}}$, which would lead to an exponential computation complexity with respect to the data dimension d . In practice, the computation cost is controlled at $\mathcal{O}(Sd)$ with the vocabulary size S , using techniques such as dimension factorization (Campbell et al., 2022) and Top-K candidate set truncation (Ye et al., 2024). Our analysis is also based on such settings.

The Picard method transforms the dependence on the previous time step in the τ -leaping process into a dependence on the entire trajectory of the previous iteration. The state of the new trajectory at time point $t_n + m\epsilon$ is equal to the initial state of the block $\hat{\mathbf{y}}_{t_n}$, plus the sum of all jump vectors from $j = 0$ to $j = m - 1$ in the last iteration which is fully known, therefore the calculation of jump for all time steps within a block can be fully parallelized with cumulative sum operation using *parallel prefix sum* algorithm, eliminating the need for a sequential loop. The overall computation process is summarized in Algorithm 1.

One may notice that there is an optional corrector added after computing one block. Such a module can compensate for the additional error introduced by the Picard iteration for a minor increase in the number of function evaluations (NFE). This, in turn, relaxes the constraints on the parameters of the parallel algorithm, such as the number of blocks N or small intervals M . A more detailed analysis will be presented in Section 4. It is important to note that the following analysis in Section 3.3 will focus exclusively on the parallel τ -leaping method itself without considering the corrector.

Algorithm 1: Parallel τ -Leaping Algorithm for Discrete Diffusion Model Sampling

Input : $\hat{y}_0 \sim q_0$, large time grid $(t_n)_{n \in [0, N]}$ with $t_0 = 0, t_N = T - \xi$, small time grid $(\tau_{n,m})_{n \in [0, N-1], m \in [0, M]}$ with $\tau_{n,0} = t_n, \tau_{n,M} = t_{n+1}$ and small step size $\epsilon = \tau_{n,m+1} - \tau_{n,m}$; Picard depth K_p ; intensity $\hat{\mu}_s^\theta$; score estimate \hat{s}_t^θ ; fixed random seeds.

Output: A sample $y_{t_N} \sim q_{t_N}$.

for $n = 0$ **to** $n = N - 1$ **do**

Initial Guess: $\hat{y}_{t_n+m\epsilon}^{(0)} \leftarrow \hat{y}_{t_n}$ for all $m \in [0, M]$

for $k = 0$ **to** $K_p - 1$ **do**

for $j = 0$ **to** $M - 1$ *Parallel do*

$\hat{\mu}_j^\theta \leftarrow \hat{\mu}_{t_n+j\epsilon}^\theta(\cdot | \hat{y}_{t_n+j\epsilon}^{(k)})$

$J_{j,y'} \sim \mathcal{P}(\hat{\mu}_j^\theta(y') \cdot \epsilon)$ for all $y' \in \mathcal{X}_{neighbor}$

$\Delta \hat{y}_j^{(k)} = \sum_{y' \in \mathcal{X}_{neighbor}} (y' - \hat{y}_{t_n+j\epsilon}^{(k)}) \cdot J_{j,y'}$

end

for $m = 1$ **to** M *Parallel do*

$\hat{y}_{t_n+m\epsilon}^{(k+1)} \leftarrow \hat{y}_{t_n} + \sum_{j=0}^{m-1} \Delta \hat{y}_j^{(k)}$

end

end

$y_{t_{n+1}} \leftarrow \hat{y}_{t_n+M\epsilon}^{(K_p)}$

$\hat{y}_{n+1} \leftarrow \text{Corrector}(y_{t_{n+1}}, \hat{s}_t^\theta)$ (Optional)

end

3.2 ASSUMPTION

We first list the assumptions that will be used in the subsequent analysis and proofs, which mainly follows the settings in (Ren et al., 2025) due to the shared CTMC and stochastic integral frameworks. These assumptions align with those established in previous theoretical works, such as those described by (Ren et al., 2025).

Assumption 3.1 (Bounded Score and Rate Matrix). The true score function $s_t(x, y)$, learned score function $\hat{s}_t^\theta(x, y)$, and the symmetric time-homogeneous rate matrix Q satisfy:

- (i) $\hat{s}_t^\theta(x, y) \in [m_s, M_s]$ and $s_t(x, y) \lesssim 1 \vee t^{-\gamma}$, for $\forall x, y \in \mathbb{X}, t > 0$, and a constant $\gamma \in [0, 1)$.
- (ii) $Q(x, y) \leq C$ and $|Q(x, x)| = |-\sum_{y \neq x} Q(x, y)| \leq \bar{D}$ for $\forall x, y \in \mathbb{X}, t > 0$ and some constants $C, \bar{D} > 0$.
- (iii) $\rho(Q)$ as the modified log-Sobolev constant of Q has a lower bound $\rho > 0$.

Assumption 3.2 (Score Estimation Accuracy). The neural network $\hat{s}_t^\theta(x_t)$ learned the true score $s_t(x_t, t)$ with δ -accuracy, which is described as

$$\sum_{n=0}^{N-1} (t_{n+1} - t_n) \mathbb{E} \left[\int_{\mathbb{X}} K \left(\frac{\hat{s}_{t_n}^\theta(x_{t_n}, y)}{\hat{s}_{t_n}^\theta(x_{t_n}, y)} \right) \hat{s}_{t_n}^\theta(x_{t_n}, y) \tilde{Q}(x_{t_n}, y) \nu(dy) \right] \leq \delta^2,$$

where $K = x - 1 - \log x$ and \tilde{Q} denotes the matrix Q with the diagonal elements set to zero.

Assumption 3.3 (Continuity of Score Function). For any $t > 0$ and $y \in \mathbb{X}$ such that $Q(x_{t-}, y) > 0$, we have

$$\left| \frac{\mu_{t^+}(y)}{\mu_t(y)} \right| := \left| \frac{p_t(x_{t-})Q(x_t, y)}{p_t(x_t)Q(x_{t-}, y)} - 1 \right| \lesssim 1 \vee t^{-\gamma},$$

for some exponent $\gamma \in [0, 1)$ and t^+, t^- denotes the right and left limits.

Remark 3.1. For simplicity, we do not consider the early-stopping in our framework, which requires $\gamma < 1$ according to Proposition C.5 in (Ren et al., 2025).

Assumption 3.4 (Lipschitz Continuity). Given the learned intensity function $\hat{\mu}_t^\theta(y|x) = \hat{s}_t^\theta(x, y) \tilde{Q}(x, y)$ and the small step size ϵ , for $\forall x_1, x_2 \in \mathbb{X}$, the following Lipschitz condition with respect to the Poisson expectation holds:

$$\sum_{y' \in \mathbb{X}} \mathbb{E}[(\mathcal{P}(\hat{\mu}_t^\theta(y'|x_1) \cdot \epsilon) - \mathcal{P}(\hat{\mu}_t^\theta(y'|x_2) \cdot \epsilon))^2] \leq (L_\mu \epsilon)^2 \|x_1 - x_2\|^2.$$

Assumption 3.4 is newly introduced which plays an important role in deriving the convergence of Picard iteration as shown in (Anari et al., 2024; Chen et al., 2024a; Zhou & Sugiyama, 2025). Since the intensity function is constructed by the learned score and the non-diagonal rate of the forward process (Ren et al., 2025), and both two components have continuity or bounded properties based on Section 2 and Assumption 3.3, the setting could be reasonable.

3.3 THEORETICAL GUARANTEES

Next, we provide an approximated error bound of our algorithm without correctors.

Theorem 3.1 (Theoretical Guarantees for Parallel τ -Leaping). *Under the assumptions above in this section, the Picard τ -leaping algorithm (Algorithm 1) achieves the following approximation error bound:*

$$D_{KL}(p_0 \|\hat{q}_T^{(K_p)}) \lesssim \underbrace{e^{-\rho T} \log |\mathbb{X}|}_{\text{truncation error}} + \underbrace{\delta^2}_{\text{score estimation error}} + \underbrace{\bar{D}^2 \epsilon T}_{\text{discretized error}} + \underbrace{C^{K_p} \text{poly}(\bar{D})}_{\text{parallization error}}$$

with constant $C < 1$. Furthermore, by choosing parameters of the following order and with assuming $\bar{D} = \mathcal{O}(d)$ (Ren et al., 2025),

$$\begin{aligned} T &= \mathcal{O}(\log(d\delta^{-2})), \quad \epsilon = \mathcal{O}(\delta^2 d^{-2} \log^{-1}(d\delta^{-2})), \quad K_p = \mathcal{O}(\log(d\delta^{-2})), \\ h &= \mathcal{O}(\delta d^{-3/2} \log^{-1/2}(d\delta^{-2})), \quad N = T/h = \mathcal{O}(\delta^{-1} d^{3/2} \log^{5/2}(d\delta^{-2})), \end{aligned}$$

the algorithm achieves $NK = \mathcal{O}(\delta^{-1} d^{3/2} \log^{5/2}(d\delta^{-2})) = \tilde{\mathcal{O}}(\delta^{-1} d^{3/2})$ approximate time complexity and $dM = dh/\epsilon = \mathcal{O}(\delta^{-1} d^{3/2} \log^{1/2}(d\delta^{-2}))$ space complexity.

The meaning of notations is shown in Algorithm 1. Compared with the sequential τ -leaping, our parallel algorithm improves the approximate time complexity from $\tilde{\mathcal{O}}(d^2)$ to $\tilde{\mathcal{O}}(d^{3/2})$ with the cost of increasing $\mathcal{O}(d^{1/2} \log^{1/2}(d))$ space complexity. The proof can be found in Appendix D.

The main obstacle to improving the parallel complexity to poly $\log(d)$ is that even under Lipschitz condition (Assumption 3.4), ensuring convergence of the Picard iteration requires the length of time slice scales as $\delta/(d^{3/2} L_p)$ rather than $1/L_p$ in continuous case where L_p denotes the smoothness of the intensity.

Remark 3.2 (Practical choices of the depth). Following the practical settings of Picard-based parallel algorithms in continuous diffusion models (Shih et al., 2023), one may choose depth from 2 to 12 depending on the complexity of the problem and See Section 5.1.1 for experimental results.

4 PARALLEL SAMPLING WITH CORRECTORS

Although Theorem 3.1 provides an error bound for the parallel sampling algorithm, the block width h cannot achieve $\mathcal{O}(1)$ with respect to d . Consequently, in practical applications, the selection of the width or the number of parallel steps is more restricted compared to parallel continuous diffusion (Chen et al., 2024a). To address this limitation, in this section, we propose two MCMC-based Plug-and-Play correctors to refine the generated distribution after the parallel algorithm completes the computation of a block. We provide a general error analysis of the corrector described in the following theorem.

Theorem 4.1 (Error Control of the MCMC Corrector). *Suppose \bar{p}_{t_n} is the true target distribution derived from the forward process, and \hat{q}_{t_n} is the generated distribution from our predictor-corrector algorithm 2 and 3. The MCMC-based corrector with stationary distribution $\bar{p}_{t_{n+1}}$ described has the following rate of error control:*

$$E_{n+1} := TV(\hat{q}_{t_{n+1}}, \bar{p}_{t_{n+1}}) \leq \gamma^{K_c} (E_n + \Delta_p) + \Delta_c, \quad (8)$$

where \bar{p}_{t_n} is the true distribution and \hat{q}_{t_n} is the generated distribution at time t_n . Δ_p denotes the discretization error from the predictor, Δ_c is the corrector approximation error from learned score, $\gamma < 1$ is the contraction coefficient related to the conductance of the problem, and K_c is the number of corrector steps.

