Exploring Non-Convex Discrete Energy Landscapes: A Langevin-Like Sampler with Replica Exchange

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

Paper under double-blind review

Abstract

Gradient-based Discrete Samplers (GDSs) are effective for sampling discrete energy landscapes. However, they often stagnate in complex, non-convex settings. To improve exploration, we introduce the Discrete Replica EXchangE Langevin (DREXEL) sampler and its variant with Adjusted Metropolis (DREAM). These samplers use two GDSs (Zhang et al., 2022b) at different temperatures and step sizes: one focuses on local exploitation, while the other explores broader energy landscapes. When energy differences are significant, sample swaps occur, which are determined by a mechanism tailored for discrete sampling to ensure detailed balance. Theoretically, we prove both DREXEL and DREAM converge asymptotically to the target energy and exhibit faster mixing than a single GDS. Experiments further confirm their efficiency in exploring non-convex discrete energy landscapes.

1 INTRODUCTION

Sampling from high-dimensional discrete distributions has been an important task for decades across 025 applications in texts (Mikolov et al., 2013; Devlin et al., 2019), images (Krizhevsky et al., 2012; 026 Ronneberger et al., 2015), signal processing (Mallat, 1989; Donoho, 2006), genome sequences 027 (Metzker, 2010; Macosko et al., 2015), etc. However, the exponential growth in the number of 028 configurations makes sampling from $\pi(\theta) \propto \exp[U(\theta)]$ computationally prohibitive. The computa-029 tional burden comes from evaluating the exact probabilities and normalizing constants, which makes exact sampling impossible in practice. Algorithms such as rejection sampling (Neumann, 1951), 031 Swendsen-Wang (Swendsen & Wang, 1987), and Hamze-Freitas (Hamze & de Freitas, 2004) leverage special structures within the problem to make global updates. In more general settings, these methods 033 may suffer from slow exploration, local dependencies, and poor convergence.

034 To make high-dimensional discrete sampling more efficient, Locally Balanced Proposals (LBPs) (Zanella, 2020; Sun et al., 2021) improved acceptance rates by adjusting 037 proposal distributions based on the likelihood ratio. Early 038 LBPs updated one coordinate at a time (Zanella, 2020; Grathwohl et al., 2021), and Grathwohl et al. (2021) developed gradient-based discrete sampler (GDS) to update 040 coordinately. Later, Zhang et al. (2022b) further extended 041 GDSs by updating all coordinates simultaneously, which 042 enhances efficiency and scalability for large-scale, high-043 dimensional computations on GPUs and TPUs. 044

Despite improvements in LBPs, how to balance the tradeoff between "global exploration" and "local exploitation" remains a challenge. High-dimensional discrete distributions are highly multi-modal, with deep and narrow wells caused by intrinsic discontinuities. Gradient-based LBPs, although effective, tend to get trapped in local modes due to their reliance on local gradients and small noise, which is insufficient for escaping these traps.



Figure 1: DREXEL & DREAM sample trajectory in discrete domains. Blue denotes a low-temperature sampler, and red hightemperature sampler. They exchange samples following a swap mechanism.

To bridge this gap, we propose two samplers: Discrete Replica EXchangE Langevin (DREXEL) and Discrete Replica Exchange with Adjusted Metropolis (DREAM). These samplers combine

GDS with the replica exchange Markov Chain Monte Carlo (reMCMC) (Chen et al., 2019) for
efficient exploration of non-convex discrete spaces. As illustrated in Figure 1, the samplers employ
two GDSs at different temperatures and step sizes: the low-temperature sampler focuses on local
exploitation, while the high-temperature sampler escapes local traps for broader exploration. Sample
swaps occur when energy differences are significant, governed by a mechanism tailored for discrete
sampling to ensure detailed balance. The combination of replica exchange and GDS makes them
particularly effective for sampling from complex discrete structures in modern applications. The
primary contributions in this work are summarized as follows:

- A novel integration of GDS with replica exchange to improve non-convex exploration;
- A swap mechanism tailored for detailed balance and sample efficiency in discrete sampling;
- Theoretical analysis of improved mixing rates over naïve discrete Langevin-like samplers;
- Superior performance in synthetic tasks, Ising models, restricted Boltzmann machines, and energy-based deep learning models for navigating non-convex discrete energy landscapes.

2 Related Work

062

063

064 065

066

067

068 069 070

071

098 099

100 101

072 Gradient-based Discrete Sampling becomes popular for complex discrete sampling tasks and its 073 original idea comes from LBPs. The concept of LBPs, as introduced by Zanella (2020), utilized local 074 information in the form of density ratios to improve the sample efficiency. Grathwohl et al. (2021) expanded LBPs by the use of first-order Taylor approximation, which further ensures computational 075 feasibility. To improve sampling in high-dimensional discrete spaces, LBPs were extended to cover 076 larger neighborhoods by performing a sequence of small moves (Sun et al., 2021). Zhang et al. 077 (2022b) further developed GDSs, which adapt the continuous Langevin MCMC to discrete spaces and allow parallel updates of all coordinates based on gradient information. Subsequently, GDSs 079 were improved through the introduction of an adaptive mechanism by which the step size can be automatically adjusted (Sun et al., 2023a). Most recently, Pynadath et al. (2024) introduced an 081 automatic cyclical scheduling approach in step sizes to better handle multi-modal distributions by 082 alternating between exploration and exploitation phases. While these methods have shown promise, 083 sampling from highly non-convex discrete distributions remains challenging, particularly when 084 dealing with strongly correlated variables or energy-based deep learning models. 085

Replica Exchange MCMC is a powerful method that enhances exploration in complex, multi-modal 086 distributions, and a variety of related algorithms build on this. For instance, unadjusted Langevin 087 MCMC (Durmus & Moulines, 2017) leverages gradient information to guide proposals but lacks 880 the exchange mechanism. Importance sampling (Wang & Landau, 2001) adjusts for the discrepancy 089 between target and proposal distributions, which offers flexibility in sampling but without temperaturebased exchanges. Simulated tempering (Lee et al., 2018) further refined the temperature-scaling 091 strategy by dynamically adjusting the temperature of a single chain. Recently, Zhang et al. (2020) proposed a cyclical step-size scheduler to maintain a balance between exploration and exploitation. 092 To enhance the exploration, reMCMC runs multiple chains at different temperatures and allows for 093 chain swaps between them. Dong & Tong (2022) analyzed its mixing by quantifying the spectral 094 gap, and Deng et al. (2020; 2022) validated its efficiency in large-scale deep learning tasks. Despite 095 its success in continuous sampling and high-dimensional settings, to the best of our knowledge, its 096 potential has not been studied in sampling from discrete distributions. 097

3 Preliminaries

The Target Distribution $\pi := \Theta \to [0, 1]$ denotes the probabilistic model we are sampling from:

$$\pi(\boldsymbol{\theta}) = \frac{1}{Z} \exp\left[\frac{U(\boldsymbol{\theta})}{\tau}\right], \quad \forall \boldsymbol{\theta} \in \Theta, \ \Theta \subset \mathbb{R}^{\mathbf{d}}.$$
 (1)

Here θ is **d**-dimensional variable, $\tau = 1.0$ denotes the respective temperature, Θ is a finite domain, *U* represents the energy function, and *Z* normalizes the distribution. Following the traditional settings in discrete sampling, we assume:

108 1. The sampling domain is coordinate-wisely factorized where $\Theta = \prod_{d=1}^{\mathbf{d}} \Theta_d$, and we primarily consider the binary cases $\Theta = \{0, 1\}^{\mathbf{d}}$ or categorical $\{0, 1, \dots, N-1\}^{\mathbf{d}}$;

2. The energy function is differentiable across \mathbb{R}^d .

118

127

147 148

159 160 161

112 The primary goal is to design an efficient sampler to approximate π within a finite sample size. 113 The empirical distribution derived from these samples converges to π , with the approximation error 114 bounded by a constant $\epsilon > 0$ under specific metrics.

Replica Exchange MCMC is a popular sampling method for non-convex exploration in continuous spaces. It updates according to the following dynamics:

$$\boldsymbol{\theta}_{i+1}^{(k)} = \boldsymbol{\theta}_i^{(k)} + \frac{\alpha_k}{2} \nabla U(\boldsymbol{\theta}_i^{(k)}) + \sqrt{\alpha_k \tau_k} \xi_k, \quad k = 1, 2; \quad i = 1, 2, \cdots, I.$$
(2)

Here $\alpha_1, \alpha_2 \in \mathbb{R}^+$ represent step sizes, $\tau_1, \tau_2 \in \mathbb{R}^+$ are temperatures, and ξ_1, ξ_2 are independent Gaussian noises drawn from $\mathcal{N}(\mathbf{0}, \mathbf{I}_{d \times d})$. The typical setup assumes $\alpha_1 < \alpha_2$ and $\tau_1 < \tau_2$, with the first chain in (2) labeled as the low-temperature chain and the second as the high-temperature chain. The gradient $\nabla U(\cdot)$ directs the algorithm toward high-probability regions. To further improve the mixing rate over Langevin MCMC, reMCMC facilitates interaction via a chain swap mechanism. Specifically, the probability to swap the *i*-th samples between $\theta_i^{(1)}$ and $\theta_i^{(2)}$ is determined by $\rho \min \{1, S(\theta_i^{(1)}, \theta_i^{(2)})\}$. The swap intensity is regulated by $\rho > 0$, and the swap function $S := \Theta \times \Theta \rightarrow \mathbb{R}^+$ is given as follows:

$$S(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) = e^{\left(\frac{1}{\tau_2} - \frac{1}{\tau_1}\right) \left[U(\boldsymbol{\theta}^{(1)}) - U(\boldsymbol{\theta}^{(2)}) \right]}.$$
(3)

Intuitively, the swap probability in reMCMC depends on the energy estimated at $\theta^{(1)}$ and $\theta^{(2)}$. When the low-temperature chain gets stuck in a local minimum and the high-temperature chain escapes to find modes with significantly lower energy, the chain will swap their samples with high probability. This enables the low-temperature chain to better characterize the newly discovered modes, while the high-temperature chain continues to search across the energy landscape. As mentioned in Chen et al. (2019), reMCMC behaves as a reversible Markov jump process due to its swap mechanism, which converges to a similar invariant distribution in (1) while the parameters can explore over \mathbb{R}^d .

135 While reMCMC is a powerful tool for non-convex exploration, its update may fail to preserve the 136 target distribution due to discretization errors introduced by the finite step size (Welling & Teh, 2011). 137 According to Roberts & Tweedie (1996), selecting an inappropriate step size can lead to a transient 138 Markov chain without a stationary distribution. To mitigate such bias, two main approaches are 139 commonly used: decaying step sizes (Vollmer et al., 2016; Teh et al., 2016) and Metropolis-Hastings (MH) corrections (Dwivedi et al., 2019; Chewi et al., 2021). While implementing decaying step sizes 140 is straightforward and does not require additional computational burden, the second approach is more 141 favorable due to its specific advantages in discrete sampling, which will be elaborated on later. 142

143 **Metropolis-Hastings Correction** is considered to correct discretization errors and ensure conver-144 gence to the target distribution. Specifically, at each iteration, a new candidate $\omega \leftarrow \theta_{i+1}^{-1}$ is first 145 generated with (2). To ensure that the resulting samples come from the target distribution, the MH 146 step determines whether to accept or reject the candidate with $\mathcal{A} := \Theta \times \Theta \rightarrow [0, 1]$:

$$\mathcal{A}(\boldsymbol{\omega}, \boldsymbol{\theta}_i) = \min\left\{1, \frac{\pi(\boldsymbol{\omega})q\left(\boldsymbol{\theta}_i \mid \boldsymbol{\omega}\right)}{\pi(\boldsymbol{\theta}_i)q\left(\boldsymbol{\omega} \mid \boldsymbol{\theta}_i\right)}\right\},\tag{4}$$

where $q := \Theta \times \Theta \rightarrow [0, 1]$ is the transition probability mapping from the current sample θ_i to the next sample θ_{i+1} . With probability $\mathcal{A}(\omega, \theta_i)$, the candidate ω is accepted in the current step; otherwise, it retains the current θ_i . This adjustment preserves the correct stationary distribution. Furthermore, because Langevin MCMC allows each sample to access any point in \mathbb{R}^d , it further ensures the Markov chain is both irreducible and ergodic (Diaconis & Freedman, 1997; Meyn & Tweedie, 2012).

