

CYCLICAL STOCHASTIC GRADIENT MCMC FOR BAYESIAN DEEP LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

The posteriors over neural network weights are high dimensional and multimodal. Each mode typically characterizes a meaningfully different representation of the data. We develop Cyclical Stochastic Gradient MCMC (SG-MCMC) to automatically explore such distributions. In particular, we propose a cyclical stepsize schedule, where larger steps discover new modes, and smaller steps characterize each mode. We prove non-asymptotic convergence theory of our proposed algorithm. Moreover, we provide extensive experimental results, including ImageNet, to demonstrate the effectiveness of cyclical SG-MCMC in learning complex multimodal distributions, especially for fully Bayesian inference with modern deep neural networks.

1 INTRODUCTION

Deep neural networks are often trained with stochastic optimization methods such as stochastic gradient descent (SGD) and its variants. Bayesian methods provide a principled alternative, which account for model uncertainty in weight space (MacKay, 1992; Neal, 1996), and achieve an automatic balance between model complexity and data fitting. Indeed, Bayesian methods have been shown to improve the generalization performance of DNNs (Hernández-Lobato & Adams, 2015; Blundell et al., 2015; Li et al., 2016a), while providing a principled representation of uncertainty on predictions which is crucial for decision making.

Approximate inference for Bayesian deep learning has typically focused on deterministic approaches, such as variational methods (Hernández-Lobato & Adams, 2015; Blundell et al., 2015). By contrast, MCMC methods are now essentially unused for inference with modern deep neural networks, despite previously providing the gold standard of performance with smaller neural networks (Neal, 1996). Stochastic gradient Markov Chain Monte Carlo (SG-MCMC) methods (Welling & Teh, 2011; Chen et al., 2014; Ding et al., 2014; Li et al., 2016a) provide a promising direction for a sampling based approach to inference in Bayesian deep learning. Indeed, it has been shown that stochastic methods, which use mini-batches of data, are crucial for finding weight parameters that provide good generalization in modern deep neural networks (Keskar et al., 2016).

However, SG-MCMC algorithms for inference with modern neural networks face several challenges: (i) In theory, SG-MCMC asymptotically converges to target distributions via a decreasing stepsize scheme, but suffers from a bounded estimation error in limited time (Teh et al., 2016; Chen et al., 2015). (ii) In practice, empirical successes have been reported by training DNNs in relatively short time (Li et al., 2016b; Chen et al., 2014; Gan et al., 2016; Neelakantan et al., 2016; Saatchi & Wilson, 2017). For example, (Saatchi & Wilson, 2017) apply SG-MCMC to generative adversarial networks (GANs) to solve the mode collapse problem and capture diverse generation styles. However, the loss surface for DNNs is highly multimodal (Auer et al., 1996; Choromanska et al., 2015). In order for MCMC to be effective for posterior inference in modern neural networks, a crucial question remains: how do we make SG-MCMC efficiently explore a highly multimodal parameter space given a practical computational budget?

Several attempts have been made to improve the sampling efficiency of SG-MCMC. Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) (Chen et al., 2014) introduces the momentum variable to the Langevin dynamics. The preconditioned stochastic gradient Langevin dynamics (pSGLD) (Li et al., 2016a) adaptively adjust the sampler’s step size according to the local geometry of parameter space. Though simple and promising, these methods are still inefficient at exploring multimodal

distributions in practice. It is our contention that this limitation arises from difficulties escaping local modes when using the small stepsizes that SG-MCMC methods typically require. Note that the stepsize in SG-MCMC controls the sampler’s behavior in two ways: the magnitude to deterministically drift towards high density regions *wrt.* the current stochastic gradient, and the level of injecting noise to randomly explore the parameter space. Therefore, a small stepsize reduces both abilities, resulting in a large numbers of iterations for the sampler to move across the modes.

In this paper, we propose to replace the traditional decreasing stepsize schedule in SG-MCMC with a cyclically changing one. To note the distinction from traditional SG-MCMC, we refer to this as the *Cyclical SG-MCMC* method (cSG-MCMC). The comparison is illustrated in Figure 1. The blue curve is the traditional decay, while the red curve shows the proposed cyclical schedule. Cyclical SG-MCMC operates in two stages: (i) *Exploration*: when the stepsize is large (dashed red curves), we consider this stage as an effective burn-in mechanism, encouraging the sampler to take large moves and leave the local mode using the stochastic gradient. (ii) *Sampling*: when the stepsize is small (solid red curves), the sampler explores one local mode. We collect samples for local distribution estimation during this stage. Further, we propose two practical techniques to improve the estimation efficiency: (1) a system temperature for exploration and exploitation; (2) A weighted combination scheme for samples collected in different cycles accommodates their relative importance.

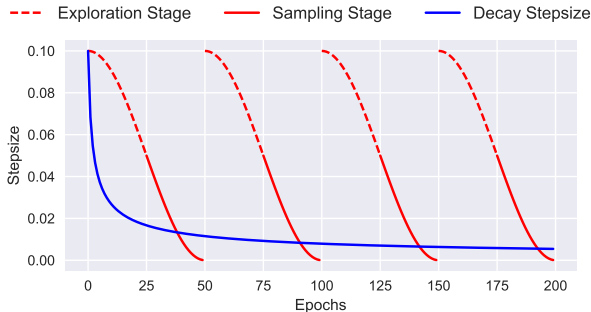


Figure 1: Illustration of the proposed cyclical stepsize schedule (red) and the traditional decreasing stepsize schedule (blue) for SG-MCMC algorithms.

This procedure can be viewed as SG-MCMC with warm restarts: the exploration stage provides the warm restarts for its following sampling stage. cSG-MCMC combines the advantages from (1) the traditional SG-MCMC to characterize the fine-scale local density of a distribution and (2) the cyclical schedule in optimization to efficiently explore multimodal posterior distributions of the parameter space. In limited time, cSG-MCMC is a practical tool to provide significantly better mixing than the traditional SG-MCMC for complex distributions. cSG-MCMC can also be considered as an *efficient* approximation to parallel MCMC; cSG-MCMC can achieve similar performance to parallel MCMC with only a fraction of cost (reciprocal to the number of chains) that parallel MCMC requires.

To support our proposal, we also prove the non-asymptotic convergence for the cyclical schedule. We note that this is the first convergence analysis of a cyclical stepsize algorithm (including work in optimization). Moreover, we provide extensive experimental results to demonstrate the advantages of cSG-MCMC in sampling from multimodal distributions, including Bayesian neural networks and uncertainty estimation on several large and challenging datasets such as ImageNet.

In short, cSG-MCMC provides a simple and automatic approach to inference in modern Bayesian deep learning, with promising results, and theoretical support. This work is a step towards enabling MCMC approaches in Bayesian deep learning. We will release the code on GitHub.

2 PRELIMINARIES: SG-MCMC WITH A DECREASING STEPSIZE SCHEDULE

SG-MCMC is a family of scalable sampling methods that enables inference with mini-batches of data. For a dataset $D = \{d_i\}_{i=1}^N$ and a θ -parameterized model, we have the likelihood $p(D|\theta)$ and prior $p(\theta)$. The posterior distribution is $p(\theta|D) \propto \exp(-U(\theta))$; where $U(\theta)$ is the potential energy given by $U(\theta) = -\log p(D|\theta) - \log p(\theta)$:

When D is too large, it is expensive to evaluate $U(\theta)$ for all the data points at each iteration. Instead, SG-MCMC methods use a minibatch to approximate $U(\theta)$: $\tilde{U}(\theta) = -\frac{N^0}{N} \sum_{i=1}^{N^0} \log p(x_{ij}|\theta)$ $\log p(\theta)$; where $N^0 \ll N$ is the size of minibatch. We recommend (Ma et al., 2015) for a general review of SG-MCMC algorithms. We describe two SG-MCMC algorithms considered in this paper.

SGLD & SGHMC (Welling & Teh, 2011) proposed Stochastic Gradient Langevin Dynamics (SGLD), which uses stochastic gradients with Gaussian noise. Posterior samples are updated at the k -th step as: $\theta_k = \theta_{k-1} - \eta \nabla \tilde{U}(\theta_{k-1}) + \sqrt{\frac{\eta}{2}} \epsilon_k$, where η is the stepsize and ϵ_k has a standard Gaussian distribution.

To improve mixing over SGLD, Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) (Chen et al., 2014) introduces an auxiliary momentum variable v . SGHMC is built upon HMC, with an additional friction term to counteract the noise introduced by a mini-batch. The update rule for posterior samples is: $\theta_k = \theta_{k-1} + v_{k-1}$, and $v_k = v_{k-1} - \eta \nabla \tilde{U}(\theta_{k-1}) - v_{k-1} + \sqrt{2\hat{\sigma}^2} \epsilon_k$, where 1 is the momentum term and $\hat{\sigma}$ is the estimate of the noise.

To guarantee asymptotic consistency with the true distribution, SG-MCMC requires that the step sizes satisfy the following assumption:

Assumption 1. The step sizes η_k are decreasing, i.e., $0 < \eta_{k+1} < \eta_k$, with 1) $\sum_{k=1}^{\infty} \eta_k = \infty$; and 2) $\sum_{k=1}^{\infty} \eta_k^2 < \infty$.

Without a decreasing step-size, the estimation error from numerical approximations is asymptotically biased. One typical decaying step-size schedule is $\eta_k = \frac{a}{b+k}$, with $a \in (0.5; 1]$ and $(a; b)$ some positive constants (Welling & Teh, 2011).

3 CYCLICAL SG-MCMC

We now introduce our *cyclical SG-MCMC* (cSG-MCMC) algorithm. cSG-MCMC consists of two stages: *exploration* and *sampling*. In the following, we first introduce the cyclical step-size schedule, and then describe the exploration stage in Section 3.1 and the sampling stage in Section 3.2. We propose an approach to combining samples for testing in Section F.

Assumption 1 guarantees the consistency of our estimation with the true distribution in asymptotic time. The approximation error in limited time is characterized as the risk of an estimator $R = B^2 + V$, where B is the bias and V is the variance. In the case of infinite computation time, the traditional SG-MCMC setting can reduce the bias and variance to zero. However, the time budget is often limited in practice, and there is always a trade-off between bias and variance. We therefore decrease the overall approximation error R by reducing the variance through obtaining more effective samples. The effective sample size can be increased if less correlated samples from different distribution modes are collected.

For deep neural networks, the parameter space is highly multimodal. SG-MCMC with the traditional decreasing stepsize schedule in practice becomes trapped in a local mode, though injecting noise may help the sampler to escape in the asymptotic regime (Zhang et al., 2017). Inspired to improve the exploration of the multimodal posteriors for deep neural networks, with a simple and automatic approach, we propose the cyclical cosine stepsize schedule for SG-MCMC. The stepsize at iteration k is defined as (Loshchilov & Hutter, 2016; Huang et al., 2017):

$$\eta_k = \frac{\eta_0}{2} \cos \left(\frac{\text{mod}(k-1; dK=Me)}{dK=Me} \pi \right) + 1; \tag{1}$$

where η_0 is the initial stepsize, M is the number of cycles and K is the number of total iterations.