The proof is demonstrated in Appendix E.1.

Remark 4.1 (Comparison with Theorem 3.1). Since the theoretical guarantees analyzed in Theorem 3.1 do not take the corrector into consideration, the error analysis in Theorem 4.1 is independent to the analysis in Theorem 3.1. However, based on the experimental results in Section 5, we believe that the condition and bound in Theorem 3.1 could be relaxed and improved after adding correctors to the proof framework, which will be explored in future work.

4.1 METROPOLIS–HASTINGS CORRECTOR

We first introduce a Metropolis–Hastings Algorithm-based MCMC corrector with symmetric proposal distribution for experiments in the next section. The computation framework is summarized in Algorithm 2.

Algorithm 2: Metropolis–Hastings Algorithm as Corrector (Corrector_{MH}($y_{t_{n+1}}$, \widehat{s}_t^θ))

Input : Sample from the predictor $y_{\text{pred}} = y_{t_{n+1}}$,
symmetric proposal distribution $Q(y'|y) = Q(y|y')$, score estimate \widehat{s}_t^θ .

Output: Corrected sample $\widehat{y}_{t_{n+1}}$

Initialization: $y_{\text{current}} \leftarrow y_{\text{pred}}$

for $k = 0$ **to** K_c **do**

$y_{\text{proposed}} \sim Q(y'|y_{\text{current}})$

$A \leftarrow \min(1, \widehat{s}_{t_{n+1}}^\theta(y_{\text{current}}, y_{\text{proposed}}))$

$u \sim \text{Unif}([0, 1])$

if $u < A$ **then**

$y_{\text{current}} \leftarrow y_{\text{proposed}}$

end

end

Given the output of the predictor after one block sampling $y_{t_{n+1}} = \widehat{y}_{t_n + M\epsilon}^{(K_p)} := y_{\text{pred}}$, the corrector runs a short, serial Metropolis-Hastings (MH) MCMC process using $\bar{p}_{t_{n+1}}$, the true distribution at t_{n+1} , as the stationary distribution. The MH corrector accepts or rejects new proposed states according to the probability distribution at the current timestep t_{n+1} provided by the learned score function. This process pulls $y_{t_{n+1}}$ from potentially low-probability regions back into high-probability ones, effectively correcting the error.

Lemma 4.1 (Acceptance rate with symmetric proposal distributions). Given the estimated score \widehat{s}_t^θ with the symmetric proposal distribution $Q(y'|y) = Q(y|y')$, the acceptance rate of changing from y to y' the Metropolis-Hastings MCMC corrector can be approximated as

$$A(y'|y) \approx \min(1, \widehat{s}_t^\theta(y, y')), \quad (9)$$

The lemma above can be easily proved according to the definition of acceptance rate and symmetric proposal distributions Q :

$$A(y'|y) = \min\left(1, \frac{\pi(y')Q(y|y')}{\pi(y)Q(y'|y)}\right) = \min\left(1, \frac{\pi(y')}{\pi(y)}\right).$$

where $\pi(\cdot) = \bar{p}_{t_{n+1}}$ is our target stationary distribution. Then, according to the definition of concrete score in (Lou et al., 2024):

$$\frac{\bar{p}_{t_{n+1}}(y')}{\bar{p}_{t_{n+1}}(y)} = \bar{s}_{t_{n+1}}(y, y'),$$

we can naturally use the learned score as the approximated acceptance rate.

The setting of symmetric proposal works well in low-dimension cases (see Section 5). However it makes the corrector difficult to work in more complex situations such as language models because such symmetry does not naturally hold, and the proposal sample requires a uniform sampling in the entire vocabulary which leads to slow mixture and acceptance rate (Levin & Peres, 2017).

4.2 TOP-K BARKER CORRECTOR

To adapt the corrector to more complex cases, with the example of absorbing diffusion, we introduce an MCMC corrector with Uniform Top-K Proposal and Barker Acceptance (Livingstone & Zanella, 2022) in Algorithm 3.

Algorithm 3: Top-K Barker Algorithm as Corrector ($\text{Corrector}_{\text{Barker}}(y_{t_{n+1}}, \hat{s}_t^\theta)$)

Input : Current time t , vocabulary $\mathcal{V} = \{1, \dots, n\}$, $x_{\text{pred}} = (x_1, \dots, x_d) \in \mathcal{V}^d$,
 score estimate \hat{s}_t^θ , candidate count K_b , subset of positions to correct $\mathcal{I} \subseteq \{1, \dots, d\}$,
 absorbing token $M' \in \mathcal{V}$

Output: Corrected sample $\hat{y}_{t_{n+1}}$

Create masked copy: $\tilde{x} = x^{(i \leftarrow M')}$ such that $\tilde{x}_i \leftarrow M'$ for $\forall i \in \mathcal{I}$

$r_{i,j} \leftarrow \hat{s}_t^\theta(x^{(i \leftarrow M)}, x^{(i \leftarrow j)})$ for $\forall i \in \mathcal{I}, j \in \mathcal{V}$

Compute unnormalized log-weights: $u_{i,j} \leftarrow \log r_{i,j}$

for $i \in \mathcal{I}$ **do**

Sort $j \in \mathcal{V}$ by $u_{i,j}$ descending and select top K_b indices as \mathcal{K}_i

if $a := x_i \notin \mathcal{K}_i$ **then**

| Replace the last element of \mathcal{K}_i by a

end

$b \sim q_i(\cdot)$, where $q_i(j) = 1/K_b, j \in \mathcal{K}_i$

$\Delta_i \leftarrow u_{i,b} - u_{i,a}$

$\alpha_i \leftarrow \frac{1}{1 + \exp(-\Delta_i)}$

$u \sim \text{Unif}[0, 1]$

if $U \leq \alpha_i$ **then**

| $x_i \leftarrow b$

end

end

The method uses the network’s outputs to score single-site replacements, but builds proposals in a MASK-anchored context so that the proposal is independent of the current token.

A Top-K truncation focuses computation on promising candidates, and a uniform or softmax proposal over this set keeps the step cheap while preserving correctness. We use the Barker acceptance, yielding a smooth decision that satisfies the detailed balance condition. So the corrector provably leaves the reverse-time marginal invariant rather than introducing a bias. The entire correction round requires only one network forward for a batch of positions and can be fully parallelized.

Definition 4.1 (Single-Coordinate Barker Kernel). Let $x = \{x_1, \dots, x_d\} \in \mathbb{X}$ be the state, and $\mathbb{X} = \mathcal{V}^d$ be a finite state space with vocabulary $\mathcal{V} = \{1, \dots, n\}$ and data dimension d . For absorbing diffusion, we set the mask token as M' . For a position $i \in \{1, \dots, d\}$ and element (token) $j \in \mathcal{V}$, we denote $x^{(i \leftarrow j)}$ as the state after replacing x_i with j . Suppose the target distribution is $\pi_t : \mathbb{X} \rightarrow (0, 1]$ with full support. Define $u_{i,j}(x) = \hat{s}_t^\theta(x, x^{(i \leftarrow j)})$ given a trained score \hat{s}_t^θ , and an anchor context $\tilde{x}^{(i)} := x^{(i \leftarrow M')}$. Construct a candidate set $\mathcal{K}_i(\tilde{x}^{(i)}) \subseteq \mathcal{V}$ of fixed size K_b , and define a proposal distribution $q_i(\cdot | x)$ supported and uniform on it. We define a single-coordinate Barker kernel $K_i(x, x')$. Given x , one first sample $b \sim q_i(\cdot | x)$ and construct $x' = x^{(i \leftarrow b)}$. Then accept x' with probability

$$\alpha_i(x \rightarrow x') = \frac{1}{\exp(-\Delta_i)}, \quad \Delta_i = u_{i,b} - u_{i,a},$$

where $a := x_i$ is the current token.

The designed kernel has the property described in Theorem 4.2.

Theorem 4.2 (Invariance of the Top-K Barker Corrector Kernel). For all $i \in \{1, \dots, d\}$, the Markov kernel K_i is π_t -reversible and has π_t as the stationary distribution satisfying the detailed balance condition.

$$\pi_t(x)K_i(x, x') = \pi_t(x')K_i(x', x). \quad (10)$$

See Appendix E.2 for the proof. With the proposition above, Algorithm 3 can then be incorporated into the error analysis framework of Theorem 4.1.

5 EXPERIMENTS

In this section, we evaluate the performance of our parallel τ -leaping algorithm against the sequential τ -leaping algorithm. The evaluation is conducted on four synthetic distributions shown in Table 2, image generation based on Imagenet, and text generation based on OpenWebText.

5.1 SYNTHETIC EXPERIMENTS

Oracle Models To directly verify the speedup of our parallel τ -leaping algorithm over its sequential counterpart, we employ an Oracle Model. This model provides the sampler with perfect, analytically tractable score information, thereby eliminating approximation errors from the score network. Consequently, the quality of the final samples serves as a direct measure of the sampling algorithm’s intrinsic precision and stability.

Controlled Variables The essence of the Parallel-in-time acceleration strategy is to partition the original total number of serial time steps U into N blocks. The U/N time steps within each block are then computed through parallel iterations, thereby reducing the total amount of serial computation and the Number of Function Evaluations (NFE).

For a fair comparison, we fix the number of correction steps K_c , parallel blocks N , and steps per block M in each experiment. The sequential baseline algorithm runs for a total of $M \times N$ steps, with a correction applied every M steps. The number of samples is controlled to avoid GPU memory bottlenecks, and all experiments utilize the Metropolis-Hastings corrector.

We conduct experiments on sampling 4 types of synthetic distributions: Chessboard, Circle, Hypercube, and Embedded Hypercube. The feature and settings of parameters are demonstrated in Table 2. One can check more detailed explanation in Appendix F.

Table 2: Characteristics and Parameter Settings of Experimental Distributions. Runtimes and KL Divergence for all experiments are averaged over 20 runs.

Distribution	Characteristics	Dimensions/Grid	N	M	K_c	Samples
Chessboard	Sparse, Multi-modal	8×8	40	50	5	8196
Circle	Non-Convex, Connected	32×32	60	50	8	4096
Hypercube	High-dimension	6-d	100	50	10	2048
Embed-Hypercube	Low-dimension manifold	6-d in 12-d	50	100	8	2048

5.1.1 PERFORMANCE

Given that the Picard iteration depth K_p is a parameter unique to the parallel sampling algorithm, we evaluate its performance at various K_p values against the sequential algorithm. The comparison is based on three key metrics: Wall-Clock Runtime, NFE, and the KL Divergence of the generated distribution. Table 3 shows the comprehensive performance comparison across four baseline distributions.

As can be seen from the table, our parallel predictor-corrector algorithm significantly reduces the Number of Function Evaluations (NFE) and clock-wall runtime while maintaining high accuracy across the vast majority of experimental settings. However, we also observe a diminishing marginal return with a continued increase in the iteration depth K_p . Furthermore, when the batch size is large, increasing the number of parallel iteration steps leads to higher data I/O and communication cost. Once the GPU bandwidth bottleneck is reached, the efficiency of the parallel algorithm will decrease. Therefore, in practical applications, a moderate iteration depth is sufficient.

