

---

# Enhancing Low-Precision Sampling via Stochastic Gradient Hamiltonian Monte Carlo

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Low-precision training has emerged as a promising low-cost technique to enhance  
2 the training efficiency of deep neural networks without sacrificing much accuracy.  
3 Its Bayesian counterpart can further provide uncertainty quantification and im-  
4 proved generalization accuracy. This paper investigates low-precision samplers  
5 via Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) with low-precision  
6 and full-precision gradients accumulators for both strongly log-concave and non-  
7 log-concave distributions. Theoretically, our results show that, to achieve  $\epsilon$ -error  
8 in the 2-Wasserstein distance for non-log-concave distributions, low-precision  
9 SGHMC achieves quadratic improvement ( $\tilde{O}(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1}))$ ) compared to  
10 the state-of-the-art low-precision sampler, Stochastic Gradient Langevin Dynam-  
11 ics (SGLD) ( $\tilde{O}(\epsilon^{-4}\lambda^{*-1}\log^5(\epsilon^{-1}))$ ). Moreover, we prove that low-precision  
12 SGHMC is more robust to the quantization error compared to low-precision SGLD  
13 due to the robustness of the momentum-based update w.r.t. gradient noise. Em-  
14 pirically, we conduct experiments on synthetic and MNIST, CIFAR-10 & CIFAR-  
15 100 datasets which successfully validate our theoretical findings. Our study high-  
16 lights the potential of low-precision SGHMC as an efficient and accurate sampling  
17 method for large-scale and resource-limited deep learning.

## 18 1 Introduction

19 In recent years, deep neural networks (DNNs) have achieved remarkable success, accompanied by  
20 an increase in model complexity [Simonyan and Zisserman, 2014, He et al., 2016, Vaswani et al.,  
21 2017, Radford et al., 2018, Chen et al., 2023]. Consequently, there is a growing interest in utilizing  
22 low-precision optimization techniques to address the computational and memory costs associated  
23 with these complex models [Wang et al., 2018, Banner et al., 2018, Wu et al., 2018, Lin et al.,  
24 2019, Sun et al., 2019, Wortsman et al., 2023]. As a counterpart of low-precision optimization, low-  
25 precision sampling is relatively unexplored but has shown promising preliminary results. Zhang  
26 et al. [2022] studied the effectiveness of Stochastic Gradient Langevin Dynamics (SGLD) [Welling  
27 and Teh, 2011] in the context of low-precision arithmetic, highlighting its superiority over the op-  
28 timization counterpart, Stochastic Gradient Descent (SGD). This superiority stems from SGLD’s  
29 inherent robustness to system noise compared with SGD.

30 Other than SGLD, Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) [Chen et al., 2014]  
31 is another popular gradient-based sampling method, closely related to the underdamped Langevin  
32 dynamics. Recently, Cheng et al. [2018], Gao et al. [2022] have shown that the SGHMC converges  
33 to its target distribution faster than the best-known convergence rate of SGLD in the 2-Wasserstein  
34 distance under both strongly log-concave and non-log-concave assumptions. Beyond this, SGHMC  
35 is analogous to stochastic gradient methods augmented with momentum, which is shown to have

36 more robust updates w.r.t. gradient estimation noise [Liu et al., 2020]. Note that the stochastic error  
 37 induced by the quantization function in the low-precision update is equivalent to an extra noise of  
 38 the stochastic gradient, causing an increase in the gradient variance. Thus, we believe the SGHMC  
 39 is particularly suited for low-precision arithmetic.

40 Our main contributions of this paper are threefold:

41 First, we conduct the first study of low-precision SGHMC. We adopt low-precision arithmetic (in-  
 42 cluding full- and low-precision gradient accumulators and variance correction (VC) version of low-  
 43 precision gradient accumulators) to SGHMC.

44 Second, we provide a comprehensive theoretical analysis of low-precision SGHMC for both strongly  
 45 log-concave and non-log-concave target distributions. All our theoretical results are summarized in  
 46 Table 3 (deferred in Appendix A), where we compare the 2-Wasserstein convergence limit and the  
 47 required gradient complexity. Our analysis exhibits the superiority of HMC-based low-precision  
 48 algorithms over SGLD counterpart w.r.t. convergence speed and robustness to quantization error,  
 49 especially under the non-log concave distributions.