The stepsize η_k varies periodically with k . In each period, η_k starts at η_0 , and gradually decreases to 0. Within one period, SG-MCMC starts with a large stepsize, resulting in aggressive exploration in the parameter space; as the stepsize is decreasing, SG-MCMC explores local regions. In the next period, the Markov chain restarts with a large stepsize, enforcing the sampler to escape from the current mode and explore a new area of the posterior.

Related work in optimization. In optimization, the cyclical cosine annealing stepsize schedule has been demonstrated to be able to find diverse solutions in multimodal objectives, though not specifically different modes, using stochastic gradient methods (Loshchilov & Hutter, 2016; Huang et al., 2017; Garipov et al., 2018). Alternatively, we adopt the technique to SG-MCMC as an effective scheme for sampling from multimodal distributions.

3.1 EXPLORATION

The first stage of cyclical SG-MCMC, *exploration*, discovers parameters near local modes of an objective function. Unfortunately, it is undesirable to directly apply the cyclical schedule in optimization to SG-MCMC for collecting samples at every step. SG-MCMC often requires a small stepsize in order to control the error induced by the noise from using a minibatch approximation. If the stepsize is too large, the stationary distribution of SG-MCMC might be far away from the true posterior distribution. To correct this error, it is possible to do stochastic Metropolis-Hastings (MH) (Korattikara et al., 2014; Bardenet et al., 2014; Chen et al., 2016b). However, stochastic MH correction is still computationally too expensive. Further, it is easy to get rejected with an aggressive large stepsize, and every rejection is a waste of gradient computation.

To alleviate this problem, we propose to introduce a system temperature T to control the sampler’s behavior: $p(jD) \propto \exp(-U(jD)/T)$. Note that setting $T = 1$ corresponds to sampling from the true posterior. When $T \rightarrow 0$, the posterior distribution becomes a point mass. Sampling from $\lim_{T \rightarrow 0} \exp(-U(jD)/T)$ is equivalent to minimizing $U(jD)$; in this context, SG-MCMC methods become stochastic gradient optimization methods.

One may increase the temperature T from 0 to 1 when the step-size is decreasing. We simply consider $T = 0$ and perform optimization as the burn-in stage, when the completed proportion of a cycle $r(k) = \frac{\text{mod}(k-1; dK=Me)}{dK=Me}$ is smaller than a given threshold: $r(k) < \tau$. Note that $\tau \in (0; 1)$ balances the proportion of the exploration and sampling stages in cSG-MCMC.

3.2 SAMPLING

The *sampling* stage corresponds to $T = 1$ of the exploration stage. When $r(k) > \tau$ or step-sizes are sufficiently small, we initiate SG-MCMC updates and collect samples until this cycle ends.

SG-MCMC with Warm Restarts One may consider the exploration stage as automatically providing warm restarts for the sampling stage. Exploration alleviates the inefficient mixing and inability to traverse the multimodal distributions of the traditional SG-MCMC methods. SG-MCMC with warm restarts explores different parts of the posterior distribution and capture multiple modes in a single training procedure.

In summary, the proposed cyclical SG-MCMC repeats the *two* stages, with three key advantages: (i) It restarts with a large stepsize at the beginning of a cycle which provides enough perturbation and encourages the model to escape from the current mode. (ii) The stepsize decreases more quickly inside one cycle than a traditional schedule, making the sampler better characterize the density of the local regions. (iii) This cyclical stepsize shares the advantage of the “super-convergence” property discussed in (Smith & Topin, 2017): cSG-MCMC can accelerate convergence for DNNs by up to an order of magnitude.

Connection to the Santa algorithm. It is interesting to note that our approach inverts steps of the Santa algorithm (Chen et al., 2016a) for optimization. Santa is a simulated-annealing-based optimization algorithm with an exploration stage when $T = 1$, then gradually anneals $T \rightarrow 0$ in a refinement stage for global optimization. In contrast, our goal is to draw samples for multimodal distributions, thus we explore with $T = 0$ and sample with $T = 1$. Another fundamental difference is that Santa adopts the traditional stepsize decay, while we use the cyclical schedule.

We visually compare the difference between cyclical and traditional step size schedules (described in Section 2) in Figure 1. The cyclical SG-MCMC algorithm is presented in Algorithm 1.

Connection to Parallel MCMC. Running parallel Markov chains is a natural and effective way to draw samples from multimodal distributions (VanDerwerken & Schmidler, 2013; Ahn et al., 2014).

Algorithm 1 Cyclical SG-MCMC.

Input: The initial stepsize σ_0 , number of cycles M , number of training iterations K and the proportion of exploration stage τ .

for $k = 1:K$ **do**

σ_k according to Eq equation 1.

if $\frac{\text{mod}(k-1; dK=Me)}{dK=Me} < \tau$ **then**

% Exploration stage

$r \sim \tilde{U}_k(\sigma_k)$

else

% Sampling stage

Collect samples using SG-MCMC methods

Output: Samples $\{f_k^g\}$

However, the training cost increases linearly with the number of chains. Cyclical SG-MCMC can be seen as an efficient way to approximate parallel MCMC. Each cycle effectively estimates a different region of posterior. Note cyclical SG-MCMC runs along a single training pass. Therefore, its computational cost is the same as single chain SG-MCMC while significantly less than parallel MCMC.

Combining Samples. In cyclical SG-MCMC, we obtain samples from multiple modes of a posterior distribution by running the cyclical step size schedule for many periods. We provide a sampling combination scheme to effectively utilize the collected samples in Section F in Appendix.

4 THEORETICAL ANALYSIS

Our algorithm is based on the SDE characterizing the Langevin dynamics: $dx_t = -\nabla U(x_t)dt + \sqrt{2\beta^{-1}}dW_t$, where $W_t \in \mathbb{R}^d$ is a d -dimensional Brownian motion. In this section, we prove non-asymptotic convergence rates for the proposed cSG-MCMC framework with a cyclical stepsize sequence β_k defined in equation 1. For simplicity, we do not consider the exploration stage in the analysis as that corresponds to stochastic optimization. Generally, there are two different ways to describe the convergence behaviours of SG-MCMC. One characterizes the sample average over a particular test function (e.g., (Chen et al., 2015; Vollmer et al., 2016)); the other is in terms of Wasserstein distance (e.g., (Raginsky et al., 2017; Xu et al., 2017)). We study both in the following.

Weak convergence Following (Chen et al., 2015; Vollmer et al., 2016), we define the posterior average of an ergodic SDE as: $\mathbb{E}_\pi \int_{\mathcal{X}} \phi(x) dx$ for some test function $\phi(\cdot)$ of interest. For the corresponding algorithm with generated samples $(x_k)_{k=1}^K$, we use the *sample average* $\hat{\mu}$ defined as $\hat{\mu} = \frac{1}{K} \sum_{k=1}^K \phi(x_k)$ to approximate $\mathbb{E}_\pi \int_{\mathcal{X}} \phi(x) dx$. We prove weak convergence of cSGLD in terms of bias and MSE, as stated in Theorem 1.

Theorem 1. *Under Assumptions 2 in the Appendix, for a smooth test function ϕ , the bias and MSE of cSGLD are bounded as:*

$$\text{BIAS: } \mathbb{E}[\hat{\mu} - \mu] = O\left(\frac{1}{\sqrt{K}} + \epsilon\right); \quad \text{MSE: } \mathbb{E}[\hat{\mu} - \mu]^2 = O\left(\frac{1}{\sqrt{K}} + \epsilon^2\right); \quad (2)$$

Convergence under Wasserstein distance Next, we consider the more general case of SGLD and characterize convergence rates in terms of a stronger metric of 2-Wasserstein distance, defined as:

$$W_2^2(\mu; \nu) := \inf_{\Gamma(\mu; \nu)} \int_{\mathcal{X} \times \mathcal{X}} \|x - y\|_2^2 d\Gamma(x; y); \quad \Gamma(\mu; \nu)$$

where $\Gamma(\mu; \nu)$ is the set of joint distributions over $(x; y)$ such that the two marginals equal μ and ν , respectively.

Denote the distribution of x_t in the SDE as μ_t . According to (Chiang & Hwang, 1987), the stationary distribution μ_∞ matches our target distribution. Let μ_K be the distribution of the sample from our proposed cSGLD algorithm at the K -th iteration. Our goal is to derive a convergence bound on $W_2(\mu_K; \mu_\infty)$. We adopt standard assumptions as in most existing work, which are detailed in Assumption 3 in Appendix. Theorem 2 summarizes our main theoretical result.

Theorem 2. *Under Assumption 3 in Appendix, there exist constants $(C_0; C_1; C_2; C_3)$ independent of the stepsizes such that the convergence rate of our proposed cSGLD with cyclical stepsize sequence equation 1 is bounded for all K satisfying $(K \bmod M = 0)$, as $W_2(\mu_K; \mu_\infty)$*

$$C_3 \exp\left(\frac{K - \epsilon_0}{2C_4}\right) + 6 + \frac{C_2 K - \epsilon_0}{2} \left[\left((C_1 \frac{3 - \epsilon_0 K}{8} + C_0 \frac{K - \epsilon_0}{2})^{\frac{1}{2}} + (C_1 \frac{3 - \epsilon_0 K}{16} + C_0 \frac{K - \epsilon_0}{4})^{\frac{1}{4}} \right) \right];$$

Particularly, if we further assume $\epsilon_0 = O(K^{-\delta})$ for $\delta > 1$, $W_2(\mu_K; \mu_\infty) \leq C_3 + 6 + \frac{C_2}{K^{1-\delta}} \left[\left(\frac{2C_1}{K^{2-\delta}} + \frac{2C_0}{K^{1-\delta}} \right)^{\frac{1}{2}} + \left(\frac{C_1}{K^{2-\delta}} + \frac{C_0}{K^{1-\delta}} \right)^{\frac{1}{4}} \right]$.

Remark 1. *i) The bound is decomposed into two parts: the first part measures convergence speed of exact solution to the stationary distribution, i.e., $\mathbb{P}_{x \sim \mu_K}$ to μ_∞ ; the second part measures the numerical error, i.e., between μ_K and $\mathbb{P}_{x \sim \mu_K}$. ii) The overall bound offers a same order of*

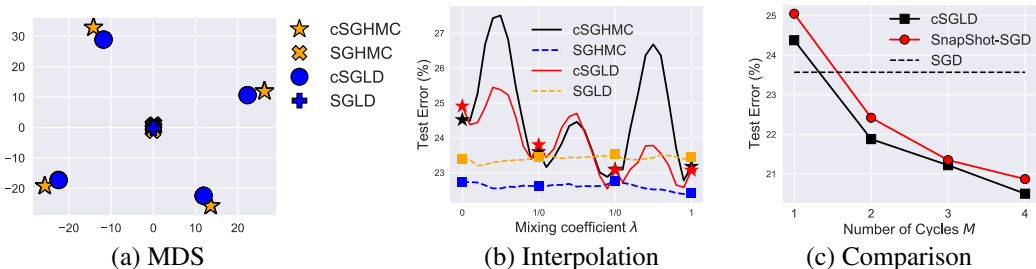


Figure 3: Results of cSG-MCMC with DNNs on CIFAR-100 dataset. (a) MDS visualization in weight space: cSG-MCMC show larger distance than traditional schedules. (b) Testing errors (%) on the path of two samples: cSG-MCMC shows more varied performance. (c) Testing errors (%) as a function of the number of cycles M : cSGLD yields consistently lower errors.

dependency on K as in standard SGLD (please see the bound for SGLD in Section E of Appendix. See also (Raginsky et al., 2017)). iii) If one imposes stricter assumptions such as in the convex case, the bound can be further improved. Specific bounds are derived in the Appendix. We did not consider this case due to the discrepancy from real applications.