Table 3: Performance comparison on four distributions. Top: Performance on Chessboard and Circle distributions. Bottom: Performance on Hypercube and Embedded Hypercube distributions.

Alg.	K_p	Chessboard			Circle		
		Runtime (s)	KL Divergence↓	NFE	Runtime (s)	KL Divergence↓	NFE
Seq.	/	2.43 ± 0.22	0.0021 ± 0.0005	2200	3.53 ± 0.05	0.0493 ± 0.0038	3480
Par.	2	0.36 ± 0.06	0.0023 ± 0.0004	280	0.77 ± 0.09	0.0500 ± 0.0023	600
	4	0.54 ± 0.11	0.0021 ± 0.0006	360	1.14 ± 0.08	0.0498 ± 0.0047	720
	6	0.70 ± 0.07	0.0019 ± 0.0005	440	1.33 ± 0.09	0.0495 ± 0.0030	840
	8	0.88 ± 0.07	0.0019 ± 0.0004	520	1.89 ± 0.06	0.0491 ± 0.0042	960
	10	1.07 ± 0.10	0.0018 ± 0.0003	600	2.22 ± 0.07	0.0488 ± 0.0043	1080
	12	1.26 ± 0.09	0.0019 ± 0.0004	680	1.59 ± 0.08	0.0482 ± 0.0050	1200

Alg.	K_p	Hypercube			Embedded Hypercube		
		Runtime (s)	KL Divergence↓	NFE	Runtime (s)	KL Divergence↓	NFE
Seq.	/	3.42 ± 0.07	1.0287 ± 0.0263	5400	3.52 ± 0.13	0.1647 ± 0.0095	6000
Par.	2	1.17 ± 0.06	1.0171 ± 0.0275	500	0.99 ± 0.03	0.1611 ± 0.0105	1200
	4	2.03 ± 0.10	1.0174 ± 0.0296	600	1.18 ± 0.05	0.1582 ± 0.0103	1400
	6	2.88 ± 0.07	1.0078 ± 0.0404	700	1.36 ± 0.05	0.1544 ± 0.0060	1600
	8	3.70 ± 0.05	1.0002 ± 0.0354	800	1.55 ± 0.04	0.1531 ± 0.0083	1800
	10	4.55 ± 0.09	1.0055 ± 0.0271	900	1.77 ± 0.06	0.1552 ± 0.0110	2000
	12	5.42 ± 0.06	0.9982 ± 0.0252	1000	1.89 ± 0.08	0.1564 ± 0.0078	2200

5.1.2 ABLATION STUDY

We conducted an analysis to independently and jointly assess the impact of two distinct algorithmic components—the Predictor and the Corrector—on the final KL divergence on circle and hypercube distributions. The result in Table 4 shows the importance of both correctors and predictors.

Table 4: Ablation study of KL Divergence under different sampling strategies.

Sampling Strategy		KL Divergence ↓	
(K_p, K_c)	Description	Ring Distribution	Hypercube Distribution
(0, 0)	Baseline	1.6076 ± 0.1551	3.4351 ± 0.1199
(2, 0)	Predictor-Only	0.1203 ± 0.0212	1.2730 ± 0.0432
(0, 1)	Corrector-Only	0.1161 ± 0.0066	0.9388 ± 0.0320
(2, 1)	Predictor + Corrector	0.0628 ± 0.0038	0.7360 ± 0.0132

5.2 REAL-WORLD DATA EXPERIMENTS

To further test the algorithm’s performance on complex high-dimensional real cases, we have added experiments on the acceleration effects in image and text generation from the perspective of the Number of Function Evaluations (NFE). The experiment was conducted on a single NVIDIA RTX4090 laptop GPU with 16GB memory. More experiments about memory cost and hyperparameter are shown in Appendix F.2.

Table 5: Parameter settings and GPU memory cost for image and text generation.

Tasks	Seq. Memeory Cost	Para. Memory Cost	M	N	K_p
Text	6.7GB	11.6GB	8	32	2
Image	3.0GB	3.7GB	10	20	2

5.2.1 IMAGE GENERATION

The experiment utilized a MaskGIT-based score model (Chang et al., 2022; Besnier & Chen, 2023) pretrained on ImageNet (Deng et al., 2009). We compared the performance of both sequential and parallel τ -leaping when generating 256×256 resolution images and evaluated the The Fréchet Inception Distance (FID) score based on 10000 samples. The parameter settings are shown in Table 5, and the results are shown in Table 6. Based on the results, when compressing the NFE from 400 to 40, our method still achieves a comparable FID score.

Table 6: Performance of sequential and parallel τ -leaping for image generation.

Imagenet	NFE	Runtime(s)	FID↓
Sequential	80	1.92	13.38
Parallel	40	1.43	12.05

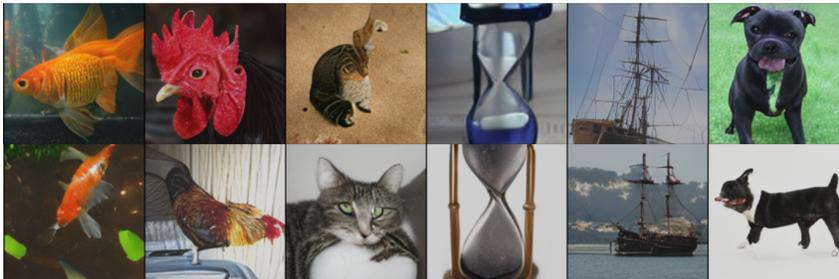


Figure 1: Generated samples from the Imagenet experiments.

5.2.2 TEXT GENERATION

The experiment utilized an RADD-based score model (Ou et al., 2024) pretrained on the OpenWebText dataset (Gokaslan & Cohen, 2019) which has GPT-2-level text generation capabilities (Radford et al., 2019). We compared the performance of both sequential and parallel τ -leaping when generating 512-token texts, and evaluated the average generative perplexity score with 1024 samples. The parameter settings are shown in Table 5, and the results are shown in Table 7. Based on the results, our method achieved a better average perplexity while reducing the NFE from 128 to 64.

Table 7: Performance of sequential and parallel τ -leaping for text generation.

Language Model	NFE	Runtime(s)	Avg. Perplexity↓
Sequential	128	2.84	52.698
Parallel	64	2.36	48.747

6 OUR CONCLUSION

In this work, we propose a parallel-in-time τ -leaping algorithm for discrete diffusion models based on Picard iteration within the Continuous-Time Markov Chain (CTMC) and stochastic integral framework. We reduce the time complexity from $\tilde{O}(d^2)$ to $\tilde{O}(d^{3/2})$ and provide proofs for both computational complexity and error analysis. Furthermore, we introduce MCMC-based predictor-corrector strategies to effectively control the additional error introduced by parallel computation. Our method achieves significant acceleration on various synthetic distribution sampling tasks, providing innovative insights and directions for the future development of discrete diffusion models.

540 **Reproducibility Statement** The source code of the experiments is provided in supplemental ma-
541 terials as the form of Jupyter Notebook.

542
543 **Ethics Statement** N/A

544
545 **REFERENCES**

546
547 Nima Anari, Sinho Chewi, and Thuy-Duong Vuong. Fast parallel sampling under isoperimetry.
548 *Proceedings of Thirty Seventh Conference on Learning Theory*, pp. 161–185, 2024.

549 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured
550 denoising diffusion models in discrete state-spaces. *Advances in neural information processing*
551 *systems*, 34:17981–17993, 2021.

552
553 Victor Besnier and Mickael Chen. A pytorch reproduction of masked generative image transformer.
554 *arXiv preprint arXiv:2310.14400*, 2023.

555 Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and
556 Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural*
557 *Information Processing Systems*, 35:28266–28279, 2022.

558 Yang Cao, Daniel T Gillespie, and Linda R Petzold. Efficient step size selection for the tau-leaping
559 simulation method. *The Journal of chemical physics*, 124(4), 2006.

560
561 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative
562 image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
563 *recognition*, pp. 11315–11325, 2022.

564 Chen-Hao Chao, Wei-Fang Sun, Hanwen Liang, Chun-Yi Lee, and Rahul G Krishnan. Be-
565 yond masked and unmasked: Discrete diffusion models via partial masking. *arXiv preprint*
566 *arXiv:2505.18495*, 2025.

567
568 Haoxuan Chen, Yinuo Ren, Lexing Ying, and Grant Rotskoff. Accelerating diffusion models with
569 parallel sampling: Inference at sub-linear time complexity. *Advances in Neural Information Pro-*
570 *cessing Systems*, 37:133661–133709, 2024a.

571 Zixiang Chen, Huizhuo Yuan, Yongqian Li, Yiwen Kou, Junkai Zhang, and Quanquan Gu. Fast
572 sampling via discrete non-markov diffusion models with predetermined transition time. *Advances*
573 *in Neural Information Processing Systems*, 37:106870–106905, 2024b.

574
575 CW Clenshaw. The numerical solution of linear differential equations in Chebyshev series. In
576 *Mathematical Proceedings of the Cambridge Philosophical Society*. Cambridge University Press,
577 1957.

578 Vincent Cohen-Addad, David P Woodruff, and Samson Zhou. Streaming Euclidean k -median and
579 k -means with $o(\log n)$ Space. In *2023 IEEE 64th Annual Symposium on Foundations of Computer*
580 *Science (FOCS)*. IEEE, 2023.

581 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
582 archical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
583 pp. 248–255. Ieee, 2009.

584
585 Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. [http://Skylion007.github.io/
586 OpenWebTextCorpus](http://Skylion007.github.io/OpenWebTextCorpus), 2019.

587 Oded Goldreich. Computational complexity: a conceptual perspective. *ACM Sigact News*, 2008.

588 Nate Gruver, Samuel Stanton, Nathan Frey, Tim GJ Rudner, Isidro Hotzel, Julien Lafrance-Vanasse,
589 Arvind Rajpal, Kyunghyun Cho, and Andrew G Wilson. Protein design with guided discrete
590 diffusion. *Advances in neural information processing systems*, 36:12489–12517, 2023.

591
592 Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and
593 Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of*
the IEEE/CVF conference on computer vision and pattern recognition, pp. 10696–10706, 2022.