Discrete Langevin Sampler (DLS) is a gradient-based approach for sampling from high-dimensional discrete distributions. Inspired by Langevin MCMC, DLS updates all coordinates in parallel from a single gradient computation to function effectively in discrete settings. Specifically, for a target distribution $\pi \propto \exp[U(\cdot)]$, DLS generates a new sample θ_{i+1} inspired by the Taylor expansion:

$$q\left(\boldsymbol{\theta}_{i+1} \mid \boldsymbol{\theta}_{i}\right) = \frac{\exp\left(-\frac{1}{2\alpha} \left\|\boldsymbol{\theta}_{i+1} - \boldsymbol{\theta}_{i} - \frac{\alpha}{2\tau} \nabla U(\boldsymbol{\theta}_{i})\right\|_{2}^{2}\right)}{Z_{\Theta}(\boldsymbol{\theta}_{i})}, \quad \boldsymbol{\theta}_{i}, \boldsymbol{\theta}_{i+1} \in \Theta.$$
(5)

¹For clarity and conciseness, we omit the chain index when there is no need to specify it.

162 Here $\nabla U(\theta)$ is the gradient of the energy function evaluated at θ , and $Z_{\Theta}(\theta)$ normalizes the distribution: 163

164

166

167

$$Z_{\Theta}(\boldsymbol{\theta}_{i}) = \sum_{\boldsymbol{\theta}_{i+1} \in \Theta} \exp\left(-\frac{1}{2\alpha} \left\|\boldsymbol{\theta}_{i+1} - \boldsymbol{\theta}_{i} - \frac{\alpha}{2\tau} \nabla U(\boldsymbol{\theta}_{i})\right\|_{2}^{2}\right).$$
(6)

This proposal distribution allows DLS to make larger, parallel updates while maintaining computational efficiency. As the dimension \mathbf{d} in parameter space grows, the cost of computing (6) becomes 168 prohibitively expensive. A key insight is that the update rule can be factorized by coordinate:

170 171

172

$$\theta_{i+1,d} \sim \text{Categorical}\left[\text{Softmax}\left(\frac{1}{2\tau}\nabla U(\boldsymbol{\theta}_i)_d(\theta_{i+1,d} - \theta_{i,d}) - \frac{(\theta_{i+1,d} - \theta_{i,d})^2}{2\alpha}\right)\right],\tag{7}$$

where d = 1, 2, ..., d is the dimension index, and $\theta_{i,d}$ represents the *i*-th sample in the *d*-th dimension. 173 This algebraic expansion, following the binomial theorem, works because $(\nabla U(\theta_i)_d)^2$ is independent 174 of $\theta_{i+1,d}$. It makes DLS scalable and computationally efficient for complex distributions (Zhang et al., 175 2022b). Furthermore, the first term in (7) biases the proposal towards low-energy regions, where the 176 gradient points towards increasing probability; the second term acts as a regularizing factor, which 177 penalizes large jumps unless they are strongly favored by the gradient. 178

179 DLS can operate with or without MH corrections. Without corrections, it is simplified to the Discrete Unadjusted Langevin Algorithm (DULA), which is computationally efficient but may introduce bias. 180 With corrections in (4), it becomes the Discrete Metropolis-adjusted Langevin algorithm (DMALA), 181 which corrects bias at an increased computational cost. Both DULA and DMALA employ non-local 182 proposals specifically for the heat kernel to enable more efficient sampling (Sun et al., 2023a). 183

184 185

187

188 189 190

191

196 197

199

DISCRETE LANGEVIN SAMPLER WITH REPLICA EXCHANGE 4

The proposed DLS variants are present here, which incorporate replica exchange and a customized sampler swap mechanism to ensure detailed balance. The complete algorithm is provided at the end.

DISCRETE SAMPLERS WITH DIFFERENT TEMPERATURES 4.1

A key challenge with the naïve DLS is the tendency to become trapped in local modes, particularly 192 in non-convex landscapes. To mitigate this, we introduce DREXEL, which incorporates replica 193 exchange to enable efficient exploration across different local modes. Specifically, we employ two 194 samplers separately with distinct step sizes and temperatures to approximate the target distribution: 195

Categorical
$$\left[\text{Softmax} \left(\frac{1}{2\tau_k} \nabla U(\theta_i^{(k)})_d \left(\theta_{i+1,d}^{(k)} - \theta_{i,d}^{(k)} \right) - \frac{\left(\theta_{i+1,d}^{(k)} - \theta_{i,d}^{(k)} \right)^2}{2\alpha_k} \right) \right], \quad k = 1, 2.$$
(8)

Here $\tau_1 < \tau_2$ and $\alpha_1 < \alpha_2$, with k = 1 being the low-temperature and k = 2 the high-temperature 200 sampler. Intuitively, larger step sizes and higher temperatures encourage more exploratory moves, 201 which allows the sampler to escape local modes through non-local jumps and explore different regions 202 of the energy landscapes. This, on the downside, raises the rejection rate, as large jumps often land in 203 low-probability regions, and introduce additional bias when approximating the target distribution. 204

To mitigate the bias, we further propose DREAM, which incorporates MH steps post-generation of 205 new samples. Once the new samples are produced through (8), the acceptance rates $\mathcal{A}(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}_{i}^{(1)})$ and 206 $\mathcal{A}(\boldsymbol{\theta}_{i+1}^{(2)}, \boldsymbol{\theta}_{i}^{(2)})$ are estimated with (4). The new samples are accepted with probability \mathcal{A} or rejected 207 with $1 - \mathcal{A}$. The acceptance rates of two samplers are independent of one another. While the 208 high-temperature sampler typically exhibits a lower acceptance rate than the low-temperature one, 209 the rejection mechanism ensures that both samplers in DREAM converge to the target asymptotically. 210

211 It should be noted that while decaying step sizes are commonly advantageous in Langevin MCMC for 212 handling big data (Teh et al., 2016), they present potential challenges in discrete sampling. In discrete 213 spaces, small steps do not equate to gradual movements as they do in continuous spaces. Instead, they tend to repeatedly propose nearly identical samples, which causes the sampler to become trapped in 214 local regions. This problem becomes severe when dealing with non-convex energy landscapes, where 215 a decaying step size worsens the issue of local traps. For this reason, the MH step is often favored as

a solution in discrete sampling. With the MH step and fixed step sizes, the sampler can make large
 jumps to facilitate global exploration. This feature is essential for navigating highly structured state
 spaces, where the sampler needs flexibility to move between distant states.

In practice, high-temperature samplers may have difficulty exploiting certain regions due to abrupt exploration, which requires excessive time to fully characterize local modes and achieve mixing.

4.2 SAMPLE SWAPS BETWEEN DISCRETE SAMPLERS

A typical solution is to implement a swap function that enables sample exchanges between samplers at different temperatures. This helps cross energy barriers by combining the exploration of hightemperature samplers with the exploitation of low-temperature ones, which improves mixing rates.

The naïve swap function (3) of reMCMC relies on energy calculations at the current samples and corresponding temperatures. However, it is not practical to handle large-scale data in mini-batch settings. Intuitively, while $\tilde{U}(\theta_{i+1}^{(1)})$ and $\tilde{U}(\theta_{i+1}^{(2)})$ are both unbiased in mini-batches, a non-linear transformation of these estimators fail to provide an unbiased estimator for $S(\theta_{i+1}^{(1)}, \theta_{i+1}^{(2)})$ (Deng et al., 2020). Under normality assumption for the energy estimate, we consider a bias correction term:

$$\tilde{S}(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}_{i+1}^{(2)}) = e^{\left(\frac{1}{\tau_2} - \frac{1}{\tau_1}\right) \left[U(\boldsymbol{\theta}_{i+1}^{(1)}) - U(\boldsymbol{\theta}_{i+1}^{(2)}) + \left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \sigma^2 \right]},\tag{9}$$

where σ^2 compensates for noise in the stochastic gradient and removes swap bias. This adjustment ensures that the swap function behaves as a Martingale and matches the expected value obtained from exact gradients. Although this correction is not strictly necessary in discrete sampling, we retain this design in practice and examine the potential need for bias correction in the experiments. The bias-corrected versions of DREXEL and DREAM are referred to as bDREXEL and bDREAM.

When reMCMC is applied to discrete spaces, a notable challenge arises: the decaying step sizes commonly employed in continuous settings are not applicable. To ensure asymptotic convergence to the target distribution with fixed step sizes, we must maintain detailed balance not only between the low-temperature and high-temperature samplers but also between the current and next output samples. The swap designs in (3) and (9), however, overlook energy and temperature differences. This potentially violates detailed balance and slows down mixing in discrete sampling tasks. To mitigate the imbalance, we propose a swap function tailored for discrete sampling:

$$\tilde{S}(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}_{i+1}^{(2)} \mid \boldsymbol{\theta}_{i}^{(1)}, \boldsymbol{\theta}_{i}^{(2)}) = e^{\left(\frac{1}{\tau_{2}} - \frac{1}{\tau_{1}}\right) \left[U(\boldsymbol{\theta}_{i+1}^{(1)}) + U(\boldsymbol{\theta}_{i}^{(1)}) - U(\boldsymbol{\theta}_{i+1}^{(2)}) - U(\boldsymbol{\theta}_{i}^{(2)}) \right]}.$$
(10)

This swap function incorporates energy estimates at the last samples, which respects the energy landscape and preserves detailed balance. Importantly, since the previous samples are treated as constants during the swap, the detailed balance between replicas remains unaffected. We will demonstrate how this correction guarantees asymptotic convergence to the target distribution in the next section.

4.3 The Proposed Algorithms

As outlined in Algorithm 1, we present DREXEL and DREAM for discrete sampling. The approaches
employ two DLSs with distinct temperatures and step sizes, which allows for sample swaps between
them. At each iteration, the current samples are updated, followed by MH steps in DREAM. The
swap mechanism exchanges samples when the high-temperature sampler locates a lower-energy
mode. After *I* iterations, the low-temperature sampler outputs samples to characterize the energy
landscape. This approach, discussed further in 5.1, improves the mixing rate over DLSs by balancing
exploration and exploitation.