50 Third, we provide promising empirical results in deep learning. We show the sampling capabilities  
 51 of HMC-based low-precision algorithms and the effectiveness of the VC function in both strongly  
 52 log-concave and non-log-concave target distributions. We also provide evidence of the superior  
 53 performance of HMC-based low-precision algorithms compared to SGLD in real-world tasks.

54 In summary, low-precision SGHMC emerges as a compelling alternative to standard SGHMC due  
 55 to its ability to enhance speed and memory efficiency without sacrificing accuracy.

## 56 2 Preliminaries

### 57 2.1 Stochastic Gradient Hamiltonian Monte Carlo

58 Given a dataset  $D$ , a model with weights (i.e., model parameters)  $\mathbf{x} \in \mathbb{R}^d$ , and a prior  $p(\mathbf{x})$ , we  
 59 are interested in sampling from the posterior  $p(\mathbf{x}|D) \propto \exp(-U(\mathbf{x}))$ , where  $U(\mathbf{x})$  is some energy  
 60 function. In order to sample from the target distribution, SGHMC [Chen et al., 2014] is proposed and  
 61 strongly related to the underdamped Langevin dynamics. Cheng et al. [2018] proposes the following  
 62 discretization of underdamped Langevin dynamics (9) with stochastic gradient:

$$\begin{aligned} \mathbf{v}_{k+1} &= \mathbf{v}_k e^{-u\gamma} - u\gamma^{-1}(1 - e^{-u\gamma})\nabla\tilde{U}(\mathbf{x}_k) + \xi_k^{\mathbf{v}} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \gamma^{-1}(1 - e^{-u\gamma})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-u\gamma} - 1)\nabla\tilde{U}(\mathbf{x}_k) + \xi_k^{\mathbf{x}}, \end{aligned} \quad (1)$$

63 where  $u, \gamma$  denote the hyperparameters of inverse mass and friction respectively,  $\nabla\tilde{U}$  is unbiased  
 64 gradient estimation of  $U$  and  $\xi_k^{\mathbf{v}}$ , and  $\eta$  is the step size.  $\xi_k^{\mathbf{x}}$  are normal distributed in  $\mathbb{R}^d$  satisfying  
 65 that :

$$\begin{aligned} \mathbb{E}\xi_k^{\mathbf{v}}(\xi_k^{\mathbf{v}})^{\top} &= u(1 - e^{-2u\gamma}) \cdot \mathbf{I}, \\ \mathbb{E}\xi_k^{\mathbf{x}}(\xi_k^{\mathbf{x}})^{\top} &= u\gamma^{-2}(2\gamma\eta + 4e^{-u\gamma} - e^{-2u\gamma} - 3) \cdot \mathbf{I}, \\ \mathbb{E}\xi_k^{\mathbf{x}}(\xi_k^{\mathbf{v}})^{\top} &= u\gamma^{-1}(1 - 2e^{-u\gamma} + e^{-2u\gamma}) \cdot \mathbf{I}. \end{aligned} \quad (2)$$

### 66 2.2 Low-Precision Quantization

67 Two popular formats to represent low-precision numbers are known as the *fixed point* (FP) and *block*  
 68 *floating point* [Song et al., 2018] (BFP). The quantization error which is defined as the gap between  
 69 two adjacent representable numbers is denoted as  $\Delta$ . Furthermore, all representable numbers are  
 70 truncated to an upper limit  $\bar{U}$  and a lower limit  $\bar{L}$ .

71 Given the low-precision number representation, a quantization function is desired to round real-  
 72 valued numbers to their low-precision counterparts. Two common quantization functions are *de-*  
 73 *terministic rounding* and *stochastic rounding*. The deterministic rounding function, denoted as  $Q^d$ ,  
 74 quantizes a number to its nearest representable neighbor. The stochastic rounding denoted as  $Q^s$   
 75 (refer to (10) of Appendix A), randomly quantizes a number to the two closest representable neigh-  
 76 bors satisfying the unbiased condition, i.e.  $\mathbb{E}[Q^s(\theta)] = \theta$ . In what follows, we use  $Q_W$  and  $Q_G$

77 to denote the stochastic rounding quantizer we used for the weights and gradients respectively, al-  
78 lowing different quantization errors. But for simplicity in the analysis and experiments, we use the  
79 same number of bits to represent the weights and gradients.