- 594 Shivam Gupta, Linda Cai, and Sitan Chen. Faster diffusion-based sampling with randomized mid-
595 points: Sequential and parallel. *arXiv e-prints*, pp. arXiv-2406, 2024.
- 596
- 597 Frank P Kelly. *Reversibility and stochastic networks*. Cambridge University Press, 2011.
- 598
- 599 Michel Ledoux. The geometry of markov diffusion generators. In *Annales de la Faculté des sciences*
600 *de Toulouse: Mathématiques*, 2000.
- 601
- 602 David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathe-
603 matical Soc., 2017.
- 604
- 605 Jose Lezama, Tim Salimans, Lu Jiang, Huiwen Chang, Jonathan Ho, and Irfan Essa. Discrete
606 predictor-corrector diffusion models for image synthesis. In *The Eleventh International Confer-*
ence on Learning Representations, 2022.
- 607
- 608 Yuchen Liang, Yingbin Liang, Lifeng Lai, and Ness Shroff. Discrete diffusion models: Novel
609 analysis and new sampler guarantees. *arXiv preprint arXiv:2509.16756*, 2025.
- 610
- 611 Samuel Livingstone and Giacomo Zanella. The barker proposal: Combining robustness and ef-
612 ficiency in gradient-based mcmc. *Journal of the Royal Statistical Society Series B: Statistical*
Methodology, 84(2):496–523, 2022.
- 613
- 614 Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the
615 ratios of the data distribution. In *Proceedings of the 41st International Conference on Machine*
616 *Learning*, pp. 32819–32848, 2024.
- 617
- 618 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan
619 Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data.
arXiv preprint arXiv:2406.03736, 2024.
- 620
- 621 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
622 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 623
- 624 Yinuo Ren, Haoxuan Chen, Grant M Rotskoff, and Lexing Ying. How discrete and continuous
625 diffusion meet: Comprehensive analysis of discrete diffusion models via a stochastic integral
626 framework. In *The Thirteenth International Conference on Learning Representations*, 2025.
- 627
- 628 Subham Sekhar Sahoo, Justin Deschenaux, Aaron Gokaslan, Guanghan Wang, Justin Chiu, and
Volodymyr Kuleshov. The diffusion duality. *arXiv preprint arXiv:2506.10892*, 2025.
- 629
- 630 Anirban Sarkar, Yijie Kang, Nirali Somia, Pablo Mantilla, Jessica Lu Zhou, Masayuki Nagai, Ziqi
631 Tang, Chris Zhao, and Peter Koo. Designing dna with tunable regulatory activity using score-
632 entropy discrete diffusion. *bioRxiv*, pp. 2024–05, 2024.
- 633
- 634 Neta Shaul, Itai Gat, Marton Havasi, Daniel Severo, Anuroop Sriram, Peter Holderrieth, Brian Kar-
635 rer, Yaron Lipman, and Ricky TQ Chen. Flow matching with general discrete paths: A kinetic-
optimal perspective. *arXiv preprint arXiv:2412.03487*, 2024.
- 636
- 637 Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of
638 diffusion models. *Advances in Neural Information Processing Systems*, 36:4263–4276, 2023.
- 639
- 640 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
preprint arXiv:2010.02502, 2020.
- 641
- 642 Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pas-
643 cal Frossard. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint*
644 *arXiv:2209.14734*, 2022.
- 645
- 646 Chenyu Wang, Masatoshi Uehara, Yichun He, Amy Wang, Tommaso Biancalani, Avantika Lal,
647 Tommi Jaakkola, Sergey Levine, Hanchen Wang, and Aviv Regev. Fine-tuning discrete diffu-
sion models via reward optimization with applications to dna and protein design. *arXiv preprint*
arXiv:2410.13643, 2024.

648 Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu, Shizhe Diao, Ligeng Zhu, Ping Luo, Song
649 Han, and Enze Xie. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache
650 and parallel decoding. *arXiv preprint arXiv:2505.22618*, 2025.
651

652 Tianyu Xie, Shuchen Xue, Zijin Feng, Tianyang Hu, Jiacheng Sun, Zhenguo Li, and Cheng Zhang.
653 Variational autoencoding discrete diffusion with enhanced dimensional correlations modeling.
654 *arXiv preprint arXiv:2505.17384*, 2025.

655 Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu.
656 Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on*
657 *Audio, Speech, and Language Processing*, 31:1720–1733, 2023.
658

659 Jiacheng Ye, Jiahui Gao, Shansan Gong, Lin Zheng, Xin Jiang, Zhenguo Li, and Lingpeng Kong.
660 Beyond autoregression: Discrete diffusion for complex reasoning and planning. *arXiv preprint*
661 *arXiv:2410.14157*, 2024.

662 Lu Yu and Arnak Dalalyan. Parallelized midpoint randomization for langevin monte carlo. *Stochastic*
663 *Processes and their Applications*, pp. 104764, 2025.
664

665 Leo Zhang. The cosine schedule is fisher-rao-optimal for masked discrete diffusion models. *arXiv*
666 *preprint arXiv:2508.04884*, 2025.

667 Lingxiao Zhao, Xueying Ding, Lijun Yu, and Leman Akoglu. Unified discrete diffusion for categor-
668 ical data. *arXiv preprint arXiv:2402.03701*, 2024.
669

670 Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. A reparameterized discrete diffusion model
671 for text generation. *arXiv preprint arXiv:2302.05737*, 2023.

672 Huanjian Zhou and Masashi Sugiyama. Parallel simulation for sampling under isoperimetry and
673 score-based diffusion models. *arXiv preprint arXiv:2412.07435*, 2024.
674

675 Huanjian Zhou and Masashi Sugiyama. Parallel simulation for log-concave sampling and score-
676 based diffusion models. In *Forty-second International Conference on Machine Learning*, 2025.

677 Yuanzhi Zhu, Xi Wang, Stephane Lathuiliere, and Vicky Kalogeiton. Dimo: Distilling masked
678 diffusion models into one-step generator. *arXiv preprint arXiv:2503.15457*, 2025.
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702	APPENDICES	
703		
704		
705	A LLM Usage Statement	14
706		
707		
708	B Related Work	15
709	B.1 Discrete Diffusion models	15
710	B.2 Acceleration for Discrete Diffusion Sampling	15
711	B.3 Picard Iteration	15
712		
713		
714		
715	C Backgrounds	15
716	C.1 Poisson Random Measure	15
717	C.2 Continuous Diffusion Models and Picard Iteration	16
718		
719		
720		
721	D Missing proof in Section 3	17
722	D.1 Useful facts	17
723	D.2 Stochastic Integral Formulation of Parallel τ -leaping	17
724	D.3 Decomposition of KL divergence	18
725	D.4 Discretization Error and Estimation Error	19
726	D.5 Convergence of Picard iteration	19
727	D.6 Proof of Theorem 3.1	21
728		
729		
730		
731		
732	E Missing Proof in Section 4	21
733	E.1 Proof of Theorem 4.1: Analysis of Algorithm 2	21
734	E.2 Proof of Proposition 4.2: Analysis of Algorithm 3	22
735		
736		
737		
738	F Experiments	23
739	F.1 Experiment Description	23
740	F.1.1 Chessboard	23
741	F.1.2 Circle	24
742	F.1.3 Hypercube	24
743	F.1.4 Embedded Hypercube	24
744	F.2 Computational Cost Analysis	25
745		
746		
747		
748		
749		
750	A LLM USAGE STATEMENT	
751		
752		

753 In preparing this manuscript, we used large language models (LLM) solely as a polishing tool to
754 improve the clarity and readability of the text. The LLM was not used for generating original ideas,
755 content, or experimental results. All conceptual contributions, analyses, and conclusions presented
in this work are entirely from the authors.

B RELATED WORK

B.1 DISCRETE DIFFUSION MODELS

Masked and Uniformed discrete diffusion models have developed from a wide range of aspects. [Austin et al. \(2021\)](#) provides foundational formalisms for discrete diffusion—multi-nomial corruption and structured transition matrices as the absorbing states. [Chao et al. \(2025\)](#) introduce intermediate token states between masked/unmasked to avoid redundant computation when sequences barely change across steps. [Vignac et al. \(2022\)](#) performs discrete denoising on graphs by noising/denoising categorical node and edge types. [Gu et al. \(2022\)](#) code sequences with discrete diffusion for text-to-image, improving quality and speed versus auto-regressive token decoders. [Gruber et al. \(2023\)](#) introduces NOS guidance to design protein sequences directly in sequence space, demonstrating antibody optimization in vitro. [Wang et al. \(2024\)](#) poptimizes discrete diffusion generators with task rewards to design biomolecular sequences.

B.2 ACCELERATION FOR DISCRETE DIFFUSION SAMPLING

So far, there have been numerous studies on accelerating sampling for discrete diffusion models. [Chen et al. \(2024b\)](#) replaces the standard Markov chain with a non-Markov schedule to skip steps and cut the number of network calls without retraining. [Sahoo et al. \(2025\)](#) adapts consistency distillation to discrete diffusion by constructing the duality connection between continuous and discrete diffusion. [Zhu et al. \(2025\)](#) proposes a token initialization strategy that injects randomness while maintaining similarity to teacher training distribution, achieving one-step distillation of masked diffusion models. [Zheng et al. \(2023\)](#); [Ou et al. \(2024\)](#) design equivalent reparameterization of discrete diffusion that yields more effective training and decoding strategies. [Wu et al. \(2025\)](#) develops confidence-aware parallel decoding to accelerate multi-token sampling while maintaining accuracy. [Zhang \(2025\)](#) gives a principled choice of discretization schedule for efficient sampling. [Shaul et al. \(2024\)](#) allows arbitrary discrete probability paths, giving more control to find shorter or easier trajectories with fewer steps for discrete generation. [Zhao et al. \(2024\)](#) derives a simple backward denoising formula, enabling exact and accelerated sampling and unifying discrete-time and continuous-time discrete diffusion.

B.3 PICARD ITERATION

Parallel sampling based on Picard recursions has been applied to various models such as first-order Markov chains for both Monte Carlo methods ([Yu & Dalalyan, 2025](#); [Anari et al., 2024](#)) and continuous generative models ([Zhou & Sugiyama, 2024](#); [Shih et al., 2023](#); [Gupta et al., 2024](#)).

C BACKGROUNDS

In this section, we introduce the basic concepts and existing convergence analysis results of discrete diffusion models in Section 2 and parallel simulation for diffusion model based on Picard iteration in Section C.2.

C.1 POISSON RANDOM MEASURE

Definition C.1 (Poisson Random Measure). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and $(\mathbb{X}, \mathcal{B}, \nu)$ be a measure space satisfying that

$$\int_{\mathbb{X}} 1 \vee |y| \vee |y|^2 \nu(dy) < \infty,$$

The random measure $N(dt, dy)$ on $\mathbb{R}^+ \times \mathbb{X}$ is called a *Poisson random measure* w.r.t. measure ν if it is a random counting measure satisfying the following properties:

- (i) For any $B \in \mathcal{B}$ and $0 \leq s < t$, $N((s, t] \times B) \sim \mathcal{P}(\nu(B)(t - s))$;
- (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.

Definition C.2 (Poisson Random Measure with Evolving Intensity). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and $(\mathbb{X}, \mathcal{B}, \nu)$ be a measure space. Suppose $\lambda_t(y)$ is a non-negative predictable process on $\mathbb{R}^+ \times \mathbb{X} \times \Omega$ satisfying that for any $0 \leq T < \bar{T}$,

$$\int_0^T \int_{\mathbb{X}} 1 \vee |y| \vee |y|^2 \lambda_t(y) \nu(dy) dt < \infty, \text{ a.s..}$$

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

(i) For any $B \in \mathcal{B}$ and $0 \leq s < t$, $N[\lambda]((s, t] \times B) \sim \mathcal{P}\left(\int_s^t \int_B \lambda_\tau(y) \nu(dy) d\tau\right)$;

(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.