70	Algo	orithm 1 DREXEL and DREAM.
71	Inpu	ut Step Sizes α_1, α_2
72	Inpi	ut Temperatures τ_1, τ_2 .
73	Inpu	at Swap Intensity $\rho > 0$.
74	Inpu	ut Initial Samples $\theta_0^{(k)} \in \Theta$, $k = 1, 2$.
75	1: I	For $i = 1, 2, \cdots, I$ do
76	2:	Sampling Steps:
77	3:	For $k = 0, 1, 2$ do:
78	4:	For $d = 1, 2, \dots, d$ do:
79	5:	Construct $q_d^{(s)}(\theta^{(s)} \theta_i^{(s)})$ following (8)
80	6:	Sample $\omega_d^{(n)} \sim q_d^{(n)}(\cdot \mid \boldsymbol{\theta}_i^{(n)})$
81	/: 8.	End For Fnd For
82	0.	MU Stone (for DDE (M))
83	9. 10·	For $k = 1.2$ do:
84	11:	Compute $\mathcal{A}(\boldsymbol{\theta}_{\cdot}^{(k)}, \boldsymbol{\theta}_{\cdot}^{(k)})$ following (4)
85	12:	Generate a number $u \sim U[0, 1]$
86	13:	Set $\boldsymbol{\theta}_{i+1}^{(k)} \leftarrow \boldsymbol{\omega}^{(k)}$ if $u \leq \mathcal{A}$ else $\boldsymbol{\theta}_{i+1}^{(k)} \leftarrow \boldsymbol{\theta}_{i}^{(k)}$
87	14:	End For
88	15:	Swapping Steps:
20	16:	Generate a number $u \sim U[0, 1]$.
09	17:	Compute $S(\theta_{i+1}^{(r)}, \theta_{i+1}^{(r)})$ following (10)
90	18:	Swap $\boldsymbol{\theta}_{i+1}^{(1)}$ and $\boldsymbol{\theta}_{i+1}^{(2)}$ if $u \leq \rho \min\left\{1, S\right\}$
91	19: I	End For
92	Out	put Samples $\{\theta_i^{(1)}\}_{i=1}^l$.
93		

5 Theoretical Analysis

In the previous section, we introduced DREXEL and DREAM, which use factorization to allow parallel updates and employ swap mechanisms to improve non-convex exploration. While these features are beneficial, the overall performance heavily relies on their convergence properties and theoretical guarantees. In this section, we provide asymptotic convergence guarantees for DREXEL (i.e. the version without the MH correction).

318

322

323

295 296

297 298

299

300

301

5.1 Asymptotic Convergence on Log-Quadratic Distributions

Our focus is first on the asymptotic behaviors of DREXEL. The analysis aims to show that as step sizes approach zero, the algorithm exhibits zero asymptotic bias, which ensures accurate sampling from the target distribution. Specifically, we focus on log-quadratic energy $\pi(\theta) \propto \exp(\theta^{T} J \theta + b^{T} \theta)$, where $J \in \mathbb{R}^{d \times d}$ is a symmetric matrix, and $b \in \mathbb{R}^{d}$ is a vector. If J is asymmetric, we apply spectral decomposition to obtain a symmetric matrix, which enables an analytically tractable solution.

Zhang et al. (2022b) showed that DLS with temperature 1 is reversible for log-quadratic energy distributions when the step size is sufficiently small. However, this result does not directly extend to the proposed algorithm, as the swap mechanism (3) introduces potential imbalances. This imbalance is due to discontinuous transitions between neighboring states in discrete space, which makes the swap acceptance rule insufficient to maintain a detailed balance. This further introduces bias during the sampling process and leads to inaccurate modeling. To address this, we carefully control the swap probability in (10) to regulate transitions between high- and low-temperature samplers.

Theorem 1. Let α_1 and α_2 be the step sizes for the low- and high-temperature samplers, and let $q(\cdot|\theta)$ be the Markov chain transition kernel. Suppose the target $\pi(\theta)$ is log-quadratic, then: **The Markov chain induced by DREXEL is reversible with respect to an intermediate distribution**

- The Markov chain induced by DREXEL is reversible with respect to an intermediate distribution $\tilde{\pi}$, i.e., for all $\theta, \theta' \in \Theta$, $\tilde{\pi}(\theta)q(\theta'|\theta) = \tilde{\pi}(\theta')q(\theta|\theta')$.
- As $\alpha_1, \alpha_2 \rightarrow 0$, the stationary distribution π' converges weakly to the target distribution π .

This analysis focuses on the state transition of the low-temperature sampler, as the high-temperature sampler only facilitates exploration and does not produce final samples. Intuitively, with probability $\rho \tilde{S}$, the next low-temperature sample is drawn from the high-temperature sampler, and with probability $\rho (1 - \tilde{S})$, it selects from the low-temperature sampler. The transition simplifies to DLS without swaps and directly maintains the detailed balance, but the swap probability becomes essential for preserving this balance once swaps are considered in discrete sampling. Our designed swap function ensures that the overall transition dynamics remain balanced, as demonstrated in Appendix D.1.

6 Experiments

To illustrate the effectiveness of our approach, we evaluate the proposed samplers across distinct discrete sampling and generative tasks. Our approach is compared against baselines including DLS (DULA and DMALA from Zhang et al. (2022b)), Any-scale Balanced sampling (AB) (Sun et al., 2023a), and the Automatic Cyclical Sampler (ACS) (Pynadath et al., 2024). More details such as experimental setups, hyper-parameters, and additional experimental results can refer to Appendix E.

6.1 SAMPLING FROM 2D SYNTHETIC PROBLEMS

We explore the challenges of sampling from 2D discrete multi-modal distributions defined over the domain $\Theta = \{1, 2, ..., N\}^d$, where N = 256 and $\mathbf{d} = 101 \times 101$. Each coordinate takes one of the discrete values. Figure 2 (top) highlights the challenges of approximating non-convex energy landscapes, where samplers often struggle to explore the landscapes effectively with limited samples.



Figure 2: **Top:** Visualization of the true energy landscapes. **Middle:** Empirical energy distributions from DMALA. **Bottom:** Empirical energy distributions from DREAM.

Metric	Sampler	wave	8gaussian 1	l 6gaussian	moon	2moon	twist	flower
	DMALA	2.419	1.337	7.690	2.397	4.848	3.767	2.765
$VI (10^{-2})$	AB	1.028	0.851	3.373	2.567	4.127	3.033	2.421
$KL(10^{-2})\downarrow$	ACS	0.930	0.521	3.145	2.059	4.207	2.154	2.479
	DREAM	0.914	0.519	3.017	1.652	4.252	2.145	2.277
	DMALA	2.085	2.084	1.977	2.095	2.019	2.049	2.100
	AB	2.036	2.057	1.910	1.911	2.070	2.059	1.983
MMD↓	ACS	2.014	2.028	1.984	1.996	2.068	2.047	1.966
	DREAM	1.969	2.007	1.922	1.908	2.014	1.976	1.913

Table 1: Experiment results with exploring 2D synthetic distributions, recorded with KL and MMD.

DREAM is compared with DMALA, AB, and ACS in these tasks. With automatic differentiation for
 gradient computation, we generate 640,000 samples to form the empirical distributions. Figure 2
 (bottom) provides a qualitative analysis, showing that DREAM can effectively capture the underlying
 complex distributions. In the wave, 8gaussians, and 16gaussians, DMALA (mid) captures only 50%

modes due to its tendency to get stuck in local minima. DREAM, by contrast, recovers all modes,
which reflects its robust exploration across different tasks. For quantitative evaluation, we report
Kullback–Leibler (KL) divergence and maximum mean discrepancy (MMD) as performance metrics
(Blessing et al., 2024). As shown in Table 1, DREAM consistently outperforms the baselines across
all distributions.

6.2 SAMPLING FROM ISING MODELS

The Ising model (Newman & Barkema, 1999) is a mathematical structure used to describe systems of interacting binary variables, which are commonly represented as spins in physical systems (MacKay, 2003). Each spin can assume binary states and interact with adjacent spins within a lattice. The interactions between these neighboring spins are governed by the energy function $U(\theta) = w\theta^{T} J\theta + b^{T}\theta$, where $\theta \in \{-1, 1\}^{d}$ is binary random variable, $J \in \{0, 1\}^{d \times d}$ is a binary adjacency matrix, $w \in \mathbb{R}^{+}$ denote the connectivity strength, and $b \in \{0, 1\}^{d}$ is the bias vector.



Figure 3: Ising model sampling results, evaluated by log RMSE. DREAM yields the best scores.

From Figure 3, for samplers without MH corrections, DREXEL shows a fast and consistent reduction in log RMSE, while DULA and ACS converge more slowly. Note that ACS exhibits periodic fluctuations due to its cyclical step size, where error initially decreases but then increases as the step size decays. For discrete sampling, it implies small step sizes do not effectively exploit local modes, and decaying step sizes may not be as effective as MH corrections in discrete settings. Among samplers with MH corrections, DREAM delivers the most efficient error reduction, which benefits from a strong exploration-exploitation balance. DMALA follows a more gradual path, while mh-ACS mirrors the periodic behavior of ACS due to its similar step-size schedule. These findings indicate that the proposed samplers generally offer better and more reliable mixing rates.

410 411 412

413

420

421

422

423

424

425

384

385 386

387

388

389

390

391

401 402

403

404

405

406

407

408

409

6.3 SAMPLING FROM RESTRICTED BOLTZMANN MACHINES

414 Restricted Boltzmann Machines (RBMs) are generative stochastic neural networks designed to model 415 complex distributions over discrete data (Fischer & Igel, 2012). RBMs typically consist of binary-416 valued hidden and visible units, where the visible units represent observed data and the hidden units 417 capture latent dependencies in the data. The energy function $U(\theta) = \log [1 + \exp (J\theta + c)] + b^{\top}\theta$, 418 where $\theta \in \{0, 1\}^d$ represents the binary state vector for the visible layer, $J \in \mathbb{R}^{m \times d}$ is the weight 419 matrix, $c \in \mathbb{R}^m$ and $b \in \mathbb{R}^d$ denote biases for hidden units and visible units correspondingly.



426 427



We trained RBMs with 500 hidden units on the MNIST dataset using contrastive divergence (Hinton, 2002). To benchmark the effectiveness of various samplers, we reported MMDs between outputs generated by each sampler and the one by the structured Block-Gibbs sampler specific to RBMs.

Figure 4 highlights DREAM as the most effective sampler, which consistently achieves the best log
MMD across all datasets and suggests superior convergence in sampling from RBMs. AB and ACS
perform nearly as well but with slightly higher MMD values. Overall, DREAM provides the most
robust sampling, followed by AB and ACS, while DMALA trails behind.

436 437

438

449

458

459

460

461

462

463

464

465

6.4 LEARNING ENERGY-BASED MODELS

439 Deep Energy-Based Models (EBMs) (Ngiam et al., 2011; Bond-Taylor et al., 2021) are a class 440 of probabilistic models where the energy function is parameterized by a ResNet (He et al., 2016). 441 Specifically, the probability of a data point x is given by $P_{\theta}(x) = \exp[E_{\theta}(x)]/Z_{\theta}$, where $E_{\theta}(x)$ is the 442 energy function parameterized by θ , and $Z_{\theta} = \mathbb{E}_{\theta \sim \Theta} \exp[E_{\theta}(x)]$ normalizes the distribution.

With DULA and DMALA as baselines, we evaluate DREXEL and DREAM² by learning Deep EBMs.
During training, the intractable likelihood gradient of the model is approximated through Persistent
Contrastive Divergence (Tieleman & Hinton, 2009), while a replay buffer (Du & Mordatch, 2019)
containing 1,000 past samples is implemented to improve both the efficiency and stability of the
training process. Each sampler runs for 40 steps per iteration. Upon completing training, Annealed
Importance Sampling (Neal, 2001) is conducted with DULA to estimate the test log-likelihoods.

Table 2: Test log-likelihoods of Deep EBMs evaluated on image datasets.