### 80 3 Low-Precision Stochastic Gradient Hamiltonian Monte Carlo

81 In this section, we investigate the convergence property of low-precision SGHMC for non-log-  
82 concave target distributions. We defer the analysis of the low-precision SGHMC under strongly  
83 log-concave target distributions, as well as the analysis of low-precision SGLD [Zhang et al., 2022]  
84 to Appendix A and B respectively. All of our theorems are based on the fixed point representation  
85 and omit the clipping effect.

86 In order to derive a convergence analysis for non-log-concave target distribution, we assume the  
87 energy function  $U(\cdot)$  is  $M$ -smooth (Assumption 1) also satisfied the dissapitiveness assumption (As-  
88 sumption 3), and the mean squared error of stochastic gradients is bounded by constant  $\sigma^2$  (As-  
89 sumption 4). Detailed assumptions and explanations are deferred in Appendix A. In the statement  
90 of theorems, the big-O notation  $\tilde{O}$  gives explicitly dependence on the quantization error  $\Delta$  and con-  
91 centration parameters  $(\lambda^*, \mu^*)$  but hides multiplicative terms that depend polynomially on the other  
92 parameters (e.g., dimension  $d$ , friction  $\gamma$ , inverse mass  $u$  and gradients variance  $\sigma^2$ ).

#### 93 3.1 Full- and low-Precision Gradient Accumulators

94 Adopting the updating rule in equations 1, we propose the low-precision SGHMC with full gradient  
95 accumulator (SGHMCLP-F) as the following:

$$\begin{aligned} \mathbf{v}_{k+1} &= \mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta})Q_G(\nabla\tilde{U}(Q_W(\mathbf{x}_k))) + \xi_k^{\mathbf{v}} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \gamma^{-1}(1 - e^{-\gamma\eta})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-\gamma\eta} - 1)Q_G(\nabla\tilde{U}(Q_W(\mathbf{x}_k))) + \xi_k^{\mathbf{x}}, \end{aligned} \quad (3)$$

96 The storage and computation costs can be further reduced by the low-precision gradient accumula-  
97 tors, i.e., the low-precision SGHMC with low-precision gradient accumulators (SGHMCLP-L):

$$\begin{aligned} \mathbf{v}_{k+1} &= Q_W\left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta})Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{v}}\right), \\ \mathbf{x}_{k+1} &= Q_W\left(\mathbf{x}_k + \gamma^{-1}(1 - e^{-\gamma\eta})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-\gamma\eta} - 1)Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{x}}\right). \end{aligned} \quad (4)$$

98 Our analysis for the above two algorithms utilizes similar techniques in Raginsky et al. [2017].

99 **Theorem 1** (Informal version of Theorem 5). *Given the smoothness, dissapitivity and assumption*  
100 *for stochastic gradients, let  $p^*$  denote the target distribution of  $\mathbf{x}$  and  $\mathbf{v}$ . Given initialization  $\mathbf{x}_0 =$*   
101  *$\mathbf{v}_0 = 0$  and  $\gamma^2 \leq 4Mu$ , for some sufficiently small  $\epsilon$  and step size  $\eta$ , the  $K$ -th iteration of the*  
102 *SGHMCLP-F update (3), i.e.,  $\mathbf{x}_K$  and  $\mathbf{v}_K$ , satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) \leq \tilde{O}\left(\epsilon + \sqrt{\Delta \log(1/\epsilon)}\right), \quad (5)$$

103 for some  $K$  satisfying

$$K = \tilde{O}\left(\frac{1}{\epsilon^2 \mu^{*2}} \log^2\left(\frac{1}{\epsilon}\right)\right),$$

104 where  $\mu^*$  is a constant w.r.t. dimension  $d$ , denoting the concentration rate of the underdamped  
105 Langevin dynamics [Zou et al., 2019].