C.2 CONTINUOUS DIFFUSION MODELS AND PICARD ITERATION

In score-based diffusion models, one considers forward process $(\mathbf{x}_t)_{t \in [0, T]}$ in \mathbb{R}^d governed by the canonical Ornstein-Uhlenbeck (OU) process (Ledoux, 2000):

$$d\mathbf{x}_t = -\frac{1}{2}\mathbf{x}_t dt + d\mathbf{B}_t, \quad \mathbf{x}_0 \sim \mathbf{q}_0, \quad t \in [0, T], \quad (11)$$

where \mathbf{q}_0 is the initial distribution over \mathbb{R}^d and \mathbf{B}_t is the Brownian Motion. The corresponding backward process $(\tilde{\mathbf{x}}_t)_{t \in [0, T]}$ in \mathbb{R}^d follows an SDE defined as

$$\begin{cases} d\tilde{\mathbf{x}}_t = \left[\frac{1}{2}\tilde{\mathbf{x}}_t + \nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t)\right] dt + d\mathbf{B}_t & t \in [0, T], \\ \tilde{\mathbf{x}}_0 \sim \mathbf{p}_0 \approx \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d) \end{cases} \quad (12)$$

where $\mathcal{N}(\cdot, \cdot)$ represents the normal distribution over \mathbb{R}^d . In practice, the score function $\nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t)$ is estimated by NN $\mathbf{s}_t^\theta : \mathbb{R}^d \mapsto \mathbb{R}^d$, where θ is the parameters of NN. The backward process is approximated by

$$\begin{cases} d\mathbf{y}_t = \left[\frac{1}{2}\mathbf{y}_t + \mathbf{s}_t^\theta(\mathbf{y}_t)\right] dt + d\mathbf{B}_t & t \in [0, T], \\ \mathbf{y}_0 \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d). \end{cases} \quad (13)$$

The main idea of parallelized simulation via Picard iteration is to regroup the discrete grids along time horizon and update all grids in same group simultaneously (Clenshaw, 1957; Anari et al., 2024; Zhou & Sugiyama, 2025). Specifically, to approximate the difference $\mathbf{x}_{t_{n+1}} - \mathbf{x}_{t_n}$ over time slice $[t_n, t_{n+1}]$ as

$$\begin{aligned} \mathbf{x}_{t_{n+1}} - \mathbf{x}_{t_n} &= \int_{t_n}^{t_{n+1}} \mathbf{s}_s^\theta(\mathbf{x}_s) ds + \sqrt{2}(\mathbf{B}_{t_{n+1}} - \mathbf{B}_{t_n}) \\ &\approx \sum_{i=1}^M \mathbf{w}_i \mathbf{s}_{t_n + \tau_{n,i}}^\theta(\mathbf{x}_{t_n + \tau_{n,i}}) + \sqrt{2}(\mathbf{B}_{t_{n+1}} - \mathbf{B}_{t_n}), \end{aligned}$$

with a discrete grid of M collocation points as $\mathbf{x}_{t_n} = \mathbf{x}_{t_n + \tau_{n,0}} \leq \mathbf{x}_{t_n + \tau_{n,1}} \leq \mathbf{x}_{t_n + \tau_{n,2}} \leq \dots \leq \mathbf{x}_{t_n + \tau_{n,M}} = \mathbf{x}_{t_{n+1}}$. We update the points in a wave-like fashion, which inherently allows for parallelization: for $m' = 1, \dots, M$,

$$\mathbf{x}_{t_n + \tau_{n,m}}^{p+1} = \mathbf{x}_{t,n} + \sum_{m=1}^{M-1} \mathbf{w}_m \mathbf{s}_{t_n + \tau_{n,i}}^\theta(\mathbf{x}_{t_n + \tau_{n,m}}^p) + \sqrt{2}(\mathbf{B}_{t_n + \tau_{n,m}} - \mathbf{B}_{t_n}).$$

With such regrouping, as long as the total time length of each group scales as $\mathcal{O}(1/L)$, the grids will converge exponentially fast. Given a sufficiently accurate starting point at time t_n , the initial error scales as $\mathcal{O}(d)$. Therefore, $K = \mathcal{O}\left(\log\left(\frac{d}{\varepsilon}\right)\right)$ steps suffice for the convergence of each group and resulting $\log d/d$ times speed-up.

864 D MISSING PROOF IN SECTION 3

865
866 In this section, we provide the detailed proof about the error bound of Picard-based parallel sam-
867 pling.
868

869 D.1 USEFUL FACTS

870
871 **Definition D.1 (Mixing Time).** We define the mixing time $t_{\text{mix}}(\epsilon)$ of the continuous-time Markov
872 chain with rate matrix \mathbf{Q} as the smallest time t such that starting from any initial distribution \mathbf{p}_0 , the
873 KL divergence $D_{\text{KL}}(\mathbf{p}_t \parallel \boldsymbol{\pi})$ is less than ϵ , i.e.

$$874 t_{\text{mix}}(\epsilon) = \inf \left\{ t \in \mathbb{R}_+ \mid D_{\text{KL}}(\mathbf{p}_t \parallel \boldsymbol{\pi}) = D_{\text{KL}}(e^{-t\mathbf{Q}}\mathbf{p}_0 \parallel \boldsymbol{\pi}) \leq \epsilon \right\}.$$

875
876
877 **Theorem D.1 (Truncation Error Theorem C.1 in Ren et al. (2025)).** The forward process equa-
878 tion 1 converges to the uniform distribution $\mathbf{p}_\infty = \mathbf{1}/|\mathbb{X}|$ exponentially fast in terms of the KL
879 divergence, i.e.

$$880 D_{\text{KL}}(\mathbf{p}_t \parallel \mathbf{p}_\infty) = D_{\text{KL}}\left(\mathbf{p}_t \parallel \frac{\mathbf{1}}{|\mathbb{X}|}\right) \lesssim e^{-\rho t} \log |\mathbb{X}|,$$

881 where $|\mathbb{X}|$ is the size of the state space, and t_{mix} is the mixing time of the continuous-time Markov
882 chain corresponding to the rate matrix \mathbf{Q} defined in Definition D.1.
883
884

885 D.2 STOCHASTIC INTEGRAL FORMULATION OF PARALLEL τ -LEAPING

886
887 **Proposition D.1 (Stochastic Integral Formulation of backward process (Ren et al., 2025, Propo-
888 sition 3.2)).** The backward process in discrete diffusion models equation 2 can be represented by
889 the following stochastic integral:
890

$$891 \tilde{x}_t = \tilde{x}_0 + \int_0^t \int_{\mathbb{X}} (y - \tilde{x}_{\tau-}) N[\mu](d\tau, dy), \quad (14)$$

892 where

$$893 \mu_\tau(y) = \tilde{s}_\tau(\tilde{x}_{\tau-}, y) \tilde{Q}_\tau(\tilde{x}_{\tau-}, y),$$

894 and X_{t-} denotes the left limit of a càdlàg process X_t at time t .

895
896 **Proposition D.2 (Stochastic Integral Formulation of parallel τ -Leaping).** In time $[t_n, t_{n+1}]$, for
897 any $k \in [K]$, $\hat{y}_{t_n+m\epsilon}$ can be interpreted by following stochastic integral equation:
898

$$899 \hat{y}_t^K = \hat{y}_{t_n} + \int_0^t \int_{\mathbb{X}} (y - \hat{y}_{[\tau]^-}^{K-1}) N[\hat{\mu}_{[\cdot]}^{K-1}](d\tau, dy), \quad (15)$$

900 where the evolving intensity $\hat{\mu}_{[\tau]}^{K-1}(y)$ is given by

$$901 \hat{\mu}_{[\tau]}^{K-1}(y) = \hat{s}_{[\tau]}(\hat{y}_{[\tau]^-}^{K-1}, y) \tilde{Q}_{[\tau]}(\hat{y}_{[\tau]^-}^{K-1}, y),$$

902 in which we used the symbol $[\tau] = \min\{t_n + l\epsilon : l \in [N]\}$ for $\tau \in [t_n, t_{n+1})$.

903
904 *Proof.* Without loss of generality, suppose the time block $[t_n, t_{n+1}]$ is uniformly divided into M
905 small steps with grid $(s_i)_{i \in [0:M]}$, then we have

$$906 \begin{aligned} 907 \hat{y}_t^K &= \hat{y}_{t_n} + \int_0^s \int_{\mathbb{X}} (y - \hat{y}_{[\tau]^-}^{K-1}) N[\hat{\mu}_{[\cdot]}^{K-1}](d\tau, dy) \\ 908 &= \hat{y}_{t_n} + \sum_{i=1}^{\lfloor (t-t_n)/\epsilon \rfloor} \int_{s_{i-1}}^{s_i} \int_{\mathbb{X}} \left(y - \hat{y}_{s_{i-1}^-}^{K-1} \right) N[\hat{\mu}_{s_{i-1}}^{K-1}](d\tau, dy) \\ 909 &= \hat{y}_{t_n} + \sum_{i=1}^{\lfloor (t-t_n)/\epsilon \rfloor} \int_{\mathbb{X}} \left(y - \hat{y}_{s_{i-1}^-}^{K-1} \right) N[\hat{\mu}_{s_{i-1}}^{K-1}]((s_{i-1}, s_i], dy), \end{aligned} \quad (16)$$

since \mathbb{X} is finite, it can be decomposed into the following sum of jumps:

$$\begin{aligned}
\hat{y}_t^K &= \hat{y}_{t_n} + \sum_{i=1}^{\lfloor (t-t_n)/\epsilon \rfloor} \int_{\mathbb{X}} \left(y - \hat{y}_{s_{i-1}}^{K-1} \right) N[\hat{\mu}_{s_{i-1}}^{K-1}]((s_{i-1}, s_i], dy) \\
&= \hat{y}_{t_n} + \sum_{i=1}^{\lfloor (t-t_n)/\epsilon \rfloor} \sum_{y \in \mathbb{X}} \left(y - \hat{y}_{s_{i-1}}^{K-1} \right) N[\hat{\mu}_{s_{i-1}}^{K-1}]((s_{i-1}, s_i], \{y\}) \\
&= \hat{y}_{t_n} + \sum_{i=1}^{\lfloor (t-t_n)/\epsilon \rfloor} \sum_{y \in \mathbb{X}} \left(y - \hat{y}_{s_{i-1}}^{K-1} \right) \mathcal{P}((s_i - s_{i-1})\hat{\mu}_{s_{i-1}}(y)), \tag{17}
\end{aligned}$$

which is exactly the Picard update formula (7) in Section 3.1. \square

D.3 DECOMPOSITION OF KL DIVERGENCE

Theorem D.2. *Let $\tilde{p}_{0:T}$ and $\hat{q}_{0:T}$ be the path measures of the backward process with the stochastic integral formulation equation 14 and the interpolating process equation 15 of τ -leaping algorithm, then it holds that*

$$D_{\text{KL}}(\tilde{p}_{0:T} \parallel \hat{q}_{0:T}) \leq D_{\text{KL}}(\tilde{p}_0 \parallel \hat{q}_0) + \mathbb{E} \left[\int_0^T \int_{\mathbb{X}} \left(\mu_t(y) \log \frac{\mu_t(y)}{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)} - \mu_t(y) + \hat{\mu}_{\lfloor t \rfloor}^{K-1}(y) \right) \nu(dy) dt \right], \tag{18}$$

where the expectation is taken w.r.t. paths generated by the backward process equation 14.