5 EXPERIMENTS

We demonstrate cSG-MCMC on several tasks, including a synthetic multimodal distribution (Section 5.1), image classification on Bayesian neural networks (Section 5.2) and uncertainty estimation in Section 5.3. We also demonstrate cSG-MCMC can improve the estimate efficiency for uni-modal distributions using Bayesian logistic regression in Section A.2 in Appendix. We choose SGLD and SGHMC as the representative baseline algorithms. Their cyclical counterpart are called cSGLD and cSGHMC, respectively.

5.1 SYNTHETIC MULTIMODAL DATA

We first demonstrate the ability of cSG-MCMC for sampling from a multi-modal distribution on a 2D mixture of 25 Gaussians. Specifically, we compare cSGLD with SGLD in two setting: (1) parallel running with 4 chains and (2) running with a single chain, respectively. Each chain runs for 50k iterations. The step-size schedule of SGLD is $\frac{\epsilon}{\sqrt{k}}$. In cSGLD, we set $M = 30$ and the initial step-size $\epsilon_0 = 0.09$. The proportion of exploration stage $\alpha = \frac{1}{4}$. Fig 2 shows the estimated density using sampling results for SGLD and cSGLD in the parallel setting. We observed that SGLD gets trapped in the local modes, depending on the initial position. In any practical time period, SGLD could only characterize partial distribution. In contrast, cSGLD is able to find and characterize all modes, regardless of the initial position. cSGLD leverages large step sizes to discover a new mode, and small step sizes to explore local modes. This result suggests cSGLD can be a significantly favourable choice in the non-asymptotic setting, for example only 50k iterations in this case. The single chain results and the quantitative results on mode coverage are reported in Section A.1 of Appendix.

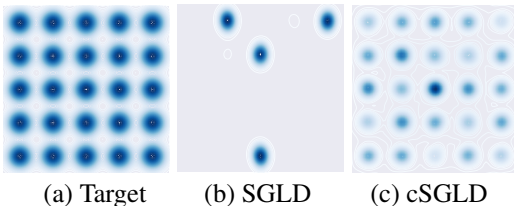


Figure 2: Sampling from a mixture of 25 Gaussians shown in (a) for the parallel setting. With a budget of $50k \times 4 = 200k$ samples, traditional SGLD in (b) has only discovered 4 of the 25 modes, while our cSGLD in (c) has fully explored the distribution.

5.2 BAYESIAN NEURAL NETWORKS

We demonstrate the effectiveness of cSG-MCMC on Bayesian neural networks for classification on CIFAR-10 and CIFAR-100. We compare with (i) traditional SG-MCMC; (ii) traditional stochastic optimization methods, including stochastic gradient descent (SGD) and stochastic gradient descent with momentum (SGDM); and (iii) *Snapshot*: a stochastic optimization ensemble method method with a the cyclical stepsize schedule (Huang et al., 2017). We use a ResNet-18 (He et al., 2016) and

run all algorithms for 200 epochs. We report the test errors averaged over 3 runs, and the standard error () from the mean predictor.

We set $M = 4$ and $\beta_0 = 0.5$ for cSGLD, cSGHMC and Snapshot. The proportion hyper-parameter $\alpha = 0.8$ and 0.94 for CIFAR-10 and CIFAR-100, respectively. We collect 3 samples per cycle. In practice, we found that the collected samples share similarly high likelihood for DNNs, thus one may simply set the normalizing term γ in Eq. equation 33 to be the same for faster testing.

For the traditional SG-MCMC methods, we found that noise injection early in training hurts convergence. To make these baselines as competitive as possible, we thus avoid noise injection for the first 150 epochs of training (corresponding to the zero temperature limit of SGLD and SGHMC), and resume SGMCMC as usual (with noise) for the last 50 epochs. This scheme is similar to the exploration and sampling stages within one cycle of cSG-MCMC. We collect 20 samples for the MCMC methods and average their predictions in testing.

Testing Performance for Image Classification

We report the testing errors in Table 1 to compare with the non-parallel algorithms. Snapshot and traditional SG-MCMC reduce the testing errors on both datasets. Performance variance for these methods is also relatively small, due to the multiple networks in the Bayesian model average. Further, cSG-MCMC significantly outperforms Snapshot ensembles and the traditional SG-MCMC, demonstrating the importance of (1) capturing diverse modes compared to traditional SG-MCMC, and (2) capturing fine-scale characteristics of the distribution compared with Snapshot ensembles.

	CIFAR-10		CIFAR-100	
SGD	5.29	0.15	23.61	0.09
SGDM	5.17	0.09	22.98	0.27
Snapshot-SGD	4.46	0.04	20.83	0.01
Snapshot-SGDM	4.39	0.01	20.81	0.10
SGLD	5.20	0.06	23.23	0.01
cSGLD	4.29	0.06	20.55	0.06
SGHMC	4.93	0.1	22.60	0.17
cSGHMC	4.27	0.03	20.50	0.11

Table 1: Comparison of test error (%) between cSG-MCMC with non-parallel algorithms. cSGLD and cSGHMC yields lower errors than their optimization counterparts, respectively.

Diversity in Weight Space. To further demonstrate our hypothesis that with a limited budget cSG-MCMC can find diverse modes, while traditional SG-MCMC cannot, we visualize the 12 samples we collect from cSG-MCMC and SG-MCMC on CIFAR-100 respectively using Multidimensional Scaling (MDS) in Figure 3 (a). MDS uses a Euclidean distance metric between the weight of samples. We see that the samples of cSG-MCMC form 4 clusters, which means they are from 4 different modes in weight space. However, all samples from SG-MCMC only form one cluster, which indicates traditional SG-MCMC gets trapped in one mode and only samples from that mode.

Diversity in Prediction. To further demonstrate the samples from different cycles of cSG-MCMC provide diverse predictions we choose one sample from each cycle and linearly interpolate between two of them (Goodfellow et al., 2014; Huang et al., 2017). Specifically, let $J(\theta)$ be the test error of a sample with parameter θ . We compute the test error of the convex combination of two samples $J(\alpha\theta_1 + (1-\alpha)\theta_2)$, where $\alpha \in [0, 1]$.

We linearly interpolate between two samples from neighboring chains of cSG-MCMC since they are the most likely to be similar. We randomly select 4 samples from SG-MCMC. If the samples are from the same mode, the test error of the linear interpolation of parameters will be relatively smooth, while if the samples are from different modes, the test error of the parameter interpolation will have a spike when α is between 0 and 1.

We show the results of interpolation for cSG-MCMC and SG-MCMC on CIFAR-100 in Figure 3 (b). We see a spike in the test error in each linear interpolation of parameters between two samples from neighboring chains in cSG-MCMC while the linear interpolation for samples of SG-MCMC is smooth. This result suggests that samples of cSG-MCMC from different chains are from different modes while samples of SG-MCMC are from the same mode.

Although the test error of a single sample of cSG-MCMC is worse than that of SG-MCMC shown in Figure 3 (c), the ensemble of these samples significantly improves the test error, indicating that samples from different modes provide different predictions and make mistakes on different data points. Thus these diverse samples can complement each other, resulting in a lower test error, and demonstrating the advantage of exploring diverse modes using cSG-MCMC.

Method	Cyclical+Parallel		Decreasing+Parallel		Decreasing+Parallel		Cyclical+Single	
Cost	200/800		200/800		100/400		200/200	
Sampler	SGLD	SGHMC	SGLD	SGHMC	SGLD	SGHMC	SGLD	SGHMC
CIFAR-10	4.09	3.95	4.15	4.09	5.11	4.52	4.29	4.27
CIFAR-100	19.37	19.19	20.29	19.72	21.16	20.82	20.55	20.50

Table 2: Comparison of test error (%) between cSG-MCMC with parallel algorithm (4 chains) on CIFAR-10 and CIFAR-100. The method is reported in the format of “step-size schedule (cyclical or decreasing) + single/parallel chain”. The cost is reported in the format of “# epoch per chain / # epoch used in all chains”. Note that a parallel algorithm with a single chain reduces to a non-parallel algorithm. Integration of the cyclical schedule with parallel algorithms provides lower testing errors.

Comparison to Parallel MCMC. cSG-MCMC can be viewed as an economical alternative to parallel MCMC. We verify how closely cSG-MCMC can approximate the performance of parallel MCMC, but with more convenience and less computational expense. We also note that we can improve parallel MCMC with the proposed cyclical stepsize schedule.

We report the testing errors in Table 2 to compare multiple-chain results. (1) Four chains used, each runs 200 epochs (800 epochs in total), the results are shown in the first 4 columns (Cyclical+Parallel vs Decreasing+Parallel). We see that cSG-MCMC variants provide lower errors than plain SG-MCMC. (2) We reduce the number of epochs (epoch) of parallel MCMC to 100 epoch each for decreasing stepsize schedule. The total cost is 400 epochs. We compare its performance with cyclical single chain (200 epochs in total) in the last 4 columns (Decreasing+Parallel vs Cyclical+Single). We see that the cyclical schedule running on a single chain performs best even with half the computational cost! All the results indicate the importance of warm re-starts using the proposed cyclical schedule. For a given total cost budget, the proposed cSGMCMC is preferable to parallel sampling.

Comparison to Snapshot Optimization. We carefully compared with Snapshot, as our cSG-MCMC can be viewed as the sampling counterpart of the Snapshot optimization method. We plot the test error wrt. various number of cycles in Fig. 3. As M increases, cSG-MCMC and Snapshot both improve. However, given a x_{MCMC} , cSG-MCMC yields substantially lower test errors than Snapshot. This result is due to the ability of cSG-MCMC to better characterize the local distribution of modes: Snapshot provides a single minimum per cycle, while cSG-MCMC fully exploits the mode with more samples, which could provide weight uncertainty estimate and avoid over-fitting.