106 **Theorem 2** (Informal version of Theorem 7). *Given the smoothness, dissapitivity and assumption*  
107 *for stochastic gradients, let  $p^*$  denote the target distribution of  $\mathbf{x}$  and  $\mathbf{v}$ . Given initialization  $\mathbf{x}_0 =$*   
108  *$\mathbf{v}_0 = 0$  and  $\gamma^2 \leq 4Mu$ , for some sufficiently small  $\epsilon$  and step size  $\eta$ , the  $K$ -th iteration of the*  
109 *SGHMCLP-L update (4), i.e.,  $\mathbf{x}_K$  and  $\mathbf{v}_K$ , satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) = \tilde{O}\left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\} \log\left(\frac{1}{\epsilon}\right) + \frac{\log^{3=2}\left(\frac{1}{\epsilon}\right)}{\epsilon^2} \sqrt{\Delta}}\right), \quad (6)$$

110 for some  $K$  satisfying

$$K = \tilde{O}\left(\frac{1}{\epsilon^2 \mu^{*2}} \log^2\left(\frac{1}{\epsilon}\right)\right).$$

111 Similar to the convergence result of full-precision SGHMC or SGLD [Raginsky et al., 2017, Gao  
112 et al., 2022], the above upper bound (5) of SGHMCLP-F contains a  $\epsilon$  term and a  $\log(\epsilon^{-1})$  term. The  
113 difference is that for the SGHMCLP-F algorithm, the quantization error  $\Delta$  affects the multiplicative  
114 constant of the  $\log(\epsilon^{-1})$  term. Without  $\Delta$ , one can choose a small  $\epsilon$  and a larger batch size (i.e., a  
115 smaller  $\sigma^2$ ) to offset  $\log(\epsilon^{-1})$  term, such that the 2-Wasserstein distance can be sufficiently small.  
116 With the same technical tools, we conduct a similar convergence analysis of SGLDLF-P for non-log-  
117 concave target distributions (refer to Theorem 10 of Appendix B). Comparing Theorems 1 and 10,  
118 we show that SGHMCLP-F can achieve lower 2-Wasserstein (i.e.  $\tilde{O}\left(\epsilon + (\log(\epsilon^{-1})\Delta)^{1=2}\right)$ ) ver-  
119 sus  $\tilde{O}\left(\epsilon + \log(\epsilon^{-1})\Delta^{1=2}\right)$ ) distance for non-log-concave target distribution within fewer iterations  
120 (i.e.,  $\tilde{O}\left(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1})\right)$  versus  $\tilde{O}\left(\epsilon^{-4}\lambda^{*-1}\log^5(\epsilon^{-1})\right)$ ).

121 We verify the advantage of SGHMCLP-F over SGLDLF-P by our simulations in section 4.

122 As for SGHMCLP-L, which additionally quantizes the weights after each update, a small stepsize  
123 can result in staying at the starting point. In such cases, ensuring convergence becomes challenging,  
124 and the output of the SGHMCLP-L has a worse convergence upper bound compared to Theorem 1.  
125 Empirically, we observe that the output  $\mathbf{x}_K$ 's distribution has an overdispersion problem (i.e. Fig-  
126 ure 1 (a) and 5 (a)). In Theorem 11, we generalize the result of the naïve SGLDLP-L in [Zhang  
127 et al., 2022] to non-log-concave target distribution. Similarly, we observe that SGHMCLP-L needs  
128 fewer iterations than SGLDLP-L in terms of the order w.r.t.  $\epsilon$  and achieves better upper bound  
129  $\tilde{O}\left(\epsilon^{-2}\log^{3=2}(\epsilon^{-1})\Delta^{1=2}\right)$  versus  $\tilde{O}\left(\epsilon^{-4}\log^5(\epsilon^{-1})\Delta^{1=2}\right)$ .