Proof. By the chain rule of KL divergence, we have

$$D_{\text{KL}}(\tilde{p}_{0:T} \parallel \hat{q}_{0:T}) = D_{\text{KL}}(\tilde{p}_0 \parallel \hat{q}_0) + \mathbb{E}[D_{\text{KL}}(\tilde{p}_{0:T} \parallel \hat{q}_{0:T} | \tilde{x}_0 = y_0 = y)]. \tag{19}$$

One can rewrite the second term as

$$\mathbb{E}[D_{\text{KL}}(\tilde{p}_{0:T} \parallel \hat{q}_{0:T} | \tilde{x}_0 = y_0 = y)] = \mathbb{E} \left[\log \frac{d\tilde{p}_{0:T}}{d\hat{q}_{0:T}} \mid \tilde{x}_0 = y_0 = y \right] = \mathbb{E} \left[\log Z_T^{-1} \left[\frac{\hat{\mu}^{K-1}}{\mu} \right] \right]. \tag{20}$$

where Z_T denotes the Change of Measure in Theorem 3.3 in (Ren et al., 2025) with the form of

$$\log Z_T[h] = \int_0^t \int_{\mathbb{X}} \log h_\tau(y) N[\lambda](d\tau \times dy) - \int_0^t \int_{\mathbb{X}} (h_\tau(y) - 1) \lambda_\tau(y) \nu(dy) d\tau. \tag{21}$$

Given $h = \frac{\hat{\mu}^{K-1}}{\mu}$, we have

$$\begin{aligned}
\log \left(Z_T \left[\frac{\hat{\mu}^{K-1}}{\mu} \right] \right) &= \int_0^T \int_{\mathbb{X}} \log \left(\frac{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)}{\mu_t(y)} \right) N[\mu](dt \times dy) \\
&\quad - \int_0^T \int_{\mathbb{X}} \left(\frac{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)}{\mu_t(y)} - 1 \right) \mu_s(y) \nu(dy) dt
\end{aligned}$$

Eventually, we can compute the expectation as

$$\begin{aligned}
\mathbb{E} \left[\log Z_T^{-1} \left[\frac{\hat{\mu}^{K-1}}{\mu} \right] \right] &= \mathbb{E} \left[\int_0^T \int_{\mathbb{X}} \left(-\log \left(\frac{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)}{\mu_t(y)} \right) + \frac{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)}{\mu_t(y)} - 1 \right) \mu_t \nu(dy) dt \right] \\
&= \mathbb{E} \left[\int_0^T \int_{\mathbb{X}} \left(\mu_t(y) \log \frac{\mu_t(y)}{\hat{\mu}_{\lfloor t \rfloor}^{K-1}(y)} - \mu_t(y) + \hat{\mu}_{\lfloor t \rfloor}^{K-1}(y) \right) \nu(dy) dt \right]
\end{aligned}$$

\square

Now, the problem remaining is reduced to bound the following discrepancy with $G(x; y) = x(\log x - \log y) - x$,

$$\begin{aligned}
& \int_{t_n}^{t_{n+1}} \int_{\mathbb{X}} \left(\mu_t(y) \log \frac{\mu_t(y)}{\widehat{\mu}_{[t]}^{K-1}(y)} - \mu_t(y) + \widehat{\mu}_{[t]}^{K-1}(y) \right) \nu(dy) dt \\
&= \int_{t_n}^{t_{n+1}} \int_{\mathbb{X}} \underbrace{\left(\mu_{[t]}(y) \log \frac{\mu_{[t]}(y)}{\widehat{\mu}_{[t]}^K(y)} - \mu_t(y) + \widehat{\mu}_{[t]}^K(y) \right)}_{:=A_{t_n}} \\
&\quad + \underbrace{G(\mu_t(y); \widehat{\mu}_{[t]}^K(y)) - G(\mu_{[t]}(y); \widehat{\mu}_{[t]}^K(y))}_{:=B_{t_n}} \\
&\quad + \underbrace{\mu_t(y) \log \frac{\widehat{\mu}_{[t]}^K(y)}{\widehat{\mu}_{[t]}^{K-1}(y)} - \widehat{\mu}_{[t]}^K(y) + \widehat{\mu}_{[t]}^{K-1}(y)}_{:=C_{t_n}} \nu(dy) ds;
\end{aligned}$$

where A_{t_n} is the estimation error of score function, B_{t_n} measures the discretization error, and C_{t_n} measures the error by Picard iteration.

D.4 DISCRETIZATION ERROR AND ESTIMATION ERROR

Proposition D.3 (Discretization error and Estimation Error (Ren et al., 2025)). For any $n \in [N]$, we have

$$\int_0^T \int_{\mathbb{X}} A_{t_n} \nu(dy) ds \leq \delta^2 \quad \text{and} \quad \int_{t_n}^{t_{n+1}} \int_{\mathbb{X}} B_{t_n} \nu(dy) ds \leq \overline{D}^2 \epsilon h.$$

D.5 CONVERGENCE OF PICARD ITERATION

Proposition D.4 (Convergence of Picard iteration). By taking $h = m\epsilon = \mathcal{O}(\frac{1}{d^{3/2}})$, for any $n \in [N]$, $j \in [K]$, $y \in \mathbb{X}$ and $t \in [t_n, t_{n+1}]$, we have

$$\max_{t \in [t_n, t_{n+1}]} \mathbb{E} \left\| \widehat{y}_{[t]}^{K-j} - \widehat{y}_{[t]}^{K-j-1} \right\| \leq C \max_{t \in [t_n, t_{n+1}]} \mathbb{E} \left\| \widehat{y}_{[t]}^{K-j-1} - \widehat{y}_{[t]}^{K-j-2} \right\|.$$

with some constant $C \in (0, 1)$.

Proof. The update step inside the block is as follows:

$$\widehat{y}_m^{(k+1)} = \widehat{y}_0^{(k+1)} + \sum_{j=0}^{m-1} \Delta \widehat{y}_j^{(k)}$$

where

$$\Delta \widehat{y}_j^{(k)} = \sum_{y' \in \mathbb{X}} (y' - \widehat{y}_{t_n+j\epsilon}^{(k)}) \cdot \text{Numjump}_j^{(k)}(y'), \quad \text{Numjump}_j^{(k)}(y') \sim \mathcal{P}(\widehat{\mu}_{t_n+j\epsilon}^k(y') | \widehat{y}_{t_n+j\epsilon}^{(k)}) \cdot \epsilon.$$

Suppose the state difference between iteration $k+1$ and k in block n at time $t_n + m\epsilon$ is defined as

$$\epsilon_{k+1}(m) = \widehat{y}_{t_n+m\epsilon}^{(k+1)} - \widehat{y}_{t_n+m\epsilon}^{(k)},$$

such difference can also be considered as the difference of accumulated jump:

$$\epsilon_k(m) = \sum_{j=0}^{m-1} \left(\Delta \widehat{y}_j^{(k)} - \Delta \widehat{y}_j^{(k-1)} \right).$$

Take square norm on both sides and use Cauchy-Schwarz Inequality, we have

$$\mathbb{E}[\|\epsilon_k(m)\|^2] = \mathbb{E} \left[\left\| \sum_{j=0}^{m-1} \left(\Delta \widehat{y}_j^{(k)} - \Delta \widehat{y}_j^{(k-1)} \right) \right\|^2 \right] \leq m \sum_{j=0}^{m-1} \mathbb{E} \left[\left\| \Delta \widehat{y}_j^{(k)} - \Delta \widehat{y}_j^{(k-1)} \right\|^2 \right].$$

Now we split the difference between two total jump vectors by using $\|A + B\|^2 \leq 2\|A\|^2 + 2\|B\|^2$:

$$\mathbb{E} \left[\left\| \Delta \hat{y}_j^{(k)} - \Delta \hat{y}_j^{(k-1)} \right\|^2 \right] \leq 2 \underbrace{\mathbb{E} \left[\left\| \sum_{y' \in \mathbb{X}} (y' - \hat{y}_j^{(k)}) (\text{NumJump}_j^{(k)}(y') - \text{NumJump}_j^{(k-1)}(y')) \right\|^2 \right]}_{I_1} + 2 \underbrace{\mathbb{E} \left[\left\| (\hat{y}_j^{(k-1)} - \hat{y}_j^{(k)}) \sum_{y' \in \mathbb{X}} \text{NumJump}_j^{(k-1)}(y') \right\|^2 \right]}_{I_2}$$

For the first term, from the Cauchy-Schwarz Inequality, we have:

$$I_1 \leq \|\mathbb{X}\| \sum_{y' \in \mathbb{X}} \mathbb{E} \left[\|y' - \hat{y}_j^{(k)}\|^2 \cdot (\text{NumJump}_j^{(k)}(y') - \text{NumJump}_j^{(k-1)}(y'))^2 \right]$$

We define $R_{\mathbb{X}}^2 = \max_{x_1, x_2 \in \mathbb{X}} \|x_1 - x_2\|^2$ as the scale of the state space. From the Lipschitz Continuity in Assumption 3.4, we have

$$I_1 \leq \|\mathbb{X}\| \cdot R_{\mathbb{X}}^2 \cdot (L_\mu \epsilon)^2 \mathbb{E} [\|\hat{y}_j^{(k)} - \hat{y}_j^{(k-1)}\|^2]$$

For the second, we utilize the Law of total expectation to decompose the form as:

$$\begin{aligned} I_2 &= \mathbb{E} \left[\mathbb{E} \left[\left\| \hat{y}_j^{(k-1)} - \hat{y}_j^{(k)} \right\|^2 \left(\sum_{y' \in \mathbb{X}} \text{NumJumps}_j^{(k-1)}(y') \right)^2 \middle| \hat{y}_j^{(k-1)}, \hat{y}_j^{(k)} \right] \right] \\ &= \mathbb{E} \left[\left\| \hat{y}_j^{(k-1)} - \hat{y}_j^{(k)} \right\|^2 \mathbb{E} \left[\left(\sum_{y' \in \mathbb{X}} \text{NumJump}_j^{(k-1)}(y') \right)^2 \middle| \hat{y}_j^{(k-1)} \right] \right] \end{aligned} \quad (22)$$

For the inner expectation, the sum of NumJump is also the sum of independent Poisson variables, which is still a Poisson variable. Thus we have

$$\lambda_{total} = \sum_{y' \in \mathbb{X}} \mathbb{E} [\text{NumJump}_j^{(k-1)}(y') | \hat{y}_j^{(k-1)}] = \sum_{y' \in \mathbb{X}} \mu_j^\theta(y' | \hat{y}_j^{(k-1)}) \cdot \epsilon$$

According to the assumption of rate matrix regularization and bounded score, we have

$$\sum_{y' \neq \hat{y}_j^{(k-1)}} \mu_j^\theta(y' | \hat{y}_j^{(k-1)}) = \sum_{y' \neq \hat{y}_j^{(k-1)}} \hat{s}^\theta(\hat{y}_j^{(k-1)}, y') \tilde{Q}(\hat{y}_j^{(k-1)}, y') \leq M_s \bar{D} := \Lambda_{max}$$

For $x \sim \mathcal{P}(\lambda)$, the second momentum is $\mathbb{E}[x^2] = \lambda^2 + \lambda$, therefore we have

$$\mathbb{E} \left[\left(\sum_{y' \in \mathbb{X}} \text{NumJump}_j^{(k-1)}(y') \right)^2 \middle| \hat{y}_j^{(k-1)} \right] \leq (\epsilon \Lambda_{max})^2 + (\epsilon \Lambda_{max}),$$

and the second term is bounded as

$$I_2 \leq 2 \mathbb{E} [\|\hat{y}_j^{(k-1)} - \hat{y}_j^{(k)}\|^2] \cdot (\epsilon^2 \Lambda_{max}^2 + \epsilon \Lambda_{max}). \quad (23)$$