-	Dataset	DULA	DMALA	bDREXEL	bDREAM	DREXEL	DREAM
-	Static MNIST	-84.579	-85.145	-85.638	-84.823	-84.509	-83.929
	Dynamic MNIST	-86.625	-84.799	-86.907	-85.104	-83.984	-82.963
	Omniglot	-118.541	-111.820	-102.405	-100.042	-101.930	-98.454
	Caltech	-108.626	-107.820	-108.199	-107.899	-93.481	-92.003

We trained Deep EBMs for 20,000 iterations on binary images from Static MNIST, Dynamic MNIST, Omniglot, and Caltech Silhouettes datasets. The test log-likelihoods for trained models across different samplers are recorded in Table 2. Among the samplers, DREAM consistently achieved the highest log-likelihoods across all datasets, with notable improvements on Omniglot and Caltech Silhouettes. For MNIST datasets, DREXEL and DREAM also showed competitive performance, particularly on Static MNIST. In contrast, bDREXEL and bDREAM generally performed worse, with DREXEL and DREAM showing clear superiority across most datasets. These findings confirm that MH steps are essential for improving performance in discrete sampling tasks. Also, the proposed swap mechanism in (10) is effective at correcting imbalance and yielding better log-likelihood estimates across diverse image datasets.

466 467 468

469

7 CONCLUSION AND DISCUSSION

In this work, we addressed the challenge of balancing global exploration and local exploitation
 in non-convex discrete energy landscapes by proposing DREXEL and DREAM. These samplers
 integrate DLS with replica exchange to overcome the limitations of traditional samplers, which tend
 to get trapped in local modes due to reliance on local gradients and small disturbances.

We theoretically prove that the proposed samplers are reversible, which guarantees the accurate preservation of the target distribution. Moreover, these samplers achieve faster mixing than the naïve DLS. The empirical evidence suggests that the proposed samplers and swap mechanism significantly improve exploration and mixing in non-convex discrete spaces. Furthermore, while DREXEL maintains detailed balance throughout the process, MH corrections are critical for optimizing performance in certain tasks.

Our current work focuses on designing a single low-temperature and a single high-temperature sampler. Future research could extend this framework by introducing multiple parallel samplers to enhance exploration. We will also study theoretical guarantees to quantify the acceleration effect of the swap mechanism in the future.

²Unless stated otherwise, DREXEL and DREAM refer to samplers based on the swap design from (10). bDREXEL and bDREAM apply the swap function described in (9).

486 ETHICS STATEMENT 487

488

489

490

491 492

493 494

495

496

497 498

499 500

501

528

We adhere to the ICLR Code of Ethics and confirm that our experiments use only public datasets. While our results are primarily based on standard benchmarks, we recognize the potential for misuse and encourage responsible application of our methods on real-world data.

Reproducibility Statement

To facilitate the reproducibility of our proposed method, the implementation of the synthetic task related to this work is made available at the following anonymous link: https://anonymous. 4open.science/r/dream-F7E6.

References

- Sungjin Ahn, Anoop Korattikara, and Max Welling. Bayesian Posterior Sampling via Stochastic Gradient Fisher Scoring. In *International Conference on Machine Learning*, pp. 1771–1778, 2012.
- Patrick Billingsley. *Probability and Measure*. John Wiley & Sons, 2017.
- Denis Blessing, Xiaogang Jia, Johannes Esslinger, Francisco Vargas, and Gerhard Neumann. Beyond
 elbos: A large-scale evaluation of variational methods for sampling. In *International Conference on Machine Learning*, 2024.
- Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks. Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7327–7347, 2021.
- Andrew Campbell, Wenlong Chen, Vincent Stimper, Jose Miguel Hernandez-Lobato, and Yichuan
 Zhang. A Gradient Based Strategy for Hamiltonian Monte Carlo Hyperparameter Optimization. In International Conference on Machine Learning, pp. 1238–1248. PMLR, 2021.
- Changyou Chen, Nan Ding, and Lawrence Carin. On the Convergence of Stochastic Gradient MCMC
 Algorithms with High-Order Integrators. *Advances in Neural Information Processing Systems*, 28, 2015.
- Yi Chen, Jinglin Chen, Jing Dong, Jian Peng, and Zhaoran Wang. Accelerating Nonconvex Learning via Replica Exchange Langevin Diffusion. In *International Conference on Learning Representation*, 2019.
- Sinho Chewi, Chen Lu, Kwangjun Ahn, Xiang Cheng, Thibaut Le Gouic, and Philippe Rigollet.
 Optimal Dimension Dependence of the Metropolis-Adjusted Langevin Algorithm. In *Conference on Learning Theory*, pp. 1260–1300. PMLR, 2021.
- Hanjun Dai, Rishabh Singh, Bo Dai, Charles Sutton, and Dale Schuurmans. Learning Discrete
 Energy-Based Models via Auxiliary-Variable Local Exploration. Advances in Neural Information
 Processing Systems, 33:10443–10455, 2020.
- Wei Deng, Qi Feng, Liyao Gao, Faming Liang, and Guang Lin. Non-convex Learning via Replica
 Exchange Stochastic Gradient MCMC. In *International Conference on Machine Learning*, pp. 2474–2483. PMLR, 2020.
- Wei Deng, Qi Feng, Georgios Karagiannis, Guang Lin, and Faming Liang. Accelerating Convergence
 of Replica Exchange Stochastic Gradient MCMC via Variance Reduction. In *International Conference on Learning Representation*, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-Training of Deep
 Bidirectional Transformers for Language Understanding. In North American Chapter of the
 Association for Computational Linguistics, 2019.
- 539 Persi Diaconis and David Freedman. On Markov Chains with Continuous State Space. *Annals of Probability*, 1997.

540 541 542	Jing Dong and Xin T Tong. Spectral Gap of Replica Exchange Langevin Diffusion on Mixture Distributions. <i>Stochastic Processes and Their Applications</i> , 151:451–489, 2022.
542 543	David L Donoho. Compressed Sensing. <i>IEEE Transactions on information theory</i> , 52(4):1289–1306, 2006.
545	
546	Yilun Du and Igor Mordatch. Implicit Generation and Modeling with Energy Based Models. Advances in Neural Information Processing Systems, 32, 2019.
547	Alain Durmus and Éric Moulines Non asymptotic Convergence Analysis for the Unadjusted
548 549	Langevin Algorithm. <i>The Annals of Applied Probability</i> , pp. 1551–1587, 2017.
550 551	Raaz Dwivedi, Yuansi Chen, Martin J Wainwright, and Bin Yu. Log-Concave Sampling: Metropolis- Hastings Algorithms Are Fast. <i>Journal of Machine Learning Research</i> , 20(183):1–42, 2019.
553 554 555	Asja Fischer and Christian Igel. An Introduction to Restricted Boltzmann Machines. In <i>Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications</i> , pp. 14–36. Springer, 2012.
556 557	Gerald B Folland. <i>Real Analysis: Modern Techniques and Their Applications</i> , volume 40. John Wiley & Sons, 1999.
558 559 560 561	Will Grathwohl, Kevin Swersky, Milad Hashemi, David Duvenaud, and Chris Maddison. Oops I Took a Gradient: Scalable Sampling for Discrete Distributions. In <i>International Conference on Machine Learning</i> , pp. 3831–3841. PMLR, 2021.
562 563	Firas Hamze and Nando de Freitas. From Fields to Trees. In <i>Conference on Uncertainty in Artificial Intelligence</i> , pp. 243–250, 2004.
564 565 566	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 770–778, 2016.
567 568	Geoffrey E Hinton. Training Products of Experts by Minimizing Contrastive Divergence. <i>Neural Computation</i> , 14(8):1771–1800, 2002.
569 570 571	Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet Classification with Deep Convolutional Neural Networks. <i>Advances in Neural Information Processing Systems</i> , 25, 2012.
572 573 574	Holden Lee, Andrej Risteski, and Rong Ge. Beyond Log-Concavity: Provable Guarantees for Sampling Multi-modal Distributions Using Simulated Tempering Langevin Monte Carlo. <i>Advances</i> <i>in Neural Information Processing Systems</i> , 31, 2018.
575 576 577	Xuechen Li, Yi Wu, Lester Mackey, and Murat A Erdogdu. Stochastic Runge-Kutta Accelerates Langevin Monte Carlo and Beyond. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
579 580	Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete Diffusion Modeling by Estimating the Ratios of the Data Distribution. In <i>International Conference on Machine Learning</i> , 2024.
581 582	David JC MacKay. Information Theory, Inference and Learning Algorithms. Cambridge university press, 2003.
583 584 585 586 587	Evan Z Macosko, Anindita Basu, Rahul Satija, James Nemesh, Karthik Shekhar, Melissa Goldman, Itay Tirosh, Allison R Bialas, Nolan Kamitaki, Emily M Martersteck, et al. Highly Parallel Genome-Wide Expression Profiling of Individual Cells Using Nanoliter Droplets. <i>Cell</i> , 161(5): 1202–1214, 2015.
588 589	Stephane G Mallat. A Theory for Multiresolution Signal Decomposition: The Wavelet Representation. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 11(7):674–693, 1989.
590 591	Enzo Marinari and Giorgio Parisi. Simulated Tempering: A New Monte Carlo Scheme. <i>Europhysics Letters</i> , 19(6):451, 1992.
592 593	Michael L Metzker. Sequencing Technologies—the Next Generation. <i>Nature Reviews Genetics</i> , 11 (1):31–46, 2010.

594 595 596	Sean P Meyn and Richard L Tweedie. <i>Markov Chains and Stochastic Stability</i> . Springer Science & Business Media, 2012.
597 598	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Represen- tations in Vector Space. In <i>International Conference on Learning Representation</i> , 2013.
599 600	Radford M Neal. Annealed Importance Sampling. Statistics and Computing, 11:125–139, 2001.
601 602	Radford M. Neal. MCMC Using Hamiltonian Dynamics. In <i>Handbook of Markov Chain Monte Carlo</i> , volume 54, pp. 113–162. Chapman and Hall/CRC, 2012.
603 604 605	Von Neumann. Various Techniques Used in Connection with Random Digits. <i>Notes by GE Forsythe</i> , pp. 36–38, 1951.
606 607	Mark EJ Newman and Gerard T Barkema. <i>Monte Carlo Methods in Statistical Physics</i> . Clarendon Press, 1999.
608 609 610	Jiquan Ngiam, Zhenghao Chen, Pang W Koh, and Andrew Y Ng. Learning Deep Energy Models. In <i>International Conference on Machine Learning</i> , pp. 1105–1112, 2011.
611 612	Patrick Pynadath, Riddhiman Bhattacharya, Arun Hariharan, and Ruqi Zhang. Gradient-Based Discrete Sampling with Automatic Cyclical Scheduling. <i>arXiv preprint arXiv:2402.17699</i> , 2024.
613 614 615	Ali Rahimi and Benjamin Recht. Random Features for Large-Scale Kernel Machines. Advances in Neural Information Processing Systems, 20, 2007.
616 617	Gareth O Roberts and Richard L Tweedie. Exponential Convergence of Langevin Distributions and Their Discrete Approximations. <i>Bernoulli</i> , 2(4):341–363, 1996.
618 619 620 621	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In <i>International Conference on Medical Image Computing and Computer-</i> <i>Assisted Intervention</i> , pp. 234–241. Springer, 2015.
622 623 624	Umut Simsekli, Roland Badeau, Taylan Cemgil, and Gaël Richard. Stochastic Quasi-Newton Langevin Monte Carlo. In <i>International Conference on Machine Learning</i> , pp. 642–651. PMLR, 2016.
625 626 627	Haoran Sun, Hanjun Dai, Wei Xia, and Arun Ramamurthy. Path Auxiliary Proposal for MCMC in Discrete Space. In <i>International Conference on Learning Representation</i> , 2021.
628 629	Haoran Sun, Bo Dai, Charles Sutton, Dale Schuurmans, and Hanjun Dai. Any-Scale Balanced Samplers for Discrete Space. In <i>International Conference on Learning Representation</i> , 2023a.
630 631 632	Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-Based Continuous-Time Discrete Diffusion Models. In <i>International Conference on Learning Representation</i> , 2023b.
633 634	Robert H Swendsen and Jian-Sheng Wang. Nonuniversal Critical Dynamics in Monte Carlo Simula- tions. <i>Physical Review Letters</i> , 58(2):86, 1987.
635 636 637	Yee Whye Teh, Alexandre Thiéry, and Sebastian J Vollmer. Consistency and Fluctuations for Stochastic Gradient Langevin Dynamics. <i>Journal of Machine Learning Research</i> , 17(7), 2016.
638 639	Tijmen Tieleman and Geoffrey Hinton. Using Fast Weights to Improve Persistent Contrastive Divergence. In <i>International Conference on Machine Learning</i> , pp. 1033–1040, 2009.
640 641 642 643	Sebastian J Vollmer, Konstantinos C Zygalakis, and Yee Whye Teh. Exploration of the (non-) asymptotic bias and variance of stochastic gradient Langevin dynamics. <i>Journal of Machine Learning Research</i> , 17(159):1–48, 2016.
644 645	Fugao Wang and David P Landau. Efficient, Multiple-Range Random Walk Algorithm to Calculate the Density of States. <i>Physical Review Letters</i> , 86(10):2050, 2001.
646 647	Max Welling and Yee W Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. In <i>International Conference on Machine Learning</i> , pp. 681–688. Citeseer, 2011.