### 130 3.2 Variance Correction

131 To resolve the overdispersion caused by the low-precision gradient accumulators, Zhang et al. [2022]  
132 propose a quantization function  $Q^{vc}$  (refer to Algorithm 1 in Appendix A) that directly samples from  
133 the discrete weight space instead of quantizing a real-valued Gaussian sample. This quantization  
134 function aims to reduce the discrepancy between the ideal sampling variance (i.e., the required vari-  
135 ance of full-precision counterpart algorithms) and the actual sampling variance in our low-precision  
136 algorithms.

137 In this work, we study the effect of  $Q^{vc}$  on low-precision SGHMC. Let  $\text{Var}_{\mathbf{v}}^{hmc} = u(1 - e^{-2})$   
138 and  $\text{Var}_{\mathbf{x}}^{hmc} = u\gamma^{-2}(2\gamma\eta + 4e^{-\gamma} - e^{-2} - 3)$ , the VC SGHMCLP-L can be done as:

$$\begin{aligned} \mathbf{v}_{k+1} &= Q^{vc}\left(\mathbf{v}_k e^{-\gamma} - u\gamma^{-1}(1 - e^{-\gamma})Q_G(\nabla\tilde{U}(\mathbf{x}_k)), \text{Var}_{\mathbf{v}}^{hmc}, \Delta\right) \\ \mathbf{x}_{k+1} &= Q^{vc}\left(\mathbf{x}_k + \gamma^{-1}(1 - e^{-\gamma})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-\gamma} - 1)Q_G(\nabla\tilde{U}(\mathbf{x}_k)), \text{Var}_{\mathbf{x}}^{hmc}, \Delta\right) \end{aligned} \quad (7)$$

139 Now, we are ready to present the convergence analysis of VC SGHMC-L.

140 **Theorem 3** (Informal version of Theorem 9). *Given the smoothness, dissipativity and assumption*  
141 *for stochastic gradients, let  $p^*$  denote the target distribution of  $\mathbf{x}$ . Given initialization  $\mathbf{x}_0 = \mathbf{v}_0 = 0$*   
142 *and  $\gamma^2 \leq 4Mu$ , for some sufficiently small  $\epsilon$  and step size  $\eta$ , the  $K$ -th iteration of the VC*  
143 *SGHMCLP-L update (4), i.e.,  $\mathbf{x}_K$ , satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K), p^*) = \tilde{O}\left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\} \log\left(\frac{1}{\epsilon}\right) + \frac{\log\left(\frac{1}{\epsilon}\right)}{\epsilon} \sqrt{\Delta}}\right), \quad (8)$$

144 for some  $K$  satisfying

$$K = \tilde{O}\left(\frac{1}{\epsilon^2\mu^{*2}} \log^2\left(\frac{1}{\epsilon}\right)\right).$$

145 Comparing with Theorem 2, the variance corrected quantization can improve the upper bound  
146 w.r.t.  $\epsilon$  from  $\tilde{O}\left(\epsilon^{-2}\log^{3=2}(\epsilon^{-1})\Delta^{1=2}\right)$  to  $\tilde{O}\left(\epsilon^{-1}\log(\epsilon^{-1})\Delta^{1=2}\right)$ . In Theorem 12, we gener-  
147 alize the result of the VC SGLDLP-L in [Zhang et al., 2022] to non-log-concave target distribu-  
148 tion. Similarly, we observe that VC SGHMCLP-L needs fewer iterations than VC SGLDLP-L in  
149 terms of the order w.r.t.  $\epsilon$  and achieves better upper bounds ( $\tilde{O}\left(\epsilon + \log(\epsilon^{-1})\epsilon^{-1}\Delta^{1=2}\right)$  versus  
150  $\tilde{O}\left(\epsilon + \log^3(\epsilon^{-1})\epsilon^{-2}\Delta^{1=2}\right)$ ).

(a) (b) (c)

Figure 1: Low-precision SGHMC on Gaussian distribution. (a): SGHMCLP-L. (b): VC SGHMCLP-L. (c): SGHMCLP-F.

(a) (b) (c)

Figure 2: Training NLL of low-precision SGHMC and SGLD on logistic model with MNIST in terms of different numbers of fractional bits. (a): Methods with full-precision gradient accumulators. (b): Methods with low-precision gradient accumulators. (c): Variance corrected quantization.