Overall, we have

$$\begin{aligned} \mathbb{E} \left[\left\| \Delta \hat{y}_j^{(k)} - \Delta \hat{y}_j^{(k-1)} \right\|^2 \right] &\leq (\|\mathbb{X}\| \cdot R_{\mathbb{X}}^2 \cdot (L_\mu \epsilon)^2 + \epsilon^2 \Lambda_{max}^2 + \epsilon \Lambda_{max}) \cdot \mathbb{E} [\|\hat{y}_j^{(k-1)} - \hat{y}_j^{(k)}\|^2] \\ &= C_0 \cdot \mathbb{E} [\|\hat{y}_j^{(k-1)} - \hat{y}_j^{(k)}\|^2], \end{aligned} \quad (24)$$

with $C_0 = \|\mathbb{X}\| \cdot R_{\mathbb{X}}^2 \cdot (L_\mu \epsilon)^2 + \epsilon^2 \Lambda_{max}^2 + \epsilon \Lambda_{max}$. Thus, we have

$$\begin{aligned} \max_{t \in [t_n, t_{n+1}]} \mathbb{E} \left\| \hat{y}_{[t]}^{K-j} - \hat{y}_{[t]}^{K-j-1} \right\| &\leq m \sum_{j=0}^{m-1} \mathbb{E} \left[\left\| \Delta \hat{y}_j^{K-j-1} - \Delta \hat{y}_j^{K-j-2} \right\|^2 \right] \\ &\leq m \cdot C_0 \cdot \sum_{j=0}^{m-1} \mathbb{E} [\|\hat{y}_j^{K-j-1} - \hat{y}_j^{K-j-2}\|^2] \end{aligned} \quad (25)$$

$$\leq m^2 C_0 \mathbb{E} [\|\hat{y}_j^{K-j-1} - \hat{y}_j^{K-j-2}\|^2] \quad (26)$$

$$\leq C \max_{t \in [t_n, t_{n+1}]} \mathbb{E} \left\| \hat{y}_{[t]}^K - \hat{y}_{[t]}^{K-j} \right\|, \quad (27)$$

the last inequality holds if we take $h = m\epsilon = \mathcal{O}(\frac{\delta}{j^{3/2}})$. \square

1080 D.6 PROOF OF THEOREM 3.1

1081
1082 Combining Theorem D.1, Lemma D.2, Proposition D.3, and Proposition D.4, we have

$$\begin{aligned}
1083 D_{\text{KL}}(\bar{p}_{0:T} \|\widehat{q}_{0:T}) &\leq D_{\text{KL}}(\bar{p}_0 \|\widehat{q}_0) + \mathbb{E} \left[\int_0^T \int_{\mathbb{X}} \left(\mu_t(y) \log \frac{\mu_t(y)}{\widehat{\mu}_{[t]}^{K-1}(y)} - \mu_t(y) + \widehat{\mu}_{[t]}^{K-1}(y) \right) \nu(dy) dt \right] \\
1084 &\leq D_{\text{KL}}(\bar{p}_0 \|\widehat{q}_0) + \sum_{n=1}^N \mathbb{E} \int_{t_n}^{t_{n+1}} \int_{\mathbb{X}} A_{t_n} + B_{t_n} + C_{t_n} \nu(dy) ds \\
1085 &\leq \lesssim e^{-\rho T} \log |\mathbb{X}| + \delta^2 + \bar{D}^2 \epsilon T + C^K \text{poly}(\bar{D}).
\end{aligned}$$

1092 E MISSING PROOF IN SECTION 4

1094 E.1 PROOF OF THEOREM 4.1: ANALYSIS OF ALGORITHM 2

1095
1096 Suppose \bar{p}_{t_n} is the true target distribution derived from the forward process, and \widehat{q}_{t_n} is the generated distribution from our predictor-corrector algorithm. We use the total variance distance $E_n = TV(\widehat{q}_{t_n}, \bar{p}_{t_n})$ as the measurement of the error before entering the n -th block computation.

1099 The predictor generates $\widehat{q}_{\text{pred}}$ from the initial distribution \widehat{q}_{t_n} . Suppose the predictor also generates a $\widetilde{q}_{\text{pred}}$ from the true distribution \bar{p}_{t_n} , then the error generated by the predictor $TV(\widehat{q}_{\text{pred}}, \bar{p}_{t_{n+1}})$ can be decomposed by the trigonometric inequality as

$$1102 TV(\widehat{q}_{\text{pred}}, \bar{p}_{t_{n+1}}) \leq TV(\widehat{q}_{\text{pred}}, \widetilde{q}_{\text{pred}}) + TV(\widetilde{q}_{\text{pred}}, \bar{p}_{t_{n+1}}). \quad (28)$$

1104 The first term can be considered as the propagation of the error. Since $\widehat{q}_{\text{pred}}$ and $\widetilde{q}_{\text{pred}}$ are generated by the same mapping (parallel τ -leaping) from the initial $(\widehat{q}_{t_n}, \bar{p}_{t_n})$, according to the Data Processing Inequality [], we have

$$1107 TV(\widehat{q}_{\text{pred}}, \widetilde{q}_{\text{pred}}) \leq TV(\widehat{q}_{t_n}, \bar{p}_{t_n}) = E_n, \quad (29)$$

1109 which means the predictor will not amplify the previous error.

1110 The second term $TV(\widetilde{q}_{\text{pred}}, \bar{p}_{t_{n+1}})$ means the new error introduced by the predictor, which consists of the discretization error from τ -leaping and Picard iteration, and the approximation error from the network learning. We denote this term as Δ_p . Therefore the bound of the total error after the predictor stage is:

$$1114 E_{\text{pred}} \leq E_n + \Delta_p. \quad (30)$$

1116 The corrector use $\widehat{q}_{\text{pred}}$ as the initial distribution, and generate $\widehat{q}_{t_{n+1}}$ after running K_c steps of Metropolis-Hastings MCMC process, which is designed with the target distribution $\bar{p}_{t_{n+1}}$ as the only stationary distribution. Since the transfer kernel of the MCMC is a contraction mapping with respect to the Total Variance distance with the factor γ , after K_c steps, we have

$$1120 E_{n+1} = TV(\widehat{q}_{t_{n+1}}, \bar{p}_{t_{n+1}}) \leq \gamma^{K_c} E_{\text{pred}} + \Delta_c = \gamma^{K_c} (E_n + \Delta_p) + \Delta_c, \quad (31)$$

1122 where Δ_c is the small approximation error of the corrector itself, since it also used \mathfrak{s}_t^θ for computing acceptance rate. Eventually, the total error after N blocks can be computed as

$$1125 E_N \leq \frac{\gamma^{K_c} \Delta_p + \Delta_c}{1 - \gamma^{K_c}} + \gamma^{N K_c} E_0 \quad (32)$$

1127 The estimation of the contract factor γ is described in the following Proposition.

1128 **Proposition E.1** (The Estimation of MCMC Contraction). *In finite state space Ω , suppose there is a reversible MCMC transfer kernel P with non-negative eigenvalues and stationary distribution π . Denote the second largest eigenvalues of P as λ_* and Φ as the conductance of the graph corresponding to P , for any distribution μ , we have*

$$1132 \|\mu P^t - \pi\|_{TV} \leq \frac{1}{2} \sqrt{\frac{1}{\pi_{\min}} - 1} \left(1 - \frac{\Phi^2}{2}\right)^t \quad (33)$$

1134 *Proof.* Let $f := \frac{d\mu}{d\pi}$ so that $\mu = f\pi$ and $\mu - \pi = (f - 1)\pi$. First, by Cauchy–Schwarz,

$$1135 \quad \|\mu - \pi\|_{\text{TV}} = \frac{1}{2} \sum_x \pi(x) |f(x) - 1| = \frac{1}{2} \|f - 1\|_{1,\pi} \leq \frac{1}{2} \|f - 1\|_{2,\pi}.$$

1138 Reversibility implies P is self-adjoint on $L^2(\pi)$ and (using detailed balance)

$$1140 \quad \frac{d(\mu P^t)}{d\pi}(y) = \sum_x f(x) \frac{\pi(x)}{\pi(y)} P^t(x, y) = \sum_x f(x) P^t(y, x) = (P^t f)(y),$$

1143 hence $\mu P^t - \pi = (P^t(f - 1))\pi$ and therefore

$$1144 \quad \|\mu P^t - \pi\|_{\text{TV}} = \frac{1}{2} \|P^t(f - 1)\|_{1,\pi} \leq \frac{1}{2} \|P^t(f - 1)\|_{2,\pi}. \quad (34)$$

1146 Next, since the chain is reversible, the spectrum of P lies in $[0, 1]$ and there is an orthonormal eigenbasis of $L_0^2(\pi) := \{g : \sum_x g(x)\pi(x) = 0\}$. Expanding $g \in L_0^2(\pi)$ in that basis gives

$$1149 \quad \|P^t g\|_{2,\pi} \leq \lambda_*^t \|g\|_{2,\pi} = (1 - \gamma)^t \|g\|_{2,\pi}.$$

1150 Applying this with $g = f - 1$ and combining with (equation 34) yields

$$1151 \quad \|\mu P^t - \pi\|_{\text{TV}} \leq \frac{1}{2} (1 - \gamma)^t \|f - 1\|_{2,\pi}. \quad (35)$$

1153 We now bound the initial L^2 deviation:

$$1155 \quad \|f - 1\|_{2,\pi}^2 = \sum_x \pi(x) \left(\frac{\mu(x)}{\pi(x)} - 1 \right)^2 = \sum_x \frac{\mu(x)^2}{\pi(x)} - 1 \leq \frac{1}{\pi_{\min}} \sum_x \mu(x)^2 - 1 \leq \frac{1}{\pi_{\min}} - 1,$$

1158 where we used $\pi(x) \geq \pi_{\min}$ and $\mu(x)^2 \leq \mu(x)$ for probabilities. Hence

$$1159 \quad \|f - 1\|_{2,\pi} \leq \sqrt{\frac{1}{\pi_{\min}} - 1}. \quad (36)$$

1162 Finally, Cheeger’s inequality for lazy reversible chains states

$$1164 \quad \gamma \geq \frac{\Phi^2}{2},$$

1166 so $(1 - \gamma)^t \leq \left(1 - \frac{\Phi^2}{2}\right)^t$. Plugging this and (equation 36) into (equation 35) gives

$$1168 \quad \|\mu P^t - \pi\|_{\text{TV}} \leq \frac{1}{2} \sqrt{\frac{1}{\pi_{\min}} - 1} \left(1 - \frac{\Phi^2}{2}\right)^t.$$

1170 □

1172 E.2 PROOF OF PROPOSITION 4.2: ANALYSIS OF ALGORITHM 3

1174 Fix $x \in \mathbb{X}$ and a token b with $b \sim q_i(\cdot|x) > 0$, and set $x' = x^{(i \leftarrow b)}$ and $a = x_i$. Use the Barker