Results on ImageNet. We further study different learning algorithms on a large-scale dataset, ImageNet. ResNet-50 is used as the architecture, and 120 epochs for each run. The results on the testing set are summarized in Table 3, including NLL, Top1 and Top5 accuracy

	NLL #	Top1 %	Top5 %
SGDM	0.9595	76.046	92.776
Snapshot-SGDM	0.8941	77.142	93.344
SGHMC	0.9308	76.274	92.994
cSGHMC	0.8882	77.114	93.524

(%), respectively. 3 cycles are considered for both cSGHMC and Snapshot, and we collect 2 samples per cycle. We see that cSGHMC yields the lowest testing NLL, indicating that the cyclical schedule is an effective technique to explore the parameter space, and diversified samples can help prevent over-fitting.

5.3 UNCERTAINTY EVALUATION

To demonstrate how predictive uncertainty benefits from exploring multiple modes in the posterior of neural network weights, we consider the task of uncertainty estimation for out-of-distribution samples (Lakshminarayanan et al., 2017). We train a three-layer MLP model on the standard MNIST train dataset until convergence using different algorithms, and estimate the entropy of the predictive distribution on the notMNIST dataset (Bulatov, 2011). Since the samples from the notMNIST dataset belong to the unseen classes, ideally the predictive distribution of the trained model should be uniform over the notMNIST digits, which gives the maximum entropy.

In Figure 4, we plot the empirical CDF for the entropy of the predictive distributions on notMNIST. We see that the uncertainty estimates from cSGHMC and cSGLD are better than the other methods,

since the probability of a low entropy prediction is overall lower. cSG-MCMC algorithms explore more modes in the weight space, each mode characterizes a meaningfully different representation of MNIST data. When testing on the out-of-distribution dataset (notMNIST), each mode can provide different predictions over the label space, leading to more reasonable uncertainty estimates. Snapshot achieves less entropy than cSG-MCMC, since it represents each mode with a single point.

The traditional SG-MCMC methods also provide better uncertainty estimation compared to their optimization counterparts, because they characterize a local region of the parameter space, rather than a single point. cSG-MCMC can be regarded as a combination of these two worlds: a wide coverage of many modes in Snapshot, and ne-scale characterization of local regions in SG-MCMC.

6 DISCUSSION

We have proposed cyclical SG-MCMC methods to automatically explore complex multimodal distributions. Our approach is particularly compelling for Bayesian deep learning, which involves rich multimodal parameter posteriors corresponding to meaningfully different representations. We have also shown that our cyclical methods explore unimodal distributions more efficiently. These results are in accordance with theory we developed to show that cyclical SG-MCMC will converge faster to samples from a stationary distribution in general settings. Moreover, we show cyclical SG-MCMC methods provide more accurate uncertainty estimation, by capturing more diversity in the hypothesis space corresponding to settings of model parameters.

While MCMC was once the gold standard for inference with neural networks, it is now rarely used in modern deep learning. We hope that this paper will help renew interest in MCMC for this purpose. Indeed, MCMC is uniquely positioned to explore the rich multimodal posterior distributions of modern neural networks, which can lead to improved accuracy, reliability, and uncertainty representation.

REFERENCES

- Sungjin Ahn, Babak Shahbaba, and Max Welling. Distributed stochastic gradient MCMC. *ICML*, 2014.
- Peter Auer, Mark Herbster, and Manfred K Warmuth. Exponentially many local minima for single neurons. *IrNIPS* 1996.
- Rémi Bardenet, Arnaud Doucet, and Chris Holmes. Towards scaling up Markov Chain Monte Carlo: an adaptive subsampling approach. *ICML*, 2014.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *ICML*, 2015.
- François Bolley and Edric Villani. Weighted chi-squared-kullback-pinsker inequalities and applications to transportation inequalities. *Annales de la Faculté des sciences de Toulouse Université Paul Sabatier*, 2005.
- Yaroslav Bulatov. Not MNIST Dataset. 2011. <http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html>
- Changyou Chen, Nan Ding, and Lawrence Carin. On the convergence of stochastic gradient MCMC algorithms with high-order integrators. *NIPS* 2015.

- Changyou Chen, David Carlson, Zhe Gan, Chunyuan Li, and Lawrence Carin. Bridging the gap between stochastic gradient MCMC and stochastic optimization. *Artificial Intelligence and Statistics* 2016a.
- Haoyu Chen, Daniel Seita, Xinlei Pan, and John Canny. An efficient minibatch acceptance test for Metropolis-Hastings. *arXiv preprint arXiv:1610.06848*, 2016b.
- Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In *ICML*, 2014.
- Tzoo-Shuh Chiang and Chii-Ruey Hwang. Diffusion for global optimization in *SIAM J. Control Optim*, pp. 737–753, 1987. ISSN 0363-0129. doi: 10.1137/0325042. <http://dx.doi.org/10.1137/0325042>
- Anna Choromanska, Mikael Henaff, Michael Mathieu, Jérôme B. Arous, and Yann LeCun. The loss surfaces of multilayer networks. *Artificial Intelligence and Statistics*, 2015.
- Arnak S. Dalalyan and Avetik Karagulyan. User-friendly guarantees for the langevin monte carlo with inaccurate gradient. *Stochastic Processes and their Applications* 2019. ISSN 0304-4149. doi: <https://doi.org/10.1016/j.spa.2019.02.016>. <http://www.sciencedirect.com/science/article/pii/S0304414918304824>
- Nan Ding, Youhan Fang, Ryan Babbush, Changyou Chen, Robert D Skeel, and Hartmut Neven. Bayesian sampling using stochastic gradient thermostats. *NIPS*, 2014.
- Zhe Gan, Chunyuan Li, Changyou Chen, Yunchen Pu, Qinliang Su, and Lawrence Carin. Scalable Bayesian learning of recurrent neural networks for language modeling. *ICML*, 2016.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of DNNs. *Advances in Neural Information Processing Systems*, pp. 8789–8798, 2018.
- Ian J Goodfellow, Oriol Vinyals, and Andrew M Saxe. Qualitatively characterizing neural network optimization problems. *arXiv preprint arXiv:1412.6544*, 2014.
- Peter J Green. Reversible jump MCMC computation and bayesian model determination. *Biometrika* 82(4):711–732, 1995.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of Bayesian neural networks. *ICML*, 2015.
- Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E Hopcroft, and Kilian Q Weinberger. Snapshot ensembles: Train 1, get m for free. *ICLR*, 2017.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Anoop Korattikara, Yutian Chen, and Max Welling. Austerity in MCMC land: Cutting the Metropolis-Hastings budget. *ICML*, 2014.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *NIPS*, 2017.
- Chunyuan Li, Changyou Chen, David E Carlson, and Lawrence Carin. Preconditioned stochastic gradient Langevin dynamics for deep neural networks. *AAAI*, 2016a.
- Chunyuan Li, Andrew Stevens, Changyou Chen, Yunchen Pu, Zhe Gan, and Lawrence Carin. Learning weight uncertainty with stochastic gradient MCMC for shape classification. *CVPR*, 2016b.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. 2016.

- Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient MCMC. In NIPS, 2015.
- David JC MacKay. A practical bayesian framework for backpropagation networks for neural computation, 1992.
- J. C. Mattingly, A. M. Stuart, and M. V. Tretyakov. Construction of numerical time-average and stationary measures via Poisson equations. SIAM J. NUMER. ANAL., 48(2):552–577, 2010.
- Radford M Neal. Bayesian learning for neural networks. New York: Springer-Verlag, 1996.
- Arvind Neelakantan, Luke Vilnis, Quoc V Le, Ilya Sutskever, Lukasz Kaiser, Karol Kurach, and James Martens. Adding gradient noise improves learning for very deep networks. ICLR workshop 2016.
- Adrian E Raftery, Michael A Newton, Jaya M Satagopan, and Pavel N Krivitsky. Estimating the integrated likelihood via posterior simulation using the harmonic mean identity. 2006.
- Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient Langevin dynamics: a nonasymptotic analysis. arXiv preprint arXiv:1702.03849, 2017.
- Yunus Saatchi and Andrew Gordon Wilson. Bayesian GANs. NIPS, 2017.
- Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of residual networks using large learning rates. arXiv preprint arXiv:1708.07120, 2017.
- Yee Whye Teh, Alexandre H Thiery, and Sebastian J Vollmer. Consistency and fluctuations for stochastic gradient Langevin dynamics. The Journal of Machine Learning Research, 2016.
- Douglas N VanDerwerken and Scott C Schmidler. Parallel Markov Chain Monte Carlo. arXiv preprint arXiv:1312.7479, 2013.
- S. J. Vollmer, K. C. Zygalakis, and Y. W. Teh. (Non-)asymptotic properties of stochastic gradient Langevin dynamics. Technical Report arXiv:1501.00438, University of Oxford, UK, January 2015. URL <http://arxiv.org/abs/1501.00438>.
- Sebastian J. Vollmer, Konstantinos C. Zygalakis, and Yee Whye Teh. Exploration of the (non-)asymptotic bias and variance of stochastic gradient langevin dynamics. Journal of Machine Learning Research, 17(159):1–48, 2016.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient Langevin dynamics. In ICML, 2011.
- Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. Global convergence of Langevin dynamics based algorithms for nonconvex optimization. arXiv preprint arXiv:1707.06618, 2017.
- Jianyi Zhang, Ruiyi Zhang, and Changyou Chen. Stochastic Particle-Optimization Sampling and the Non-Asymptotic Convergence Theorem. arXiv e-prints art. arXiv:1809.01293, Sep 2018.
- Y. Zhang, P. Liang, and M. Charikar. A hitting time analysis of stochastic gradient Langevin dynamics. In COLT, 2017.

Supplementary Material: Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning

A EXPERIMENTAL RESULTS

A.1 SYNTHETIC MULTIMODAL DISTRIBUTION

The density of the distribution is

$$F(x) = \prod_{i=1}^5 N(x; \mu_i, \Sigma_i);$$

where $\mu = \frac{1}{25}$, $\Sigma = \text{diag}(4, 2, 0, 2, 4)$, $\mu_i = \mu + \frac{h}{0.03} \cdot \frac{i}{0.03}$.

In Figure 5, we show the estimated density for SGLD and cSGLD in the non-parallel setting.

(a) Target (b) SGLD (c) cSGLD

Figure 5: Sampling from a mixture of 25 Gaussians in the non-parallel setting. With a budget of 50K samples, traditional SGLD has only discovered one of the 25 modes, while our proposed cSGLD has explored significantly more of the distribution.

To quantitatively show the ability of different algorithms to explore multi-modal distributions, we define the mode-coverage metric: when the number of samples falling within the radius of a mode center is larger than a threshold, we consider this mode covered. On this dataset, we choose $r = 0.25$ and $n = 100$. Table 4 shows the mode-coverage for several algorithms, based on 10 different runs.

Algorithm	Mode coverage
SGLD	1.8 0.13
cSGLD	6.7 0.52
Parallel SGLD	18 0.47
Parallel cSGLD	24.4 0.22

Table 4: Mode coverage over 10 different runs, standard error.