648 649 650	Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Hard Prompts Made Easy: Gradient-Based Discrete Optimization for Prompt Tuning and Discovery. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
651 652 653	Ulli Wolff. Collective Monte Carlo Updating for Spin Systems. <i>Physical Review Letters</i> , 62(4):361, 1989.
654 655	Giacomo Zanella. Informed Proposals for Local MCMC in Discrete Spaces. <i>Journal of the American Statistical Association</i> , 115(530):852–865, 2020.
656 657 658 659	Dinghuai Zhang, Nikolay Malkin, Zhen Liu, Alexandra Volokhova, Aaron Courville, and Yoshua Ben- gio. Generative Flow Networks for Discrete Probabilistic Modeling. In <i>International Conference</i> <i>on Machine Learning</i> , pp. 26412–26428. PMLR, 2022a.
660 661 662	Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning. In <i>International Conference on Learning Representation</i> , 2020.
663 664 665	Ruqi Zhang, Xingchao Liu, and Qiang Liu. A Langevin-Like Sampler for Discrete Distributions. In <i>International Conference on Machine Learning</i> , pp. 26375–26396. PMLR, 2022b.
666 667 668	
669 670	
671 672	
673 674 675	
676 677	
678 679	
680 681 682	
683 684	
685 686	
688 689	
690 691	
692 693	
695 696	
697 698	
699 700 701	

A FUTHER DISCUSSION ON RELATED WORK

704 Modeling Non-Convex Discrete Distributions is important in modern machine learning tasks. 705 Swendsen & Wang (1987) designed cluster updates in Monte Carlo simulations to efficiently navigate 706 complex energy landscapes by updating groups of variables simultaneously. Wolff (1989) further 707 improved sampling efficiency by flipping single large clusters of spins, thereby reducing autocorrelation times near critical points in non-convex distributions. Marinari & Parisi (1992) considered 708 temperature as a dynamic variable and allowed the system to overcome energy barriers and explore multiple modes of a non-convex distribution effectively. While these models might work for specific 710 discrete structures, they rely on random walk or Gibbs sampling. Their extensions to gradient-based 711 discrete sampling may significantly improve their efficiency in non-convex exploration. 712

Related Discrete Methods have been made in modeling and optimizing discrete distributions. Zhang
et al. (2022a) introduced energy-based generative flow networks to amortize expensive MCMC
exploration into a fixed number of actions. Discrete Diffusion Models (Sun et al., 2023b; Lou et al., 2024) extended continuous-time diffusion models to discrete spaces with well-defined score functions for discrete variables. Wen et al. (2024) proposed an efficient gradient-based discrete optimization method for generative models. These approaches to discretizing continuous methods and handling complex discrete data offer valuable insights for developing DLSs.

720 Stochastic Gradient Langevin MCMC (Welling & Teh, 2011) has become a favored MCMC method in big data due to its effective transition from optimization to sampling. However, its lack 721 of adaptive step sizes to the energy curvature limits the use of this crucial information. To further 722 leverage curvature information, Quasi-Newton methods (Ahn et al., 2012; Simsekli et al., 2016) 723 exploit curvature information by adjusting step sizes, while Hamiltonian Monte Carlo (Neal, 2012; 724 Campbell et al., 2021) and higher-order approaches (Chen et al., 2015; Li et al., 2019) employ larger 725 step sizes to improve stability. These approaches, however, still encounter difficulties in avoiding 726 local traps, which is where advanced reMCMC methods help balance exploration and exploitation 727 when navigating non-convex energy landscapes.

728 729 730

B DREXEL AND DREAM WITH BINARY VARIABLES

⁷³¹ ⁷³² When the variable domain Θ is binary $\{0, 1\}^d$, Algorithm 1 can be further simplified. In this binary ⁷³³ setting, the Hadamard product, denoted by \odot , simplifies several operations. This streamlined version ⁷³⁴ demonstrates that both DREXEL and DREAM can be efficiently parallelized across CPUs and GPUs, ⁷³⁵ which leads to reduced computational cost.

Binary variables are particularly advantageous in this context. The binary domain facilitates the use of efficient bitwise operations, which not only speed up the computation but also enable the algorithm to scale better in high-dimensional spaces. Moreover, the simplicity of the binary domain reduces algorithmic complexity, which can facilitate computational efficiency.

740 741

745 746

C DREXEL AND DREAM WITH CATEGORICAL VARIABLES

We further examine how DREXEL and DREAM can be formulated for categorical variables using
one-hot vectors and ordinal integers.

In **one-hot encoding**, each categorical variable θ_i is represented as a vector in $\{0, 1\}^N$ where exactly one element is 1, and the rest are 0. The update rule for one-hot encoded variables is given by:

$$\operatorname{Categorical}\left[\operatorname{Softmax}\left(\frac{1}{2\tau_{k}}\nabla U(\boldsymbol{\theta}_{i}^{(k)})_{d}\left(\boldsymbol{\theta}_{i+1,d}^{(k)}-\boldsymbol{\theta}_{i,d}^{(k)}\right)-\frac{\left\|\boldsymbol{\theta}_{i+1,d}^{(k)}-\boldsymbol{\theta}_{i,d}^{(k)}\right\|_{2}^{2}}{2\alpha_{k}}\right)\right], \quad k = 1, 2.$$

⁷⁵² In this setting, the difference $\theta_{i+1}^{(k)} - \theta_i^{(k)}$ results in a vector with exactly two non-zero elements, which reflects a transition between categories.

For **ordinal variables**, where categories have a natural ordering, θ_i can be represented as integers in $\{0, 1, \dots, N-1\}$. The update rule becomes:

756
757
758
759
760
Categorical
$$\left[\text{Softmax} \left(\frac{1}{2\tau_k} \nabla U(\boldsymbol{\theta}_i^{(k)})_d \left(\theta_{i+1,d}^{(k)} - \theta_{i,d}^{(k)} \right) - \frac{\left(\theta_{i+1,d}^{(k)} - \theta_{i,d}^{(k)} \right)^2}{2\alpha_k} \right) \right], \quad k = 1, 2.$$

Here, the scalar difference $(\theta_{i+1,d}^{(k)} - \theta_{i,d}^{(k)})$ captures the magnitude and direction of the transition between ordered categories. This representation leverages the ordering information to inform the proposal distribution more precisely.

764 765 Algorithm 2 DREXEL or DREAM with Binary Variables. 766 **Input** Step Sizes α_1 , α_2 , Temperatures τ_1 , τ_2 , and Swap Intensity $\rho > 0$. **Input** Initial Samples $\theta_0^{(k)} \in \Theta$, k = 1, 2. 767 768 1: For $i = 1, 2, \dots, I$ do 769 2: **Sampling Steps:** 770 For k = 0, 1, 2 do: 3: Compute $P_k(\boldsymbol{\theta}_i^{(k)}) = \frac{\exp\left(-\frac{1}{2\tau_k}\nabla U(\boldsymbol{\theta}_i^{(k)})\odot(2\boldsymbol{\theta}_i^{(k)}-1)-\frac{1}{2\sigma_k\tau_k}\right)}{\exp\left(-\frac{1}{2\tau_k}\nabla U(\boldsymbol{\theta}_i^{(k)})\odot(2\boldsymbol{\theta}_i^{(k)}-1)-\frac{1}{2\sigma_k\tau_k}\right)+1}$ 771 4: 772 773 5: Sample $\boldsymbol{u} \sim U([0, 1]^d)$ 774 Set $I_k \leftarrow \dim(u \le P_k(\theta_i^{(k)}))$ 6: 775 Set $\omega^{(k)} \leftarrow \text{flipdim}(I_k)$ 7: 776 8: **End For** MH Steps (for DREAM): 777 9: 10: **For** k = 1, 2 **do:** 778 Compute $q_k(\boldsymbol{\omega}^{(k)}|\boldsymbol{\theta}_i^{(k)}) = \prod_{d \in \boldsymbol{I}_k} \mathbb{P}(\boldsymbol{\theta}_i^{(k)})_d \cdot \prod_{d \notin \boldsymbol{I}_k} \left(1 - \mathbb{P}(\boldsymbol{\theta}_i^{(k)})_d\right)$ 11: 779 Compute $q_k(\boldsymbol{\theta}_i^{(k)}|\boldsymbol{\omega}^{(k)}) = \prod_{d \in \boldsymbol{I}_k} \mathbb{P}(\boldsymbol{\omega}^{(k)})_d \cdot \prod_{d \notin \boldsymbol{I}_k} \left(1 - \mathbb{P}(\boldsymbol{\omega}^{(k)})_d\right)$ 12: Compute $P_k(\boldsymbol{\omega}^{(k)}) = \frac{\exp\left(-\frac{1}{2\tau_k}\nabla U(\boldsymbol{\omega}^{(k)}) \odot (2\boldsymbol{\omega}^{(k)} - 1) - \frac{1}{2\alpha_k\tau_k}\right)}{\exp\left(-\frac{1}{2\tau_k}\nabla U(\boldsymbol{\omega}^{(k)}) \odot (2\boldsymbol{\omega}^{(k)} - 1) - \frac{1}{2\alpha_k\tau_k}\right) + 1}$ 781 13: 782 783 14: Compute $\mathcal{A}(\boldsymbol{\omega}^{(k)}, \boldsymbol{\theta}_i^{(k)})$ follows (4) 784 Generate a number $u \sim U[0, 1]$ Set $\theta_{i+1}^{(k)} \leftarrow \omega^{(k)}$ if $u \leq \mathcal{A}$ else $\theta_{i+1}^{(k)} \leftarrow \theta_i^{(k)}$ 15: 785 16: 786 17: End For 787 18: **Swapping Steps:** 788 19: Generate a number $u \sim U[0, 1]$. Compute $\tilde{S}(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}_{i+1}^{(2)})$ follows (10) 20: 789 Swap $\boldsymbol{\theta}_{i+1}^{(1)}$ and $\boldsymbol{\theta}_{i+1}^{(2)}$ if $u \leq \rho \min\{1, \tilde{S}\}$ 790 21: 791 22: End For **Output** Samples $\{\boldsymbol{\theta}_i^{(1)}\}_{i=1}^{I}$. 792 793 794

D **THEORETICAL ANALYSIS**

In this section, we first validate the assumptions of smoothness, dissipativity, and the fine grid used in our analysis. Next, we provide an asymptotic analysis to show the weak convergence of DREXEL to the target distribution. We conclude with a detailed examination of its non-asymptotic behavior.