151 Interestingly, the naïve SGHMCLP-L has similar dependence on the quantization error with VC  
152 SGLDLP-L but saves more computation resources since the variance corrected quantization requires  
153 sampling discrete random variables. We verify our finding in Table 2.

## 154 4 Experiments

155 We assess the performance of the proposed low-precision SGHMC algorithms through sampling a  
156 Gaussian distribution and implementing a Bayesian logistic regression to the MNIST dataset (Sec-  
157 tion 4.1), and training a Bayesian ResNet-18 on the CIFAR-10 and CIFAR-100 datasets (Section  
158 4.2). We compare our proposed algorithms with their SGLD counterparts. Details and additional ex-  
159 periment results (e.g., sampling Gaussian mixture distribution and MLP training on MNIST dataset)  
160 can be found in Appendix F. In all experiments, torch [Zhang et al., 2019] is employed for Low-  
161 Precision sampling with the same quantization.

### 162 4.1 Sampling Gaussian distributions & MNIST

163 We use a Gaussian distribution to represent the log-concave distribution. The simulation results  
164 are shown in Figure 1. It shows that the SGHMCLP-F samples the true Gaussian distribution  
165 well. Regarding the naïve SGHMCLP-L, we observe an overdispersion problem and the variance  
166 corrected function solves this problem.

167 We further examine the sampling performance of low-precision SGHMC and SGLD on real-world  
168 data. We use logistic models to represent the class of strongly log-concave distributions. The results  
169 are in Figure 2. We use fixed point numbers with 2 integer bits and vary the number of fractional  
170 bits which corresponds to varying the quantization gap. We report train negative log-likelihood  
171 (NLL) with different numbers of fractional bits in Figure 2. From the results on MNIST, we can  
172 see that when adopted to full-precision gradient accumulators low-precision SGHMC are robust to  
173 the quantization error. Even when we use only 2 fractional bits, SGHMCLP-F can still converge  
174 to a good distribution but with more iteration. As the precision error increases, both SGHMCLP-  
175 L and SGLDLP-L have a worse convergence pattern compared to SGHMCLP-F and SGLDLP-F.

(a) (b)

Figure 3: Log of training NLL of low-precision SGHMC and SGLD on ResNet-18 with CIFAR100 and constant step sizes. (a): 8-bit Fixed Point. (b): 8-bit Block Float Point.

Table 1: Test errors (%) of full-precision gradient accumulators on CIFAR with ResNet-18.

	32-bit Floating						8-bit Fixed Point						8-bit Block Floating Point					
	SGD		SGLD		SGHMC		SGD		SGLD		SGHMC		SGD		SGLD		SGHMC	
CIFAR-10	4.73	0.10	4.52	0.07	4.78	0.08	5.19	0.09	5.07	0.04	5.08	0.08	4.75	0.21	4.58	0.07	4.93	0.09
CIFAR-100	22.34	0.22	22.40	0.04	22.37	0.04	23.71	0.18	23.36	0.10	23.54	0.10	22.86	0.14	22.70	0.22	22.39	0.11

176 We showed empirically that SGHMCLP-L and VC SGHMCLP-L outperform SGLDLP-L and VC  
 177 SGLDLP in Figure 2, showing low-precision SGHMC is more robust to the quantization error.

## 178 4.2 CIFAR-10 & CIFAR-100

179 We consider computer vision tasks CIFAR10 and CIFAR100 on the ResNet-18. We use 8-bit num-  
 180 ber representation as it becomes increasingly popular and powered by new chips. We report the  
 181 average test errors over 3 runs in Tables 1 and 2. We use 8-bit fixed point (FP) and block floating  
 182 point (BFP) representing weights and gradients. SGHMCLP-F is comparable with SGDLP-F and the  
 183 naïve SGHMCLP-L significantly outperforms the SGLDLP-L and SGDLP-L across datasets. Fur-  
 184 thermore, from the result in Figure 3, we empirically show that the convergence speed of SGHMC  
 185 is way better than the SGLD. Besides the variance corrected quantization function can bring some  
 186 gain on the test accuracy, the performance of SGHMCLP-L is good enough and comparable with  
 187 the performance of VC SGLDLP-L. By using BFP, the performance of all low-precision methods  
 188 improves over fixed point, and we observe similar results as the FP.