A.2 BAYESIAN LOGISTIC REGRESSION

We consider Bayesian logistic regression (BLR) on three real-world datasets from the UCI repository: Australian (15 covariates, 690 data points), German (25 covariates, 1000 data points) and Heart (14 covariates, 270 data points). For all experiments, we collect 5000 samples with 1000 burn-in iterations. Following the settings in (Li et al., 2016a), we report mean effective sample size (ESS) in Table 5.

Note that BLR is unimodal in parameter space. We use this experiment as an adversarial situation for cSG-MCMC, which we primarily designed to explore multiple modes. We note that even in the unimodal setting, cSG-MCMC more effectively explores the parameter space than popular alternatives. We can also use these experiments to understand how samplers respond to varying parameter dimensionality and training set sizes.

Overall, cSG-MCMC dramatically outperforms SG-MCMC, which demonstrates the fast mixing rate due to the warm restarts. On the small dataset heart, SGHMC and cSGHMC achieve the same results, because the posterior of BLR on this dataset is simple. However, in higher dimensional spaces (e.g., Australian and German), cSG-MCMC shows significantly higher ESS; this result means that each cycle in cSG-MCMC can characterize a different region of the posteriors, combining multiple cycles yields more accurate overall approximation.

	Australian	German	Heart
SGLD	1676	492	2199
cSGLD	2138	978	2541
SGHMC	1317	2007	5000
cSGHMC	4707	2436	5000

Table 5: Effective sample size for samples for the unimodal posteriors in Bayesian linear regression, obtained using cyclical and traditional SG-MCMC algorithms, respectively.

B ASSUMPTIONS

B.1 ASSUMPTIONS IN WEAK CONVERGENCE ANALYSIS

In the analysis, we define a functional that solves the following Poisson Equation

$$L(\varphi_k) = -\varphi_k; \quad \text{or equivalently, } \frac{1}{K} \sum_{k=1}^K L(\varphi_k) = -\varphi_k. \quad (3)$$

The solution functional φ_k characterizes the difference between φ_k and the posterior average for every k , thus would typically possess a unique solution, which is at least as smooth as the elliptic or hypoelliptic settings (Mattingly et al., 2010). Following (Chen et al., 2015; Vollmer et al., 2016), we make certain assumptions on the solution functional of the Poisson equation equation 3.

Assumption 2. φ_k and its up to 3rd-order derivatives $\mathcal{D}^k \varphi_k$, are bounded by a function φ , i.e., $\|\mathcal{D}^k \varphi_k\| \leq H_k V^p$ for $k = (0; 1; 2; 3)$, $H_k; p_k > 0$. Furthermore, the expectation of φ on $f_k g$ is bounded: $\sup_{\varphi} \mathbb{E} V^p(\varphi_k) < 1$, and V is smooth such that $\sup_{s \in \mathcal{S}_2(0;1)} V^p(s + (1-s)\varphi) \leq C(V^p(\varphi) + V^p(\varphi_0))$, $8; \varphi_0; p \max_{\varphi} 2p_k g$ for some $C > 0$.

B.2 ASSUMPTIONS IN CONVERGENCE UNDER WASSERSTEIN DISTANCE

Following existing work in (Raginsky et al., 2017), we adopt the following standard assumptions summarized in Assumption 3.

Assumption 3. There exists some constants $A \geq 0$ and $B \geq 0$, such that $\mathbb{E} U(0) \leq A$ and $\mathbb{E} U(0) \leq B$.

The function U is L_U -smooth: $\|U(w) - U(v)\| \leq L_U \|w - v\|$.

The function U is $(m_U; b)$ -dissipative, which means for some $m_U > 0$ and $b > 0$ $\|U(w) - U(v)\| \leq m_U \|w - v\|^2 + b$.

There exists some constant $\beta \in [0; 1)$, such that $\mathbb{E}[\|U_k(w) - U(w)\|^2] \leq 2(M_U^2 \|w\|^2 + B^2)$.

We can choose β_0 which satisfies the requirement: $\beta_0 := \log \int_{\mathcal{W}} e^{\beta_0 \|w\|^2} \varphi_0(w) dw < 1$.

C PROOF OF THEOREM 1

To prove the theorem, we borrow tools developed by (Chen et al., 2015; Vollmer et al., 2015). We first rephrase the stepsize assumptions in general SG-MCMC in Assumption 4.

Assumption 4. The algorithm adopts an n -th order integrator. The step sizes h_k are such that $0 < h_{k+1} < h_k$, and satisfy 1) $\prod_{k=1}^K h_k = 1$; and 2) $\lim_{K \rightarrow \infty} \frac{\prod_{k=1}^K h_k^{N+1}}{S_K^2} = 0$.

Our prove can be derived by the following results from (Chen et al., 2015).

Lemma 1 ((Chen et al., 2015)) Let $S_K = \prod_{k=1}^K h_k$. Under Assumptions 2 and 4, for a smooth test function ϕ , the bias and MSE of a decreasing-step-size SG-MCMC with n -th order integrator at time S_K are bounded as:

$$\text{BIAS: } E \sim \phi = O \left(\frac{1}{S_K} + \frac{\prod_{k=1}^K h_k^{N+1}}{S_K} \right) \quad (4)$$

$$\text{MSE: } E \sim \phi^2 = O \left(\sum_{l=1}^n \frac{h_k^2}{S_K^2} E_k |V_l|^2 + \frac{1}{S_K} + \frac{(\prod_{k=1}^K h_k^{N+1})^2}{S_K^2} \right) \quad (5)$$

Note that Assumption 4 is only required if one wants to prove the asymptotically unbiased of an algorithm. Lemma 1 still applies even if Assumption 4 is not satisfied. In this case one would obtain a biased algorithm, which is the case of cSGLD.

Proof of Theorem 1 Our results is actually a special case of Lemma 1. To see that, first note that our cSGLD adopts a first order integrator, thus $n=1$. To proceed, note that $S_K = \prod_{k=1}^K h_k = O(\frac{1}{\sqrt{K}})$, and

$$\begin{aligned} \sum_{j=0}^{K-1} \frac{1}{4} \left[\cos\left(\frac{\text{mod}(j+1; [K=M])}{[K=M]}\right) + 1 \right]^2 \\ = \sum_{j=0}^{K-1} \frac{1}{4} \left[\cos^2\left(\frac{\text{mod}(j+1; K=M)}{K=M}\right) + 1 \right]^2 \\ = \frac{2}{4} \frac{K}{M} \left(\frac{M}{2} + M \right) = \frac{3}{8} \frac{2K}{8} \end{aligned} \quad (6)$$

As a result, for the bias, we have

$$\begin{aligned} E \sim \phi &= O \left(\frac{1}{S_K} + \frac{\prod_{k=1}^K h_k^{N+1}}{S_K} \right) = O \left(\frac{1}{\sqrt{K}} + \frac{3}{8} \frac{2K}{8} \right) \\ &= O \left(\frac{1}{\sqrt{K}} \right) \end{aligned}$$

For the MSE, note the first term $\sum_{l=1}^n \frac{h_k^2}{S_K^2} E_k |V_l|^2$ has a higher order than other terms, thus it is omitted in the big-O notation, i.e.,

$$\begin{aligned} E \sim \phi^2 &= O \left(\frac{1}{\sqrt{K}} + \left(\frac{3}{8} \frac{2K}{8} \right)^2 \right) \\ &= O \left(\frac{1}{\sqrt{K}} \right) \end{aligned}$$

This completes the proof. □

D PROOF OF THEOREM 2

Proof of the bound for $W_2(\pi_K; \pi_1)$ in cSGLD. Firstly, we introduce the following SDE

$$d_t = -r U(t) dt + \sqrt{2} dW_t; \quad (7)$$

Let μ_t denote the distribution of μ_t , and the stationary distribution of equation (34), which means $\mu_1 = p(\cdot|D)$.

$$\mu_{k+1} = \mu_k \circ U_k(\mu_k) + \frac{P}{2} \mu_{k+1} \quad (8)$$

Further, let μ_k denote the distribution of μ_k .

Since

$$W_2(\mu_k; \mu_1) = W_2(\mu_k; P_{k=1}^{\mu_k}) + W_2(P_{k=1}^{\mu_k}; \mu_1) \quad (9)$$

, we need to give the bounds for these two parts respectively.

D.1 $W_2(\mu_k; P_{k=1}^{\mu_k})$

For the first part, $W_2(\mu_k; P_{k=1}^{\mu_k})$, our proof is based on the proof of Lemma 3.6 in (Raginsky et al., 2017) with some modifications. We first assume $\mu(w) = \mu \circ U(w); \mu \in \mathbb{R}^d$; which is a general assumption according to the way we choose the minibatch. And we note which will be used in the following proof:

$$p(t) = \prod_{i=1}^k Z_j \mu_{i-1}(t) < \prod_{i=1}^{k-1} \mu_i \quad (10)$$

Then we focus on the following continuous-time interpolation of

$$\mu_t = \int_0^t \mu_{k=1}^{\mu_t}(s) ds + \frac{P}{2} \int_0^t dW_s^{(d)} \quad (11)$$

where $\mu_{k=1}^{\mu_t}(s) = \mu_{k=1}^{\mu_t}(s)$ for $t \in [0, 1]$. And for each k , $\mu_{k=1}^{\mu_t}(s)$ and μ_k have the same probability law μ_k .

Since μ_t is not a Markov process, we define the following process which has the same one-time marginals as μ_t

$$V(t) = \int_0^t G_s(V(s)) ds + \frac{P}{2} \int_0^t dW_s^{(d)} \quad (12)$$

with

$$G_t(x) := E \left[\mu_{k=1}^{\mu_t}(x) \right] \quad (13)$$

Let $P_V^t := L(V(s); 0 \leq s \leq t)$ and $P^t := L(\mu(s); 0 \leq s \leq t)$ and according to the proof of Lemma 3.6 in (Raginsky et al., 2017), we can derive a similar result for the relative entropy of P_V^t and P^t :

$$\begin{aligned} D_{KL}(P_V^t \parallel P^t) &= \int_0^t dP_V^t \log \frac{dP_V^t}{dP^t} \\ &= \frac{1}{4} \int_0^t E_k \mu(U(V(s))) G_s(V(s)) k^2 ds \\ &= \frac{1}{4} \int_0^t E_k \mu(\mu(s)) G_s(\mu(s)) k^2 ds \end{aligned}$$

The last line follows the fact that $\mathbb{P}(\cdot|s) = L(V(s)); 8s$.