1175 acceptance with the standard proposal-ratio correction:

$$1176 \quad \alpha_i(x \rightarrow x') = \begin{cases} \frac{r_i(x, x')}{1 + r_i(x, x')}, & \text{if } q_i(a|x') > 0, \\ 0, & \text{otherwise,} \end{cases} \quad \text{where } r_i(x, x') = \frac{\pi_t(x') q_i(a|x')}{\pi_t(x) q_i(b|x)}.$$

1180 Setting $\alpha_i = 0$ when the reverse proposal has zero support will enforce symmetric support and

1181 avoids irreversibility. Such a setting defines a single-coordinate Markov kernel K_i as

$$1182 \quad K_i(x, x') = \begin{cases} q_i(b|x) \alpha_i(x \rightarrow x'), & \text{if } x' = x^{(i \leftarrow b)} \neq x, \\ 1 - \sum_b q_i(b|x) \alpha_i(x \rightarrow x^{(i \leftarrow b)}), & \text{if } x' = x, \\ 0, & \text{others.} \end{cases}$$

1187 We first prove two following lemma.

Lemma E.1 (Anchor Cancellation). For any $x, i, \tilde{x}^{(i)} \equiv x^{(i \leftarrow M')}$, and $a, b \in \mathcal{V}$, we have

$$\log \frac{\pi_t(x^{(i \leftarrow b)})}{\pi_t(x)} = \log \underbrace{\frac{\pi_t(x^{(i \leftarrow b)})}{\pi_t(x^{(i \leftarrow M')})}_{u_{i,b}(\tilde{x}^{(i)})} - \log \underbrace{\frac{\pi_t(x)}{\pi_t(x^{(i \leftarrow M')})}_{u_{i,a}(\tilde{x}^{(i)})} = u_{i,b}(\tilde{x}^{(i)}) - u_{i,a}(\tilde{x}^{(i)}).$$

Lemma E.2 (Proposal on Anchored Context). If $q_i(\cdot|x)$ is defined on $\mathcal{K}_i(\tilde{x}^{(i)})$, then for $x' = x^{(i \leftarrow b)}$, we have

$$q_i(\cdot|x') \equiv q_i(\cdot|x)$$

whenever both are defined, because $\tilde{x}^{(i)}$ is identical for x and x' . In particular, if q_i is *uniform* on a fixed-size set \mathcal{K}_i , then

$$q_i(b|x) = q_i(a|x') = \frac{1}{K}$$

whenever both are positive.

Now we prove the main result. Consider any $x' \neq x$ such that $K_i(x, x') > 0$; then $x' = x^{(i \leftarrow b)}$ for some b with $q_i(b|x) > 0$. Let $a = x_i$ again. If $q_i(a|x') = 0$, our definition sets $\alpha_i(x \rightarrow x') = 0$, hence $K_i(x, x') = 0$ and detailed balance holds trivially because $K_i(x', x) = 0$ as well since the reverse move is never proposed.

Otherwise $q_i(a|x') > 0$. By the acceptance rule,

$$\begin{aligned} \pi_t(x) K_i(x, x') &= \pi_t(x) q_i(b|x) \frac{r_i(x, x')}{1 + r_i(x, x')}, \\ \pi_t(x') K_i(x', x) &= \pi_t(x') q_i(a|x') \frac{1}{1 + r_i(x, x')}, \end{aligned}$$

since the reverse acceptance uses $r_i(x', x) = 1/r_i(x, x')$ under Barker. Taking the ratio,

$$\frac{\pi_t(x) K_i(x, x')}{\pi_t(x') K_i(x', x)} = \frac{\pi_t(x) q_i(b|x) \frac{r_i}{1+r_i}}{\pi_t(x') q_i(a|x') \frac{1}{1+r_i}} = \frac{\pi_t(x) q_i(b|x) r_i(x, x')}{\pi_t(x') q_i(a|x')} = 1,$$

by the very definition

$$r_i(x, x') = \frac{\pi_t(x') q_i(a|x')}{\pi_t(x) q_i(b|x)}.$$

Hence $\pi_t(x) K_i(x, x') = \pi_t(x') K_i(x', x)$ for all off-diagonal pairs with positive mass. Including the diagonal case, which follows by summing both sides over x' , detailed balance condition holds. Therefore K_i is π_t -reversible and, in particular, π_t -invariant.

F EXPERIMENTS

We provide extra experiment information in this section.

F.1 EXPERIMENT DESCRIPTION

F.1.1 CHESSBOARD

The chessboard distribution has several key characteristics that make it an excellent test benchmark: (1) Discrete: The distribution is defined on a finite set of points on a two-dimensional grid; (2) Sparse: Nearly half of the grid points have a probability of exactly zero. A successful sampler must learn to restrict its generated samples strictly to the points that have non-zero probability mass (i.e., the support of the distribution); (3) Multi-modal: The probability mass is distributed across multiple, disconnected modes rather than being concentrated in a single region. The sampler is required to capture all of these distinct modes; (4) Structured: It exhibits a distinct, non-random geometric structure. This challenges the sampler’s ability to reproduce the correct global pattern, rather than merely matching general statistical moments.

We conduct the experiment on a 8×8 chessboard distribution with varying Picard iteration depths K_p . We fix $N = 40$, $M = 50$ and corrector step $K_c = 5$ with totally 8196 samples. Runtimes and KL Divergence are averaged over 20 runs.

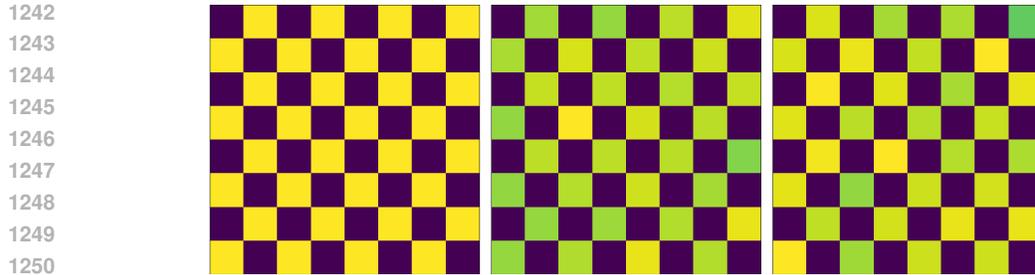


Figure 2: Visualization of the chessboard experiments. (Left) The target distribution. (Middle) The Picard sampling result. (Right) The sequential sampling result.

F.1.2 CIRCLE

The ring distribution on a 2D discrete grid concentrates its entire probability mass on grid points located within an annulus defined by an inner radius r_{in} and an outer radius r_{out} . A key characteristic is its Non-Convexity; the high-probability region encloses a central "hole" of zero probability, which leads to the distribution's support non-convex. Another defining feature is its Connectivity. Unlike the disjoint, multi-modal structure of the checkerboard distribution, the support of the ring distribution forms a single connected component, meaning any point on the ring can traverse to any other point through a series of steps to adjacent locations.

We conduct the experiment on a circle distribution at 32×32 2D grid. with varying Picard iteration depths K_p . We fix $N = 60$, $M = 50$ and corrector step $K_c = 8$ with totally 4096 samples. Runtime and KL Divergence are averaged over 20 runs.

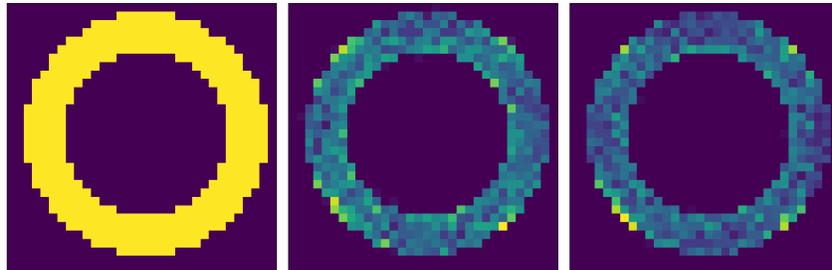


Figure 3: Visualization of the circle experiments. (Left) The target distribution. (Middle) The Picard sampling result. (Right) The sequential sampling result.

F.1.3 HYPERCUBE

The d -dimensional hypercube distribution is a uniform distribution over 2^d vertices of a hypercube, where each vertex corresponds to a unique binary vector in $\{0, 1\}^d$. With its well-defined graph structure, exponential state space, and large diameter, this distribution provides a rigorous benchmark for evaluating the high-dimensional scalability and mixing efficiency of sampling algorithms.

We conduct the experiment on a 6-dimensional hypercube distribution. We fix $N = 100$, $M = 50$ and corrector step $K_c = 10$ with totally 2048 samples. Runtime and KL Divergence are averaged over 20 runs.

F.1.4 EMBEDDED HYPERCUBE

We design an experiment with a controlled intrinsic dimension by embedding a k -dimensional subcube within a d -dimensional ambient space $k < d$. While the ambient space is the full hypercube, the probability mass is confined to a uniform distribution on the subcube. This is achieved by fixing the final $d - k$ bits to 0. The objective is to test the sampler's ability to identify the low-

dimensional manifold containing the data support. A proficient sampler should learn to set the last $d - k$ dimensions to 0 and generate uniform bit combinations for the first k dimensions.

We conduct the experiment on a 6-dimensional subcube distribution embedded in a 12-dimensional hypercube. We fix $N = 50, M = 100$ and corrector step $K_c = 8$ with totally 2048 samples. Runtime and KL Divergence are averaged over 20 runs.

F.2 COMPUTATIONAL COST ANALYSIS

To demonstrate the scalability of our algorithm with respect to data dimension d , we conducted additional experiments on sampling efficiency using d -dimensional independent Bernoulli distributions ($N = 40, M = 10, K_p = 3$). We observed that as d increases from 32 to 2048, the accuracy remains stable.

Table 8: Performance of sequential and parallel tau-leaping with different data dimension d .

Dim	Sequential (s)	Parallel (s)	Sequential Error	Parallel Error	GPU Memory (MB)
8	0.4341	0.0808	0.0042	0.0026	326
16	0.4344	0.0805	0.0037	0.0045	341
32	0.4242	0.0813	0.0049	0.0047	359
64	0.4155	0.0837	0.0035	0.0031	395
128	0.4153	0.0882	0.0038	0.0025	469
256	0.4203	0.1235	0.0044	0.0053	614
512	0.4097	0.2202	0.0040	0.0034	906
1024	0.5463	0.4144	0.0038	0.0033	1487
2048	1.1652	0.8041	0.0046	0.0042	2647

Although the runtime of our parallel algorithm begins to rise noticeably after $d = 128$, it maintains a performance advantage over the serial algorithm. As d increases further, the runtime of the serial algorithm spikes significantly after $d = 1024$, with a rate of growth that eventually surpasses that of the parallel approach. This empirically validates that, despite the additional logarithmic increase in complexity, the reduction in the polynomial term ensures the parallel algorithm remains advantageous in large d cases.

To demonstrate the scalability with respect to the number of small steps M , we tested on 256-dimensional independent Bernoulli distributions sampling with M from 10 to 40. ($N = 10, K_p = 3$)

Table 9: Performance of sequential and parallel tau-leaping with different data dimension M .

M	Parallel (s)	Parallel Error	GPU Memory (MB)
10	0.1235	0.0053	614
20	0.2668	0.0048	894
30	0.3610	0.0029	1176
40	0.4572	0.0031	1481