## 189 5 Conclusion

190 We provide the first comprehensive investigation for low-precision SGHMC in both strongly log-  
 191 concave and non-log-concave target distributions with several variants of low-precision training.  
 192 In particular, we prove that for non-log-concave distributions, low-precision SGHMC with full-  
 193 precision, low-precision, and variance-corrected gradient accumulators, all achieve an acceleration  
 194 in iterations and have a better convergence upper bound w.r.t the quantization error compared to the  
 195 low-precision SGLD counterpart. Moreover, we study the improvement of variance-corrected quan-  
 196 tization applied to low-precision SGHMC under different cases. Under certain conditions, the naïve  
 197 SGHMCLP-L can replace the VC SGLDLP-L to get comparable results saving more computation

Table 2: Test errors (%) of low-precision gradient accumulators on CIFAR with ResNet-18.

	8-bit Fixed Point												8-bit Block Floating Point							
	SGD		SGLD		VC SGLD		SGHMC		VC SGHMC		SGD		SGLD		VC SGLD		SGHMC		VC SGHMC	
CIFAR-10	8.50	0.22	7.81	0.07	7.03	0.23	6.63	0.01	6.60	0.06	5.86	0.18	5.75	0.05	5.51	0.01	5.38	0.06	5.15	0.08
CIFAR-100	28.42	0.35	27.15	0.35	26.73	0.12	26.57	0.10	26.43	0.19	26.75	0.11	26.11	0.38	25.14	0.11	25.29	0.03	24.45	0.16

198 resources. We conduct empirical experiments on Gaussian, Gaussian mixture distribution, logistic  
199 regression, and Bayesian deep learning tasks to justify our theoretical findings.