Then we will let $\mathbb{P} = \mathbb{P}_{k=1}^K$ and we can use the martingale property of the integral to derive:

$$\begin{aligned}
 & D_{KL}(\mathbb{P}_{V_{k=1}^K} \|\mathbb{P}_{k=1}^K) \\
 &= \frac{1}{4} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) G_s(\cdot|s) k^2 ds \\
 & \quad + \frac{1}{2} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) r U(\cdot|s) k^2 ds \\
 & \quad + \frac{1}{2} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) G_s(\cdot|s) k^2 ds \\
 & \quad + \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 ds \tag{14}
 \end{aligned}$$

$$\begin{aligned}
 & \quad + \frac{1}{2} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) G_s(\cdot|s) k^2 ds \tag{15}
 \end{aligned}$$

For the first part (14), we consider some $\mathbb{P}_{k=1}^j$ and $\mathbb{P}_{k=1}^{j+1}$, for which the following holds:

$$\begin{aligned}
 & \mathbb{P}(\cdot|s) - \mathbb{P}(\cdot|s) \\
 &= (s) \mathbb{P}_{k=1}^j(r U_k(\cdot|s) + \frac{1}{2} (W_s^{(d)} - W_{\mathbb{P}_{k=1}^j}^{(d)})) \\
 &= (s) \mathbb{P}_{k=1}^j(r U(\cdot|s)) + (s) \mathbb{P}_{k=1}^j(r U(\cdot|s) - r U_k(\cdot|s)) + \frac{1}{2} (W_s^{(d)} - W_{\mathbb{P}_{k=1}^j}^{(d)}) \tag{16}
 \end{aligned}$$

Thus, we can use Lemma 3.1 and 3.2 in (Raginsky et al., 2017) for the following result:

$$\begin{aligned}
 & \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 \leq 3 \sum_{j+1}^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 + 3 \sum_{j+1}^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) r U_j(\cdot|s) k^2 + 6 \sum_{j+1}^2 d \\
 & \quad + 12 \sum_{j+1}^2 (L_U^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 + B^2) + 6 \sum_{j+1}^2 d
 \end{aligned}$$

Hence we can bound the first part, (choosing $\epsilon = 1$),

$$\begin{aligned}
 & \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^1 \mathbb{E} \int_{k=1}^K \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 ds \\
 & \quad + \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^1 12 \sum_{j+1}^2 (L_U^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 + B^2) + 6 \sum_{j+1}^2 d \\
 & \quad + L_U^2 \max_{j \in \{K-1\}} 6(L_U^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 + B^2) + 3d \sum_{j=0}^{K-1} \sum_{j+1}^2 \\
 & \quad + L_U^2 \max_{j \in \{K-1\}} 6(L_U^2 \mathbb{E} \int_{k=1}^K U(\cdot|s) k^2 + B^2) + 3d \frac{3}{8} \sum_{j=0}^{K-1} \tag{17}
 \end{aligned}$$

The last line (17) follows from equation 6. The second part (15) can be bounded as follows:

$$\begin{aligned}
& \frac{1}{2} \sum_{j=0}^{K-1} \int_{k=1}^{j+1} P_{k=1}^{j+1} E_k U_{\left(\frac{\cdot}{i}\right)} G_s \left(\frac{\cdot}{i}\right) k^2 ds \\
&= \frac{1}{2} \sum_{j=0}^{K-1} E_k U_{\left(\frac{\cdot}{j}\right)} r_{\left(\frac{\cdot}{j}\right)} k^2 \\
& \leq \max_{j=0}^{K-1} (L_U^2 E_k j k^2 + B^2) \sum_{j=0}^{K-1} \\
& \leq \max_{j=0}^{K-1} (L_U^2 E_k j k^2 + B^2) \left(\frac{0}{2} \sum_{j=0}^{K-1} \left(\cos\left(\frac{\text{mod}(j; K=M)}{K=M}\right) + 1 \right) \right) \\
& \leq \max_{j=0}^{K-1} (L_U^2 E_k j k^2 + B^2) \left(\frac{K}{2} \right)
\end{aligned}$$

Due to the data-processing inequality for the relative entropy, we have

$$\begin{aligned}
D_{KL}(P_{k=1}^k | P_{k=1}^k) & \leq D_{KL}(P_V^t | P^t) \\
& \leq \frac{L_U^2}{2} \sum_{j=0}^{K-1} \int_{k=1}^{j+1} P_{k=1}^{j+1} E_k(s) \left(\frac{\cdot}{i}\right) k^2 ds \\
& \quad + \frac{1}{2} \sum_{j=0}^{K-1} \int_{k=1}^{j+1} P_{k=1}^{j+1} E_k U_{\left(\frac{\cdot}{i}\right)} G_s \left(\frac{\cdot}{i}\right) k^2 ds \\
& \leq L_U^2 \max_{j=0}^{K-1} 6(L_U^2 E_k j k^2 + B^2) + 3d \frac{3}{8} K \\
& \quad + \max_{j=0}^{K-1} (L_U^2 E_k j k^2 + B^2) \left(\frac{K}{2} \right)
\end{aligned}$$

According to the proof of Lemma 3.2 in (Raginsky et al., 2017), we can bound the term

$$E_k k^2 \leq (1 - 2 \min m_U + 4 \frac{2}{\min} M_U^2) E_k k^2 + 2 \min b + 4 \frac{2}{\min} B^2 + \frac{2 \min d}{4M_U^2}$$

Similar to the statement of Lemma 3.2 in (Raginsky et al., 2017), we can bound $\sum_{k=1}^j (0; 1 \wedge \frac{m_U}{4M_U^2})$. Then, we can know that

$$E_k k^2 \leq (1 - 2 \min m_U + 4 \frac{2}{\min} M_U^2) E_k k^2 + 2 \min b + 4 \frac{2}{\min} B^2 + \frac{2 \min d}{4M_U^2} \quad (18)$$

, where \min is defined as $\min \left(\frac{0}{2} \cos \frac{\text{mod}(dK=M e^{-1}; dK=M e)}{dK=M e} + 1 \right)$.

There are two cases to consider.

If $1 - 2 \min m_U + 4 \frac{2}{\min} M_U^2 > 0$, then from equation 18 it follows that

$$\begin{aligned}
E_k k^2 & \leq 2 \min b + 4 \frac{2}{\min} B^2 + \frac{2 \min d}{4M_U^2} \\
& \leq E_k k^2 + 2 \left(b + 2B^2 + \frac{d}{4M_U^2} \right)
\end{aligned}$$

If $0 < 1 - 2 \min m_U + 4 \frac{2}{\min} M_U^2 < 1$, then iterating equation 18 gives

$$E_k k^2 \leq (1 - 2 \min m_U + 4 \frac{2}{\min} M_U^2)^k E_k k^2 + \frac{0b + 2 \frac{2}{\min} B^2 + \frac{d}{4M_U^2}}{\min m_U - 2 \frac{2}{\min} M_U^2} \quad (19)$$

$$E_k k^2 \leq \frac{2 \min}{m_U \min} (b + 2B^2 + \frac{d}{4M_U^2}) \quad (20)$$

¹Note: we only focus on the case when $\text{mod } M = 0$.

Now, we have

$$\max_{\theta} \sum_{j=1}^K (L_U^2 E_{\theta} k_j^2 + B^2) \\ (L_U^2 (\theta + 2(1 - \frac{\theta}{m_U \min})) (b + 2B^2 + d) + B^2) := C_0$$

Due to the expression of $\frac{\theta}{\min}$, C_0 is independent of θ . Then we denote $\frac{C_0}{L_U^2} (C_0 + d)$ as C_1 and we can derive

$$D_{KL}(\mu_{\theta} \parallel P_{k=1}^K \mu_k) \leq C_1 (\frac{3}{8} \frac{2K}{\theta}) + C_0 (\frac{K}{2} \frac{\theta}{\theta})$$

Then according to Proposition 3.1 in (Bolley & Villani, 2005) and Lemma 3.3 in (Raginsky et al., 2017), if we denote $\theta + 2b + 2d$ as C_2 , we can derive the following result:

$$W_2(\mu_{\theta}; P_{k=1}^K \mu_k) \leq (12 + C_2 (\frac{K}{\theta}))^{\frac{1}{2}} [D_{KL}(\mu_{\theta} \parallel P_{k=1}^K \mu_k)^{\frac{1}{2}} + D_{KL}(\mu_{\theta} \parallel P_{k=1}^K \mu_k)^{\frac{1}{4}}] \\ (12 + \frac{C_2 K}{2} \frac{\theta}{\theta})^{\frac{1}{2}} [(\frac{3C_1}{8} \frac{2K}{\theta} + \frac{K C_0}{2} \frac{\theta}{\theta})^{\frac{1}{2}} + (\frac{3C_1}{16} \frac{2K}{\theta} + \frac{K C_0}{4} \frac{\theta}{\theta})^{\frac{1}{4}}]$$

D.2 $W_2(P_{k=1}^K \mu_k; \mu_1)$

We can directly get the following results from (3.17) in (Raginsky et al., 2017) that there exist some positive constants $(C_3; C_4)$,

$$W_2(P_{k=1}^K \mu_k; \mu_1) \leq C_3 \exp(\sum_{k=1}^K C_4)$$

Now combining the bounds for $W_2(\mu_{\theta}; P_{k=1}^K \mu_k)$ and $W_2(P_{k=1}^K \mu_k; \mu_1)$, substituting $\theta = O(1/K)$, and noting $W_2(P_{k=1}^K \mu_k; \mu_1)$ decreases w.r.t K , we arrive at the bound stated in the theorem. □

E RELATION WITH SGLD

For the standard polynomially-decay-stepsize SGLD, the convergence rate is bounded as

$$W_2(\mu_K; \mu_1) \leq W_2(\mu_K; P_{k=1}^K \mu_{h_k}) + W_2(P_{k=1}^K \mu_{h_k}; \mu_1) \quad (21)$$

where $W_2(\mu_K; P_{k=1}^K \mu_{h_k}) \leq (6 + h_0 \sum_{k=1}^K \frac{1}{k})^{\frac{1}{2}}$

$$[(D_1 h_0^2 \frac{2}{6} + D_0 h_0 \sum_{k=1}^K \frac{1}{k})^{\frac{1}{2}} + (D_1 h_0^2 \frac{2}{16} + D_0 \frac{h_0}{2} \sum_{k=1}^K \frac{1}{k})^{\frac{1}{4}}]$$

and $W_2(P_{k=1}^K \mu_{h_k}; \mu_1) \leq C_3 \exp(\frac{\sum_{k=1}^K h_k}{C_4})$.

Proof of the bound of $W_2(\mu_K; \mu_1)$ in the standard SGLD Similar to the proof of $W_2(\mu_K; \mu_1)$ in cSGLD, we get the following update rule for SGLD with the stepsize following a polynomial decay i.e., $h_k = \frac{h_0}{k}$,

$$\mu_{k+1} = \mu_k + \epsilon_k \mathcal{U}_k(\mu_k) h_{k+1} + \sqrt{2h_{k+1}} \mu_{k+1} \quad (22)$$

Let μ_k denote the distribution of μ_k .

Since

$$W_2(\mu_K; \mu_1) \leq W_2(\mu_K; P_{k=1}^K \mu_{h_k}) + W_2(P_{k=1}^K \mu_{h_k}; \mu_1) \quad (23)$$

, we need to give the bounds for these two parts respectively.