D.1 PROOF OF THEOREM 1

To explore the reversibility of the discrete replica exchange Langevin sampler, we follow the proof of Zhang et al. (2022b). For the discrete replica exchange Langevin sampler, we denote the update rule as follows:

$$q_1\left(\boldsymbol{\theta}' \mid \boldsymbol{\theta}^{(1)}\right) \propto \text{Categorical}\left(\text{Softmax}\left(\frac{1}{\tau_1} \nabla U(\boldsymbol{\theta}^{(1)})_d (\theta_d' - \theta_d^{(1)}) - \frac{(\theta_d' - \theta_d^{(1)})^2}{2\alpha_1}\right)\right)$$

(1) 0 >>>

$$q_2\left(\boldsymbol{\theta}' \mid \boldsymbol{\theta}^{(2)}\right) \propto \text{Categorical}\left(\text{Softmax}\left(\frac{1}{\tau_2}\nabla U(\boldsymbol{\theta}^{(2)})_d(\boldsymbol{\theta}'_d - \boldsymbol{\theta}^{(2)}_d) - \frac{(\boldsymbol{\theta}'_d - \boldsymbol{\theta}^{(2)}_d)^2}{2\alpha_2}\right)\right),$$

where $d = 1, 2, \dots, \mathbf{d}, \forall \boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)} \in \Theta, \Theta \in \mathbb{R}^{\mathbf{d}}$. $\boldsymbol{\theta}^{(1)}$ denotes the current sample from the lowtemperature sampler, and $\theta^{(2)}$ is the sample from the high-temperature sampler.

Proof. We consider the transition probability $q(\theta' | \theta^{(1)})$. Different from the direct transition in DLS, we consider two scenarios to transition from $\theta^{(1)}$ to θ' : with probability 1 - S, there is no chain swap, and the model parameter change from $\theta^{(1)}$ to θ' in the low-temperature sampler; with probability \hat{S} , there is a chain swap, and the high-temperature sampler generate new sample from $\theta^{(2)}$ to θ' . We recall the definition of the proposed swap function $S(\cdot, \cdot)$ in (10) as follows:

$$S\left(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}_{i+1}^{(2)} \mid \boldsymbol{\theta}_{i}^{(1)}, \boldsymbol{\theta}_{i}^{(2)}\right) := e^{\left(\frac{1}{\tau_{2}} - \frac{1}{\tau_{1}}\right) \left[U\left(\boldsymbol{\theta}_{i+1}^{(1)}\right) + U\left(\boldsymbol{\theta}_{i}^{(1)}\right) - U\left(\boldsymbol{\theta}_{i+1}^{(2)}\right) - U\left(\boldsymbol{\theta}_{i}^{(2)}\right) \right]}.$$

Following this, we rewrite the transition probability $q(\theta' | \theta)$ of the discrete replica exchange Langevin sampler as:

$$q\left(\boldsymbol{\theta}' \mid \boldsymbol{\theta}_{i}^{(1)}\right) = \sum_{\boldsymbol{\theta}_{i}^{(2)}} \sum_{\boldsymbol{\theta}_{i+1}^{(2)}} \pi\left(\boldsymbol{\theta}_{i}^{(2)}\right) q_{2}\left(\boldsymbol{\theta}_{i+1}^{(2)} \mid \boldsymbol{\theta}_{i}^{(2)}\right) \left[1 - S\left(\boldsymbol{\theta}', \boldsymbol{\theta}_{i+1}^{(2)}\right)\right] q_{1}\left(\boldsymbol{\theta}' \mid \boldsymbol{\theta}_{i}^{(1)}\right) + \sum_{\boldsymbol{\theta}_{i}^{(2)}} \sum_{\boldsymbol{\theta}_{i+1}^{(1)}} \pi\left(\boldsymbol{\theta}_{i}^{(2)}\right) q_{2}\left(\boldsymbol{\theta}' \mid \boldsymbol{\theta}_{i}^{(2)}\right) S\left(\boldsymbol{\theta}_{i+1}^{(1)}, \boldsymbol{\theta}'\right) q_{1}\left(\boldsymbol{\theta}_{i+1}^{(1)} \mid \boldsymbol{\theta}_{i}^{(1)}\right),$$
(11)

where the first term on the right-hand side of (11) is the probability change from θ to θ' in the low-temperature sampler, and the second term is the probability that the low-temperature sampler starts with $\theta_i^{(1)}$, the high-temperature chain ends on θ' , with a probability S to have the sample swap. To further demonstrate the reversibility of the proposed discrete replica exchange Langevin sampler, we multiply $\pi_{\alpha}(\theta)$ from both sides³:

$$\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q\left(\theta'\mid\theta_{i}^{(1)}\right) = \sum_{\theta_{i}^{(2)}}\sum_{\theta_{i+1}^{(2)}}\pi_{\alpha}\left(\theta_{i}^{(2)}\right)q_{2}\left(\theta_{i+1}^{(2)}\mid\theta_{i}^{(2)}\right)\left[1-S\left(\theta',\theta_{i+1}^{(2)}\right)\right]\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta'\mid\theta_{i}^{(1)}\right) + \sum_{\theta_{i}^{(2)}}\sum_{\theta_{i+1}^{(1)}}\pi_{\alpha}\left(\theta_{i}^{(2)}\right)q_{2}\left(\theta'\mid\theta_{i}^{(2)}\right)S\left(\theta_{i+1}^{(1)},\theta'\right)\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta_{i+1}^{(1)}\mid\theta_{i}^{(1)}\right)$$

$$= \sum_{\theta_{i}^{(2)}}\sum_{\theta_{i+1}^{(2)}}\pi_{\alpha}\left(\theta_{i}^{(2)}\right)q_{2}\left(\theta_{i+1}^{(2)}\mid\theta_{i}^{(2)}\right)\left[1-e^{\left(\frac{1}{\tau_{2}}-\frac{1}{\tau_{1}}\right)\left[U(\theta')+U(\theta_{i}^{(1)})-U(\theta_{i+1}^{(2)})-U(\theta_{i}^{(2)})\right]}\right]}\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta'\mid\theta_{i}^{(1)}\right)$$

$$+ \sum_{\theta_{i}^{(2)}}\sum_{\theta_{i+1}^{(1)}}\pi_{\alpha}\left(\theta_{i}^{(2)}\right)q_{2}\left(\theta'\mid\theta_{i}^{(2)}\right)\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta_{i+1}^{(1)}\mid\theta_{i}^{(1)}\right)e^{\left(\frac{1}{\tau_{2}}-\frac{1}{\tau_{1}}\right)\left[U(\theta_{i+1}^{(1)})+U(\theta_{i}^{(1)})-U(\theta_{i}^{(2)})\right]},$$

$$(12)$$

³We ignore the subscript of the step size from α_1 to α for the stationary distribution π_{α} for simplicity.

where (a) replace the swap function with (10).

Recall the assumption that the target distribution is defined as log-quadratic $\pi(\theta) = \exp(\theta^{T} J \theta + b^{T} \theta)/Z$, where $J \in \mathbb{R}^{d \times d}$ is a symmetric matrix, $b \in \mathbb{R}^{d}$ is a vector, and Z normalizes the distribution. We then have $\nabla U(\theta) = 2J^{T}\theta + b$ and $\nabla^{2} U(\theta) = 2J$ is a constant.

We further denote $Z_{\alpha}(\theta) = \sum_{\mathbf{x}} \exp\left[\frac{1}{2}(U(\mathbf{x}) - U(\theta)) - (\mathbf{x} - \theta)^{\top}\left(\frac{1}{2\alpha}\mathbf{I} + \frac{1}{2}\mathbf{J}\right)(\mathbf{x} - \theta)\right]$ and $\pi_{\alpha}(\theta) = Z_{\alpha}(\theta)\pi(\theta) / \sum_{\mathbf{x}} Z_{\alpha}(\mathbf{x})\pi(\mathbf{x})$. According to Theorem 1 in Zhang et al. (2022b), we have the transition from $\theta_i^{(2)}$ to $\theta_{i+1}^{(2)}$ multiplying the stationary distribution $\pi_{\alpha}\left(\theta_i^{(2)}\right)$ as:

$$\pi_{\alpha}\left(\boldsymbol{\theta}_{i}^{(2)}\right)q_{1}\left(\boldsymbol{\theta}_{i+1}^{(2)} \mid \boldsymbol{\theta}_{i}^{(2)}\right) = \frac{Z_{\alpha}\left(\boldsymbol{\theta}_{i}^{(2)}\right)\pi_{\alpha}\left(\boldsymbol{\theta}_{i}^{(2)}\right)}{\sum_{\boldsymbol{x}}Z(\boldsymbol{x})\pi_{\alpha}(\boldsymbol{x})} \cdot \frac{\exp\left[\frac{1}{2}\left(U\left(\boldsymbol{\theta}_{i+1}^{(2)}\right) - U(\boldsymbol{\theta}_{i}^{(2)})\right) - \left(\boldsymbol{\theta}_{i+1}^{(2)} - \boldsymbol{\theta}_{i}^{(2)}\right)^{\mathsf{T}}\left(\frac{1}{2\alpha}\mathbf{I} + \frac{1}{2}\boldsymbol{J}\right)\left(\boldsymbol{\theta}_{i+1}^{(2)} - \boldsymbol{\theta}_{i}^{(2)}\right)\right]}{Z_{\alpha}\left(\boldsymbol{\theta}_{i}^{(2)}\right)} \\ = \frac{\exp\left[\frac{1}{2}\left(U\left(\boldsymbol{\theta}_{i+1}^{(2)}\right) + U(\boldsymbol{\theta}_{i}^{(2)})\right) - \left(\boldsymbol{\theta}_{i+1}^{(2)} - \boldsymbol{\theta}_{i}^{(2)}\right)^{\mathsf{T}}\left(\frac{1}{2\alpha}\mathbf{I} + \frac{1}{2}\boldsymbol{J}\right)\left(\boldsymbol{\theta}_{i+1}^{(2)} - \boldsymbol{\theta}_{i}^{(2)}\right)\right]}{Z \cdot \sum_{\boldsymbol{x}}Z(\boldsymbol{x})\pi_{\alpha}(\boldsymbol{x})},$$

which is symmetric. Similarly, we can expand the following distributions according to the definition of log-quadratic targets and the transition probabilities $\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta' \mid \theta_{i}^{(1)}\right), \pi_{\alpha}\left(\theta_{i}^{(2)}\right)q_{2}\left(\theta' \mid \theta_{i}^{(2)}\right)$, and $\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q_{1}\left(\theta_{i+1}^{(1)} \mid \theta_{i}^{(1)}\right)$. Similarly, all of them are symmetric, which indicates that $\pi_{\alpha}\left(\theta_{i}^{(1)}\right)q\left(\theta' \mid \theta_{i}^{(1)}\right)$ is also symmetric. Therefore, we conclude that $q\left(\theta' \mid \theta_{i}^{(1)}\right)$ give in (12) is reversible and the stationary distribution is $\pi_{\alpha}\left(\theta_{i}^{(1)}\right)$.