E.1 $W_2(\tilde{\mu}_K; \mathbb{P}_{k=1}^K \mu_k)$

We first assume $\mathbb{E}(r \cup(w)) = r \cup(w)$; $\forall w \in \mathbb{R}^d$; which is a general assumption according to the way we choose the minibatch. Following the proof in (Raginsky et al., 2017) and the analysis of the SPOS method in (Zhang et al., 2018), we define the following $p(t)$ which will be used in the following proof:

$$p(t) = \frac{1}{K} \sum_{i=1}^K \mathbb{1}_{\{t < \frac{i}{K}\}} \mu_i \quad (24)$$

Then we focus on the following continuous-time interpolation of

$$\tilde{\mu}(t) = \int_0^t \mathbb{E} r \cup_{\tilde{\mu}(s)} \left(\frac{1}{K} \sum_{k=1}^K \mu_k \right) ds + \sqrt{\frac{1}{2}} \int_0^t dW_s^{(d)}; \quad (25)$$

$$(26)$$

where $\tilde{\mu} = \mathbb{E} r \cup_k$ for $t \in [\frac{i-1}{K}, \frac{i}{K}]$. And for each k , μ_k and μ_k have the same probability law μ_k .

Since $\tilde{\mu}(t)$ is not a Markov process, we define the following process which has the same one-time marginals as $\tilde{\mu}(t)$

$$V(t) = \int_0^t G_s(V(s)) ds + \sqrt{\frac{1}{2}} \int_0^t dW_s^{(d)} \quad (27)$$

with

$$G_t(x) := \mathbb{E} r \cup_{\tilde{\mu}(t)} \left(\frac{1}{K} \sum_{i=1}^K \mu_i \right) \cup_j(t) = x^5 \quad (28)$$

Let $P_V^t := \mathbb{L}(V(s) : 0 \leq s \leq t)$ and $P^t := \mathbb{L}(\tilde{\mu}(s) : 0 \leq s \leq t)$ and according to the proof of Lemma 3.6 in (Raginsky et al., 2017), we can derive the similar result for the relative entropy of P_V^t and P^t :

$$\begin{aligned} D_{KL}(P_V^t \parallel P^t) &= \int_0^t dP_V^t \log \frac{dP_V^t}{dP^t} \\ &= \frac{1}{4} \int_0^t \mathbb{E} k r \cup(V(s)) G_s(V(s)) k^2 ds \\ &= \frac{1}{4} \int_0^t \mathbb{E} k r \cup(\tilde{\mu}(s)) G_s(\tilde{\mu}(s)) k^2 ds \end{aligned}$$

The last line follows the fact that $U(s) = L(V(s)); 8s$.

Then we will let $\mathbb{P} = \mathbb{P}_{k=1}^K h_k$ and we can use the martingale property of integral to derive:

$$\begin{aligned}
& D_{KL}(\mathbb{P}_V^{\mathbb{P}_{k=1}^K h_k} \parallel \mathbb{P}_{k=1}^{\mathbb{P}_{k=1}^K h_k}) \\
&= \frac{1}{4} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} \text{kr} U(s) \quad G_s(s) k^2 ds \\
& \quad \frac{1}{2} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} \text{kr} U(s) \quad r \quad U(s) \quad h_i k^2 ds \\
& \quad + \frac{1}{2} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} \text{kr} U(s) \quad h_i \quad G_s(s) \quad h_i k^2 ds \\
& \quad \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} k(s) \quad h_i k^2 ds \tag{29}
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} \text{kr} U(s) \quad h_i \quad G_s(s) \quad h_i k^2 ds \tag{30}
\end{aligned}$$

For the first part (29), we consider some $\mathbb{P}_{k=1}^j h_k; \mathbb{P}_{k=1}^{j+1} h_k$, the following equation holds:

$$\begin{aligned}
& U(s) \quad U(s) \quad h_k \\
&= (s) \quad h_k r \quad U_k(s) + \mathbb{P}_{k=1}^j (W_s^{(d)} \quad W_{\mathbb{P}_{k=1}^j h_k}^{(d)}) \\
&= (s) \quad h_k r \quad U(s) + (s) \quad h_k (r \quad U(s) \quad r \quad U_k(s)) + \mathbb{P}_{k=1}^j (W_s^{(d)} \quad W_{\mathbb{P}_{k=1}^j h_k}^{(d)}) \tag{31}
\end{aligned}$$

Thus, we can use Lemma 3.1 and 3.2 in (Raginsky et al., 2017) for the following result:

$$\begin{aligned}
& \mathbb{E} k(s) \quad h_k k^2 \\
& \quad 3h_{j+1}^2 \mathbb{E} \text{kr} U(s) k^2 + 3h_{j+1}^2 \mathbb{E} \text{kr} U(s) \quad r \quad U_k(s) k^2 + 6h_{j+1} d \\
& \quad 12h_{j+1}^2 (L_U^2 \mathbb{E} k_j k^2 + B^2) + 6h_{j+1} d
\end{aligned}$$

Hence we can bound the first part, (choosing $\mathbb{P} = 1$),

$$\begin{aligned}
& \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^{K-1} \int_{k=1}^{j+1} \mathbb{E} k(s) \quad h_k k^2 ds \\
& \quad \frac{L_U^2}{2} \mathbb{E} \int_{j=0}^{K-1} 12h_{j+1}^3 (L_U^2 \mathbb{E} k_j k^2 + B^2) + 6h_{j+1}^2 d \\
& \quad L_U^2 \max_{j=0}^{K-1} 6(L_U^2 \mathbb{E} k_j k^2 + B^2) + 3d \left(\sum_{j=0}^{K-1} h_{j+1}^2 \right) \\
& \quad L_U^2 \max_{j=0}^{K-1} 6(L_U^2 \mathbb{E} k_j k^2 + B^2) + 3d \frac{2}{6} h_0^2 \tag{32}
\end{aligned}$$

where the last line follows from the fact that

$$\mathbb{E} \int_{j=0}^{K-1} \frac{1}{(j+1)^3} \quad \mathbb{E} \int_{j=0}^{K-1} \frac{1}{(j+1)^2} \quad \mathbb{E} \int_{j=0}^{K-1} \frac{1}{(j+1)^2} = \frac{2}{6}$$

The second part (30) can be bounded as follows:

$$\begin{aligned}
& \frac{1}{2} \int_{j=0}^K \int_{k=1}^{j+1} \mathbb{E}_{k \sim P_{k=1}^{j+1}} \mathbb{E}_{s \sim U(\cdot | h_i)} G_s(\cdot | h_i) k^2 ds \\
&= \frac{1}{2} \int_{j=0}^K h_{j+1} \mathbb{E}_{k \sim U(\cdot | r_{\cdot} | \cdot)} k^2 \\
& \leq \max_{j=0}^{K-1} (L_U^2 \mathbb{E}_{k \sim P_{k=1}^{j+1}} k^2 + B^2) h_{j+1} \\
& \leq \max_{j=1}^K (L_U^2 \mathbb{E}_{k \sim P_{k=1}^j} k^2 + B^2) (h_0 + \frac{1}{j})
\end{aligned}$$

Due to the data-processing inequality for the relative entropy, we have

$$\begin{aligned}
D_{KL}(\tilde{\kappa} \prod_{k=1}^K h_k) & \leq D_{KL}(P_V^t \prod_{k=1}^K P^t) \\
& \leq \frac{L_U^2}{2} \int_{j=0}^K \int_{k=1}^{j+1} \mathbb{E}_{k \sim P_{k=1}^{j+1}} \mathbb{E}_{s \sim U(\cdot | h_i)} k^2 ds \\
& \quad + \frac{1}{2} \int_{j=0}^K \int_{k=1}^{j+1} \mathbb{E}_{k \sim P_{k=1}^{j+1}} \mathbb{E}_{s \sim U(\cdot | h_i)} G_s(\cdot | h_i) k^2 ds \\
& \leq L_U^2 \max_{j=0}^{K-1} (6(L_U^2 \mathbb{E}_{k \sim P_{k=1}^{j+1}} k^2 + B^2) + 3d) \frac{h_0^2}{6} \\
& \quad + \max_{j=1}^K (L_U^2 \mathbb{E}_{k \sim P_{k=1}^j} k^2 + B^2) (h_0 + \frac{1}{j})
\end{aligned}$$

Similar to the proof of cSGLD, we have

$$\max_{j=0}^{K-1} (L_U^2 \mathbb{E}_{k \sim P_{k=1}^{j+1}} k^2 + B^2) \leq D_0$$

Then we denote $\frac{L_U^2}{6}(D_0 + d)$ as D_1 and we can derive

$$D_{KL}(\tilde{\kappa} \prod_{k=1}^K h_k) \leq D_1 h_0^2 \frac{2}{6} + D_0 h_0 \sum_{j=1}^K \frac{1}{j}$$

Then according to Proposition 3.1 in (Bolley & Villani, 2005) and Lemma 3.3 in (Raginsky et al., 2017), if we denote $D_0 + 2b + 2d$ as D_2 , we can derive the following result,

$$\begin{aligned}
W_2(\tilde{\kappa}; \prod_{k=1}^K h_k) & \leq [12 + D_2 \sum_{k=1}^K h_k]^{1=2} [(D_{KL}(\tilde{\kappa} \prod_{k=1}^K h_k))^{1=2} + (D_{KL}(\tilde{\kappa} \prod_{k=1}^K h_k) = 2)^{1=4}] \\
&= [12 + D_2 \sum_{j=1}^K \frac{1}{j}]^{1=2} [(D_1 h_0^2 \frac{2}{6} + D_0 h_0 \sum_{j=1}^K \frac{1}{j})^{1=2} + (D_1 h_0^2 \frac{2}{12} + D_0 h_0 \sum_{j=1}^K \frac{1}{2j})^{1=4}]
\end{aligned}$$

Now we derive the bound for $W_2(\tilde{\kappa}; \prod_{k=1}^K h_k)$.

$$E.2 \quad W_2(\prod_{k=1}^K h_k; 1)$$

We can directly get the following results from (3.17) in (Raginsky et al., 2017) that there exist some positive constants $(C_3; C_4)$,

$$W_2(\prod_{k=1}^K h_k; 1) \leq C_3 \exp(\sum_{k=1}^K h_k = C_4)$$

□

Based on the convergence error bounds, we discuss an informal comparison with standard SGLD. Consider the following two cases. We must emphasize that since the $W_2(\sim K; \prod_{k=1}^K h_k)$ in the equation 9 increases w.r.t. K , our h_0 must be set small enough in practice. Hence, in this informal comparison, we also set h_0 small enough to make $W_2(\sim K; \prod_{k=1}^K h_k)$ less important.