To further prove the stationary distribution π_{α} converges weakly to the target π as the step sizes are close to zero, we first observe that for any $\theta \in \Theta$,

$$Z_{\alpha}(\boldsymbol{\theta}) = \sum_{\boldsymbol{x} \in \boldsymbol{\theta}} \exp\left(\frac{1}{2} \left(U(\boldsymbol{x}) - U(\boldsymbol{\theta})\right) - (\boldsymbol{x} - \boldsymbol{\theta})^{\top} \left(\frac{1}{2\alpha} \mathbf{I} + \frac{1}{2} \boldsymbol{J}\right) (\boldsymbol{x} - \boldsymbol{\theta})\right).$$

As $\alpha \to 0$, the term involving $\frac{1}{2\alpha}$ I dominates the quadratic form unless $x = \theta$. Specifically, for $x \neq \theta$,

$$(\boldsymbol{x} - \boldsymbol{\theta})^{\top} \left(\frac{1}{2\alpha} \mathbf{I} + \frac{1}{2} \boldsymbol{J} \right) (\boldsymbol{x} - \boldsymbol{\theta}) \geq \frac{1}{2\alpha} \| \boldsymbol{x} - \boldsymbol{\theta} \|^{2},$$

which tends to infinity as $\alpha \to 0$. Therefore, the terms in the sum for $x \neq \theta$ become negligible, which means $\exp\left(-(x-\theta)^{\top}\left(\frac{1}{2\alpha}\mathbf{I}+\frac{1}{2}\boldsymbol{J}\right)(x-\theta)\right) \to 0$ as $\alpha \to 0$. For $x = \theta$, the exponent simplifies to zero: $\exp\left(\frac{1}{2}(U(\theta)-U(\theta))-0\right) = 1$. Thus, we can conclude that $\lim_{\alpha\to 0} Z_{\alpha}(\theta) = 1$.

We denote the denominator of $\pi_{\alpha}(\theta)$ is $D_{\alpha} = \sum_{x \in \Theta} Z_{\alpha}(x)\pi(x)$. Following the above derivation, we can easily find that $Z_{\alpha}(x) \to 1$ as $\alpha \to 0$ for each $x \in \Theta$. Therefore, we also have $\lim_{\alpha \to 0} D_{\alpha} = \sum_{x \in \theta} \pi(x) = 1$.

Combining the above results, we have $\lim_{\alpha \to 0} \pi_{\alpha}(\theta) = \lim_{\alpha \to 0} \frac{Z_{\alpha}(\theta)\pi(\theta)}{D_{\alpha}} = \frac{\pi(\theta)}{1} = \pi(\theta)$. Thus, we derive a conclusion that $\pi_{\alpha}(\theta)$ converges point-wisely to $\pi(\theta)$ as $\alpha \to 0$.

Since π_{α} and π are probability mass functions on a discrete space θ , and $\pi_{\alpha}(\theta) \to \pi(\theta)$ point-wisely, according to Scheffé's Lemma (Billingsley, 2017), it further implies that $\lim_{\alpha \to 0} \sum_{\theta \in \Theta} |\pi_{\alpha}(\theta) - \pi(\theta)| = 0$.

This convergence implies weak convergence of π_{α} to π : from Dominated Convergence Theorem (Folland, 1999), we have for any bounded function $f : \Theta \to \mathbb{R}$,

$$\lim_{\alpha \to 0} \sum_{\theta \in \Theta} f(\theta) \pi_{\alpha}(\theta) = \sum_{\theta \in \theta} f(\theta) \pi(\theta)$$

911 This completes the proof that π_{α} converges weakly to π as $\alpha \to 0$.

909 910

912

913 914

915 916

870

871

877 878

879

880

881

882 883

884

885

890 891 892

E Additional Experimental Results

917 The experiments were run on a server featuring an Intel(R) Core(TM) i9-14900K processor, RTX 4090 GPUs, and 128 GB DDR4 memory.

918 E.1 SAMPLING FROM 2D SYNTHETIC ENERGIES 919

924

925

926

932 933

934 935

944

945

946

947 948 949

955

961 962

971

To evaluate the efficiency of MCMC algorithms in exploring non-convex discrete energy landscapes,
we consider a set of energy functions that present varying degrees of complexity and multimodality.
These functions are designed to test the algorithms' ability to navigate challenging landscapes
characterized by multiple local minima, sharp ridges, and disconnected modes.

Wave energy function is a periodic sinusoidal surface with alternating peaks and valleys:

 $U(x, y) = \sin(3x)\sin(3y).$

Its highly rugged landscape arises from repeated oscillations, which makes it a suitable test for an algorithm's ability to navigate sharp and oscillating energy contours without getting trapped in local minima.

Eight Gaussian energy function (Dai et al., 2020) consists of eight equally spaced Gaussian components arranged in a circular pattern:

$$U(x, y) = \sum_{i=1}^{8} \frac{-\left[(x - x_i)^2 + (y - y_i)^2\right]}{2\sigma^2}$$

where the centers (x_i, y_i) are located at $(\pm 1, 0)$, $(0, \pm 1)$, and $\left(\pm \frac{\sqrt{2}}{2}, \pm \frac{\sqrt{2}}{2}\right)$, with $\sigma = 1.0$ for all modes. This function presents multiple well-separated basins of attraction, testing the algorithm's ability to explore disconnected modes effectively.

Moon energy function creates a complex, asymmetric landscape resembling a crescent shape:

$$U(x,y) = -\frac{1}{10}y^4 - \frac{1}{2}\left(4x - y^2 + \frac{24}{5}\right)^2.$$

With deep valley and steep ridge features, this non-convex structure evaluates the algorithm's capacity to explore non-uniform, curved regions and traverse narrow channels between high-energy barriers.

Two Moons energy function describes a landscape with two prominent crescent-shaped modes separated by a low-energy region:

$$U(x,y) = -\frac{2}{25} \left(x^2 + y^2 - 2 \right)^2 + \log \left[e^{-\frac{1}{2} \left(\frac{5x-4}{4} \right)^2} + e^{-\frac{1}{2} \left(\frac{5x+4}{4} \right)^2} \right]$$

This function challenges MCMC methods to jump between distinct modes, and tests their efficiencyin exploring multi-modal distributions where modes are not directly connected.

Twist energy function represents a twisted sinusoidal landscape where energy levels change smoothly along a sinusoidal curve: $1 \int (\pi x)^2$

$$U(x, y) = -\frac{1}{2} \left[y - \sin\left(\frac{\pi x}{2}\right) \right]^2$$

The narrow, twisted valleys require careful gradient following, which tests the algorithm's ability tosample from highly structured, nonlinear regions of the energy landscape.

Flower energy function combines radial symmetry with angular variations, forming a complex landscape with a petal-like structure:

$$U(x, y) = \sin\left(\sqrt{x^2 + y^2}\right) + \cos\left(5\tan\left(\frac{y}{x}\right)\right).$$

With multiple local minima and ridges radiating from the center, this intricate landscape tests the
 algorithm's capability to explore multi-modal, rotationally symmetric energy landscapes with sharp
 transitions between regions.

In our experiments, we evaluate different samplers using Kullback-Leibler (KL) divergence and
 Maximum Mean Discrepancy (MMD).

KL divergence measures the difference between two probability distributions. Given two distributions π and $\tilde{\pi}$, the KL divergence is defined as:

$$D_{KL}(\pi \parallel ilde{\pi}) = \sum_{m{ heta} \in \Theta} \pi(m{ heta}) \log rac{\pi(m{ heta})}{ ilde{\pi}(m{ heta})}$$



Figure 5: Visualization of the true energy function and the empirical energy function yield by DMALA, AB, ACS, and DREAM. Seven energy functions are tested here, which include wave, eight Gaussians, sixteen Gaussians, moon, two moons, twist, and flower energy functions. Red colors denote high-density regions, and blue colors represent low-density regions.

1001 where $\pi(\theta)$ represents the probability of θ under the target distribution, and $\tilde{\pi}(\theta)$ represents the 1002 probability of θ under the empirical distribution from the samplers. This metric quantifies how much 1003 information is lost when $\tilde{\pi}$ is used to approximate π , with lower values indicating better performance.

1004 **MMD** is a kernel-based test used to compare distributions. It is computed as:

$$\mathrm{MMD}^{2}(\pi,\tilde{\pi}) = \mathbb{E}_{\boldsymbol{x},\boldsymbol{x}'\sim\pi}[k(\boldsymbol{x},\boldsymbol{x}')] + \mathbb{E}_{\boldsymbol{y},\boldsymbol{y}'\sim\tilde{\pi}}[k(\boldsymbol{y},\boldsymbol{y}')] - 2\mathbb{E}_{\boldsymbol{x}\sim\pi,\boldsymbol{y}\sim\tilde{\pi}}[k(\boldsymbol{x},\boldsymbol{y})],$$

where k(x, y) is a positive-definite kernel function. MMD measures the similarity between the empirical distributions of the generated and target samples.

In practice, however, directly computing MMD is computationally expensive. Therefore, we use an approximation based on Random Fourier Features (RFF) (Rahimi & Recht, 2007).

For two distributions π and $\tilde{\pi}$, we first map the data samples $X \sim \pi$ and $Y \sim \tilde{\pi}$ to a new feature space 1012 using the random Fourier transformation: $\phi(X) = \sqrt{\frac{2}{D}}\cos(WX^T + \tilde{b})$, where $W \in \mathbb{R}^{D \times \mathbf{d}}$ are random 1013 1014 Gaussian variables sampled from $\mathcal{N}(0, 1/\tilde{\sigma}^2)$, and \tilde{b} are random uniform variables in the range $[0, 2\pi]$. 1015 The parameter $\tilde{\sigma}$ controls the kernel bandwidth, and D is the number of random features. Once 1016 mapped, the empirical mean feature embeddings for X and Y are computed for both distributions 1017 $\mu_X = \frac{1}{n} \sum_{i=1}^n \phi(X_i), \ \mu_Y = \frac{1}{m} \sum_{i=1}^m \phi(Y_i)$. Finally, the MMD is approximated by the squared difference 1018 of the mean embeddings: 1019 $\text{MMD}^2(\pi, \tilde{\pi}) \approx \|\mu_X - \mu_Y\|^2$.

1020

996

997

998

999 1000

This approach allows us to efficiently compute the MMD between two distributions using RFFs.

1022 To evaluate the effectiveness of the proposed sampler, we explore its performance on a set of non-1023 convex discrete energy landscapes that vary in complexity and multimodality. We further compare it 1024 with baselines such as DMALA, ACS, and AB. Unless specified otherwise, the default temperature 1025 for each sampler is set at 1.0. DMALA is implemented with a step size of 0.15. AB is used with 1026 parameters $\sigma = 0.10$ and $\alpha = 0.50$. For ACS, a cyclical step size scheduler with an initial step size 1026 of 0.60 across 10 cycles is applied. DREAM uses small and large step sizes of 0.15 and 0.60 and temperatures of 1.0 and 5.0.

Figure 5 illustrates a comparison between the empirical distributions obtained from different samplers 1029 and the ground truth (top row). A clear distinction can be observed between the discrete samplers 1030 in terms of their ability to capture the full complexity of the landscape. From the figure, DREAM 1031 produces the most balanced and comprehensive empirical distribution, which captures all significant 1032 modes of the energy landscape. The improvements are significant in tasks of approximating wave 1033 and multi-Gaussian energy functions. While other discrete samplers (such as DMALA, AB, ACS) 1034 fail to capture all modes, DREAM exhibits the most comprehensive exploration, as reflected in the 1035 uniformity of the empirical distribution across all modes. Its ability to distribute samples effectively, 1036 even in the presence of disconnected modes and sharp energy barriers, demonstrates its robustness in navigating complex discrete energy landscapes. By contrast, DMALA shows a heavy concentration 1037 of samples around certain modes, which indicates that it struggles to escape local minima. This 1038 leads to poor coverage of the landscape and a lack of diversity in the sampled regions. ACS and AB 1039 perform better in terms of covering multiple modes but still show uneven sample distributions. Some 1040 modes are under-sampled, while others are over-sampled, particularly in regions with shallow energy 1041 gradients. 1042

1042

E.2 SAMPLING FROM ISING MODELS

We sampled from the Ising model using multiple samplers with a default temperature of 1.0, unless otherwise specified. DULA utilized a step size of 0.20, while DMALA had a step size of 0.40. For ACS, we employed a cyclical step size scheduler with 10 cycles and an initial step size of 0.30. ACS with MH corrections used an initial step size of 5.0. For DREXEL, small and large step sizes were set at 0.15 and 0.50, with temperatures at 1.0 and 5.0. DREAM followed a similar temperature schedule, with a small step size of 0.35 and a large size of 0.50.

¹⁰⁵¹ In this study, different samplers are evaluated with log Root Mean Square Error (log RMSE). Log RMSE evaluates prediction accuracy for comparing the true values from π and the predictions from $\tilde{\pi}$, which can be adapted as:

1054 1055

1056

log RMSE = log $\left(\sqrt{\frac{1}{n}\sum_{i=1}^{n}(\pi(\boldsymbol{x}_{i})-\tilde{\pi}(\boldsymbol{x}_{i}))^{2}}\right)$,

1058 where $\pi(\mathbf{x}_i)$ is the true value under the target distribution, and $\tilde{\pi}(\mathbf{x}_i)$ is the corresponding approximation 1059 from the empirical distribution.

We further examine the influence of bias correction terms in DREXEL and DREAM with bias corrections (bDREXEL and bDREAM) on the experimental results of Ising models. For comparison, we also evaluate DREXEL and DREAM, which incorporate historical energy corrections. Specifically, we modify the swap function (10) by directly adding bias corrections and adjusting σ^2 in (9). Each experiment is repeated 20 times with different random seeds, and the average and standard deviation of log RMSE are reported. The results are shown in Figure 6.

The results demonstrate that DREAM achieves the lowest log RMSE values when corrections are small, which indicates the minimal influence of bias correction on its performance. In contrast, the performance of bDREXEL and bDREAM becomes better with increasing σ^2 , which suggests its performance may heavily rely on a good selection of bias correction. bDREAM and DREAM show comparable log RMSEs, with log RMSE across different bias correction magnitudes. However, DREAM yields the lowest log RMSE with no correction, while bDREAM has the lowest one when the correction increases.

1073

1074 E.3 SAMPLING FROM RESTRICTED BOLTZMANN MACHINES

We trained RBMs with the Adam optimizer with a learning rate of 0.001, over 1,000 iterations, and a batch size of 128. For training, we used contrastive divergence (CD), which approximates the loglikelihood gradient by performing k = 10 Gibbs sampling steps. The gradient of the log-likelihood for an RBM is given by:

$$\nabla_{\theta} \log P_{\theta}(\boldsymbol{x}) = \mathbb{E}_{P_{\text{data}}}[\nabla_{\theta} \log P_{\theta}(\boldsymbol{x})] - \mathbb{E}_{P_{\theta}}[\nabla_{\theta} \log P_{\theta}(\boldsymbol{x})],$$



Figure 6: Log RMSE results for bDREXEL, bDREAM, DREXEL, and DREAM under varying
 correction terms. Solid and dashed lines denote the log RMSE average values, and a lighter shade
 represents 95% confidence intervals. Lower log RMSE values indicate better performance.

1100

where x is the visible layer, and θ are the model parameters. The first term corresponds to the data distribution, while the second term is the expectation under the model's distribution. It should be noted that direct computation of the model distribution expectation is expensive, and CD approximates the second term by running *k*-step Gibbs sampling to obtain samples from the model. This approximation enables efficient training of RBMs by focusing on the contrast between the observed and modeled distributions.

To sample from the RBMs, we employed several discrete samplers at a default temperature of 1.0 unless otherwise noted. DMALA used a step size of 0.15, and ACS applied a cyclical step size scheduler with 10 cycles, starting at a step size of 0.50. The Any-scale Balanced (AB) sampler was configured with $\sigma = 0.10$ and $\alpha = 0.50$. For DREAM, the small and large step sizes were set between 0.15–0.20 and 0.40–0.50, with temperatures set at 1.0 and 2.0.

1112 1113

1114

1121

1131

1132

E.4 LEARNING DEEP ENERGY-BASED MODELS

We trained Deep EBMs using a ResNet-64 backbone and optimized the model with the Adam optimizer at a fixed learning rate of 0.001 without gradient clipping. The model was trained for 50,000 iterations with a batch size of 256. we employed Persistent Contrastive Divergence (PCD) to approximate the intractable likelihood gradient, which builds upon standard contrastive divergence by maintaining persistent Markov chains throughout training. It allows for more stable and accurate sampling. Specifically, the model's log-likelihood gradient is given by:

$$\nabla_{\theta} \log P_{\theta}(\mathbf{x}) = \mathbb{E}_{P_{\text{data}}} [\nabla_{\theta} \log P_{\theta}(\mathbf{x})] - \mathbb{E}_{P_{\theta}} [\nabla_{\theta} \log P_{\theta}(\mathbf{x})]$$

where the second term (the model expectation) is intractable. PCD approximates this by updating samples across training iterations via Gibbs sampling, which ensures that the Markov chain does not restart after each parameter update. Additionally, a replay buffer containing 1,000 past samples is used to further stabilize training. The buffer stores past model samples and reuses them to reduce variance, thus improving both the efficiency and stability of the learning process.

To evaluate Deep EBMs, we applied Annealed Importance Sampling (AIS) with DULA to estimate the test log-likelihoods. AIS is a technique used to estimate partition functions by smoothly interpolating between a known distribution and the target distribution. This is achieved by introducing a sequence of intermediate distributions:

$$P_t(\boldsymbol{x}) = \frac{1}{Z_t} \exp(-\beta_t E(\boldsymbol{x})),$$

where *t* denotes the current annealing step, $E(\mathbf{x})$ is the energy function, Z_t is the partition function, and β_t is a temperature that gradually transitions between 0 and 1 over the course of the annealing 1134 process. When $\beta_t = 0$, the intermediate distribution is identical to the proposal distribution (which we 1135 can sample from easily). When $\beta_t = 1$, the intermediate distribution becomes the target distribution, 1136 which is more complex and generally intractable to sample directly. To adjust β_t , we typically choose 1137 a monotonic schedule that increases smoothly from 0 to 1 over the course of the AIS process. A 1138 common choice is a linear interpolation ($\beta_t = t/T, t = 0, 1, 2, ..., T$) or an exponential schedule $(\beta_t = (t/T)^2)$, where β_t increases evenly across T annealing steps. AIS computes an estimate of the 1139 partition function by sampling from these intermediate distributions and adjusting the importance 1140 weights over time: 1141

- 1142
- 1143

1150

 $\tilde{Z}_{\theta} = Z_0 \prod_{t=1}^{T} \frac{P_t(\boldsymbol{x})}{P_{t-1}(\boldsymbol{x})},$

1144 where Z_0 is an initial distribution that is easy to sample from, which serves as a starting point for the 1145 annealing process. Typically, Z₀ is chosen to be the partition function of a simple proposal distribution 1146 $p_0(x)$, which is often a uniform distribution or a Gaussian distribution with parameters that are easy 1147 to compute. In our experiments, we used AIS with 40 samples and 30,000 annealing steps. DULA 1148 was configured with a step size of 0.08 and a temperature of 1.00. Detailed hyperparameters for 1149 training Deep EBMs are listed in Table 3.

Static MNIST	DULA	DMALA	bDREXEL	bDREAM	DREXEL	DREAN
Step size	0.08	0.10	0.05 0.15	0.05 0.15	0.11 0.25	0.10 0.30
Temperature	1.0	1.0 -	1.0 5.0	1.0 5.0	1.0 5.0	1.0 5.0
Correction	-	-	0.00	0.00	0.00	0.00
Dunamia MNIST			LDDEVEI	LDDEAM	DDEVEI	
	DULA	DWALA	UDREALL	UDREAM	DREALL	DREAD
Step size	0.08	0.10	$0.05 \\ 0.15$	0.05 0.15	0.11 0.25	0.11 0.25
Temperature	1.0	1.0	1.0 5.0	1.0 5.0	1.0 5.0	1.0 5.0
Correction	-	-	0.00	0.00	0.00	0.00
Omniglot	DULA	DMALA	bDREXEL	bDREAM	DREXEL	DREAN
Step size	0.08	0.10	0.05 0.15	0.05 0.15	0.08 0.15	0.08 0.15
Temperature	1.0	1.0	1.0 5.0	1.0 5.0	1.0 5.0	1.0 5.0
Correction	-	-	0.00	0.00	1.00	0.00
Caltech	DULA	DMALA	bDREXEL	bDREAM	DREXEL	DREAN
Step size	0.08	0.10	0.05 0.15	0.05 0.15	0.08 0.20	0.08 0.20
Temperature	1.0	1.0	1.0 5.0	1.0 5.0	1.0 5.0	1.0 5.0



¹¹⁸³ Table 3: Hyper-parameters used in learning Deep EBMs. From top to bottom, hyper-parameters in Static MNIST, Dynamic MNIST, Omniglot, and Caltech Silhouettes are recorded. 1184

1185

For a consistent comparison with previous works (Zhang et al., 2022b), we record its log-likelihood 1186 on the test set after 50,000 iterations. In general, Table 4 yields consistently lower log-likelihood 1187 than the results in Table 2 since they train with more iterations. But similar trends are shown in Table 2 as well: DREAM consistently achieved the highest log-likelihoods across all datasets, with significant improvements on Omniglot and Caltech Silhouettes. This demonstrates the proposed swap mechanism in (10) effectively corrects imbalance, which leads to improved log-likelihood estimates across diverse image datasets. For the MNIST datasets, both DREXEL and DREAM showed competitive performance. bDREXEL and bDREAM generally perform worse than DREAM, with consistently lower log-likelihoods, which further confirms the advantage of historical energy corrections. These findings suggest that incorporating MH steps is crucial for enhancing performance in discrete sampling tasks. DULA and DMALA exhibit the lowest performance overall, which emphasizes the benefits of MH steps and the need to consider DREXEL and DREAM to enhance exploration.

Dataset	DULA	DMALA	bDREXEL	bDREAM	DREXEL	DREAM
Static MNIST	-79.672	-77.581	-77.212	-76.840	-75.685	-74.883
Dynamic MNIST	-81.144	-79.411	-81.273	-81.043	-71.091	-70.905
Omniglot	-114.203	-109.095	-94.382	-90.807	-89.971	-89.643
Caltech Silhouettes	-102.546	-98.554	-96.073	-93.969	-89.764	-86.624

Table 4: EBM learning results (log-likelihood) on the test set after 50,000 iterations.

Here we provide the generated results (Figure 7) from DREAM across Static MNIST, Dynamic MNIST, Omniglot, and Caltech Silhouettes. These images demonstrate the ability of trained deep EBMs to capture the underlying data distribution. The deep EBM is capable of producing high-quality samples that visually resemble the training data, which indicates that the learned energy function effectively models the complex, high-dimensional structure of the data.