i) If the initial stepsizes satisfy $h_0 \ll h_k$, our algorithm cSGLD runs much faster than the standard SGLD in terms of the amount of “diffusion time”, i.e., the “t” indexing μ in the continuous-time SDE mentioned above. This result follows from $\sum_{k=1}^K h_k = \frac{K h_0}{2}$ and $\prod_{k=1}^K h_k = \frac{K h_0}{K} = O(h_0 \log K)$. In standard SGLD, since the error described by $W_2(\sim K; \prod_{k=1}^K h_k)$ increases w.r.t. K , h_0 needs to be set small enough in practice to reduce the error. Following the general analysis of SGLD in (Raginsky et al., 2017; Xu et al., 2017), the dominant term in the decomposition equation 21 will be $W_2(\prod_{k=1}^K h_k; 1)$ since it decreases exponentially fast with the increase of $\prod_{k=1}^K h_k$ and $W_2(\sim K; \prod_{k=1}^K h_k)$ is small due to the setting of small h_0 . Since $\sum_{k=1}^K h_k$ increases much faster in our algorithm than the term $\prod_{k=1}^K h_k$ in standard SGLD, our algorithm thus endows less error for K iterations, i.e., $W_2(\prod_{k=1}^K h_k; 1) \ll W_2(\sim K; \prod_{k=1}^K h_k)$. Hence, our algorithm outperforms standard SGLD, as will be verified in our experiments.

ii) Instead of setting the h_0 small enough, one may consider increasing h_0 to make standard SGLD run as “fast” as our proposed algorithm, i.e., $\sum_{k=1}^K h_k = \frac{K h_0}{2}$. Now the $W_2(\prod_{k=1}^K h_k; 1)$ in equation 21 is almost the same as $W_2(\prod_{k=1}^K h_k; 1)$ in equation 9. However, in this case, it is worth noting that h_0 scales as $O(\frac{1}{\log K})$. We can notice that h_0 is much larger than the h_k and thus the $W_2(\sim K; \prod_{k=1}^K h_k)$ cannot be ignored. Now the h_0^2 term in $W_2(\sim K; \prod_{k=1}^K h_k)$ would scale as $O(\frac{1}{\log^2 K})$, which makes $W_2(\sim K; \prod_{k=1}^K h_k)$ in equation 21 much larger than our $W_2(\prod_{k=1}^K h_k; 1)$ defined in equation 9 since $O(\frac{1}{\log^2 K}) \gg O(\frac{1}{\log K})$. Again, our algorithm cSGLD achieves a faster convergence rate than standard SGLD.

F COMBINING SAMPLES

In cyclical SG-MCMC, we obtain samples from multiple modes of a posterior distribution by running the cyclical step size schedule for many periods. We now show how to effectively utilize the collected samples. We consider each cycle exploring different part of the target distribution on a metric space \mathcal{D} . As we have M cycles in total, the m th cycle characterizes a local region \mathcal{D}_m , defining the “sub-posterior” distribution $p_m(\cdot | \mathcal{D}_m) = \frac{p(\cdot | \mathcal{D}_m)}{w_m}$; with $w_m = \int_{\mathcal{D}_m} p(\cdot | \mathcal{D}_m) d\mu$; where w_m is a normalizing constant. For a testing function $f(\cdot)$, we are often interested in its true posterior expectation $\mathbb{E}[f(\cdot)] = \int_{\mathcal{D}} f(\cdot) p(\cdot | \mathcal{D}) d\mu$. The sample-based estimation is

$$\hat{f} = \sum_{m=1}^M \frac{K_m}{K} \hat{f}_m \quad \text{with} \quad \hat{f}_m = \frac{1}{K_m} \sum_{j=1}^{K_m} f(D_j^{(m)}); \quad (33)$$

where K_m is the number of samples from the m th cycle, and $(D_j^{(m)})_{j=1}^{K_m}$.

The weight for each cycle w_m is estimated using the harmonic mean method (Green, 1995; Raftery et al., 2006): $\hat{w}_m = \left[\frac{1}{K_m} \sum_{j=1}^{K_m} \frac{1}{p(D_j^{(m)})} \right]^{-1}$: This approach provides a simple and consistent estimator, where the only additional cost is to traverse the training dataset to evaluate the likelihood $p(D_j^{(m)})$ for each sample $D_j^{(m)}$. We evaluate the likelihood once off-line and store the result for testing.

If \mathcal{D}_m are not disjoint, we can assume new sub-regions \mathcal{D}_m^* which are disjoint and compute the estimator as following

$$\hat{f}_m = \frac{1}{\bar{n}_m} \sum_{m=1}^M \sum_{j=1}^{K_m} f(D_j^{(m)}) 1_{\mathcal{D}_m^*}(D_j^{(m)})$$

where

$$\bar{n}_m = \prod_{j=1}^m 1_{\sim_m(j)}$$

and $1_{\sim_m(j)}$ equals 1 only when $j \in \sim_m$. By doing so, our estimator still holds even if \sim_m are not disjoint.

G THEORETICAL ANALYSIS UNDER CONVEX ASSUMPTION

Firstly, we introduce the following SDE

$$d_t = -r U(t)dt + \sqrt{p} dW_t; \quad (34)$$

Let π_t denote the distribution of t , and the stationary distribution of equation 34 (see D), which means $\pi = p(jD)$.

However, the exact evaluation of the gradient is computationally expensive. Hence, we need to adopt noisy evaluations of U . For simplicity, we assume that at any point, we can observe the value

$$r \mathbf{u}_k = r U(\mathbf{x}_k) + \mathbf{k}_k$$

where $\mathbf{k}_k : k = 0; 1; 2; \dots$ is a sequence of random (noise) vectors. Then the algorithm is defined as:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k r \mathbf{u}_k + \sqrt{2\eta_k} \mathbf{k}_{k+1} \quad (35)$$

Further, let π_k denote the distribution of \mathbf{x}_k .

Following the existing work in (Dalalyan & Karagulyan, 2019), we adopt the following standard assumptions summarized in Assumption 5,

Assumption 5.

For some positive constants m and M , it holds

$$|U(\mathbf{x}) - U(\mathbf{y})| \leq r |U(\mathbf{x}) - U(\mathbf{y})| \leq (m/2)k_2 \|\mathbf{x} - \mathbf{y}\|_2$$

$$|kr U(\mathbf{x}) - r U(\mathbf{y})| \leq M k_2 \|\mathbf{x} - \mathbf{y}\|_2$$

for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$

(bounded bias) $\mathbb{E}[k(\mathbf{x}_k - \mathbf{x}_k^*)] \leq \eta^2 d$

(bounded variance) $\mathbb{E}[k_k - \mathbb{E}(k_k)] \leq \eta^2 d$

(independence of updates) \mathbf{k}_{k+1} in equation 35 is independent of $\{\mathbf{k}_1; \mathbf{k}_2; \dots; \mathbf{k}_k\}$

G.1 THEOREM

Under Assumption 5 in Appendix and $\eta \in (0; \frac{1}{m} \wedge \frac{2}{M})$, if we define the η_{\min} as $\frac{\eta}{2} \cos \frac{\pi}{dK=M} + 1$, we can derive the the following bounds.

If $m \eta_{\min} + M \eta \leq 2$, then $W_2(\pi_{k+1}; \pi)$

$$(1 - m \eta_{\min})^k W_2(\pi_0; \pi) + \frac{(1.65M \eta^3 + \eta^2) d^{1.5}}{m \eta_{\min}} + \frac{\eta^2 \eta_{\min} d^{1.5}}{1.65M \eta^2 + \eta^2 m \eta_{\min}}; \quad (36)$$

If $m \eta_{\min} + M \eta > 2$, then $W_2(\pi_{k+1}; \pi)$

$$(1 - (2 - M \eta))^{k+1} W_2(\pi_0; \pi) + \frac{(1.65M \eta^3 + \eta^2) d^{1.5}}{2 - M \eta} + \frac{\eta^2 \eta_{\min} d^{1.5}}{1.65M \eta^2 + \eta^2 \frac{2 - M \eta}{2}}; \quad (37)$$

where the $M; m; \eta$ are some positive constants defined in Assumption 5

G.2 PROOF

Proof. According to the equation 1, we can find that the stepsize η_k varies from η_0 to η_{\min} , where η_{\min} is defined as $\eta_{\min} = \frac{\eta_0}{2} \cos \frac{\text{mod}(dK-Me-1; dK-Me)}{dK-Me} + 1$. When $0 < \eta_0 < \min(2=M; 1=m)$, it is easy for us to know that $0 < \eta_k < \min(2=M; 1=m)$ for every $k > 0$. Then we can derive that all the η_k , $\max(1-m; M-k-1)$ will satisfy $0 < \eta_k < 1$. Now according to the Proposition 2 in the (Dalalyan & Karagulyan, 2019), we can derive the result that

$$W_2(\eta_{k+1}; \gamma)^2 \leq \eta_{k+1} W_2(\eta_k; \gamma) + 1.65M(\frac{3}{k+1}d)^{1=2} + \eta_{k+1} \bar{\rho}g^2 + \frac{2}{k+1}d \quad (38)$$

Then we will use another lemma derived from (Dalalyan & Karagulyan, 2019).

Lemma 2. *If A, B, C are non-negative numbers such that $A \geq (0, 1)$ and the sequence of non-negative numbers y_k satisfies the following inequality*

$$y_{k+1}^2 \leq [(1-A)y_k + C]^2 + B^2$$

for every integer $k > 0$. Then,

$$y_k \leq (1-A)^k y_0 + \frac{C}{A} + \frac{B^2}{C + \frac{B^2}{AB}}$$

Using Lemma 2, we can finish our proof now.

If $m_{\min} + M_0 \leq 2$, the η_k will satisfy $\eta_k \leq 1 - m_{\min}$ for every $k > 0$. Then the equation 38 will turn into

$$W_2(\eta_{k+1}; \gamma)^2 \leq (1 - m_{\min}) W_2(\eta_k; \gamma) + 1.65M(\frac{3}{k+1}d)^{1=2} + \eta_0 d^{1=2} g^2 + (\eta_0 d^{1=2})^2$$

for every $k > 0$. Then we can set $A = m_{\min}$, $C = 1.65M(\frac{3}{k+1}d)^{1=2} + \eta_0 d^{1=2}$, $B = \eta_0 d^{1=2}$ and we can get the result.

If $m_{\min} + M_0 > 2$, the η_k will satisfy $\eta_k \leq M_0 - 1$ for every $k > 0$. Then the equation 38 will turn into

$$W_2(\eta_{k+1}; \gamma)^2 \leq [1 - (2 - M_0)] W_2(\eta_k; \gamma) + 1.65M(\frac{3}{k+1}d)^{1=2} + \eta_0 d^{1=2} g^2 + (\eta_0 d^{1=2})^2$$

for every $k > 0$. Then we can set $A = 2 - M_0$, $C = 1.65M(\frac{3}{k+1}d)^{1=2} + \eta_0 d^{1=2}$, $B = \eta_0 d^{1=2}$ and we can get the result.

□

G.3 HYPERPARAMETERS SETTING DISCUSSION

There are several hyperparameters in Algorithm 1. We now discuss how to set them in practice. Given the training budget K , there is a trade-off between the number of cycles M and the cycle length. We find that $M \geq [3; 6]$ usually works well in practice. η needs tuning for different tasks by cross-validation. Generally, if the length of each cycle is relatively short, η needs to be large (e.g. 0.8) so that the sampler has enough time to reach a good region before starting sampling.