## 200 References

- 201 R. Banner, I. Hubara, E. Hoffer, and D. Soudry. Scalable methods for 8-bit training of neural  
202 networks. *Advances in neural information processing systems*, 2018.
- 203 F. Bolley and C. Villani. Weighted chi-squared-kullback-pinsker inequalities and applications to trans-  
204 portation inequalities. *Annales de la Faculté des sciences de Toulouse: Mathématiques*, vol-  
205 ume 14, pages 331–352, 2005.
- 206 T. Chen, E. Fox, and C. Guestrin. Stochastic gradient hamiltonian monte carlo. *International*  
207 *conference on machine learning*, pages 1683–1691. PMLR, 2014.
- 208 X. Chen, C. Liang, D. Huang, E. Real, K. Wang, Y. Liu, H. Pham, X. Dong, T. Luong, C.-J. Hsieh,  
209 et al. Symbolic discovery of optimization algorithms. *arXiv preprint arXiv:2302.06675*, 2023.
- 210 X. Cheng, N. S. Chatterji, P. L. Bartlett, and M. I. Jordan. Underdamped langevin mcmc: A non-  
211 asymptotic analysis. *Conference on learning theory*, pages 300–323. PMLR, 2018.
- 212 A. S. Dalalyan and A. Karagulyan. User-friendly guarantees for the langevin monte carlo with  
213 inaccurate gradients. *Stochastic Processes and their Applications*, 129(12):5278–5311, 2019.
- 214 C. De Sa, M. Leszczynski, J. Zhang, A. Marzoev, C. R. Aberger, K. Olukotun, and C. H. High-  
215 accuracy low-precision training. *arXiv preprint arXiv:1803.03383*, 2018.
- 216 X. Gao, M. Gürbüzbalaban, and L. Zhu. Global convergence of stochastic gradient hamiltonian  
217 monte carlo for nonconvex stochastic optimization: Nonasymptotic performance bounds and  
218 momentum-based acceleration. *Operations Research*, 70(5):2931–2947, 2022.
- 219 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *Proceedings*  
220 *of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 221 Z. Li and C. M. De Sa. Dimension-free bounds for low-precision training. *Advances in Neural*  
222 *Information Processing Systems*, 32, 2019.
- 223 P.-C. Lin, M.-K. Sun, C. Kung, and T.-D. Chiueh. Floatsd: A new weight representation and as-  
224 sociated update method for efficient convolutional neural network training. *IEEE Journal on*  
225 *Emerging and Selected Topics in Circuits and Systems*, 9(2):267–279, 2019.
- 226 Y. Liu, Y. Gao, and W. Yin. An improved analysis of stochastic gradient descent with momentum.  
227 *Advances in Neural Information Processing Systems*, 33:18261–18271, 2020.
- 228 A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by  
229 generative pre-training. 2018.
- 230 M. Raginsky, A. Rakhlin, and M. Telgarsky. Non-convex learning via stochastic gradient langevin  
231 dynamics: a nonasymptotic analysis. *Conference on Learning Theory*, pages 1674–1703.  
232 PMLR, 2017.
- 233 K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recogni-  
234 tion. *arXiv preprint arXiv:1409.1556*, 2014.
- 235 Z. Song, Z. Liu, and D. Wang. Computation error analysis of block floating point arithmetic ori-  
236 ented convolution neural network accelerator design. *Proceedings of the AAAI Conference on*  
237 *Artificial Intelligence*, volume 32, 2018.
- 238 X. Sun, J. Choi, C.-Y. Chen, N. Wang, S. Venkataramani, V. V. Srinivasan, X. Cui, W. Zhang, and  
239 K. Gopalakrishnan. Hybrid 8-bit floating point (hfp8) training and inference for deep neural  
240 networks. *Advances in neural information processing systems*, 32, 2019.
- 241 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, . Kaiser, and I. Polo-  
242 sukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- 243 N. Wang, J. Choi, D. Brand, C.-Y. Chen, and K. Gopalakrishnan. Training deep neural networks  
244 with 8-bit floating point numbers. *Advances in neural information processing systems*, 2018.
- 245 M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient langevin dynamics. In  
246 *ceedings of the 28th international conference on machine learning (ICML)*, pages 681–688,  
247 2011.
- 248 M. Wortsman, T. Dettmers, L. Zettlemoyer, A. Morcos, A. Farhadi, and L. Schmidt. Stable and  
249 low-precision training for large-scale vision-language models. *arXiv preprint arXiv:2304.13013*,  
250 2023.
- 251 S. Wu, G. Li, F. Chen, and L. Shi. Training and inference with integers in deep neural networks.  
252 *arXiv preprint arXiv:1802.04680*, 2018.
- 253 R. Zhang, A. G. Wilson, and C. De Sa. Low-precision stochastic gradient langevin dynamics. In  
254 *International Conference on Machine Learning*, pages 26624–26644. PMLR, 2022.
- 255 T. Zhang, Z. Lin, G. Yang, and C. De Sa. Qpytorch: A low-precision arithmetic simulation frame-  
256 work. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-*  
257 *NeurIPS Edition (EMC2-NIPS)*, pages 10–13. IEEE, 2019.
- 258 D. Zou, P. Xu, and Q. Gu. Stochastic gradient hamiltonian monte carlo methods with recursive  
259 variance reduction. *Advances in Neural Information Processing Systems*, 2019.



260 **A Additional Results for Low-precision Stochastic Gradient Hamiltonian**  
 261 **Monte Carlo**

262 The underdamped Langevin dynamics has a continuous-time diffusion form:

$$\begin{aligned} dv_t &= -v_t dt - \nabla U(x_t) dt + \sqrt{\frac{2}{\mu}} dB_t \\ dx_t &= v_t dt \end{aligned} \tag{9}$$

263 And we formally define the stochastic rounding quantization function as:

$$Q^s(\cdot) = \begin{cases} \lfloor \cdot \rfloor & ; \text{ w.p. } \frac{1}{2} \\ \lceil \cdot \rceil & ; \text{ w.p. } \frac{1}{2} \end{cases} \tag{10}$$

264 Before diving into the theorems, we introduce some necessary assumptions.

265 **Assumption 1 (Smoothness)** The energy function  $U$  is  $M$ -smooth, i.e., there exists a positive constant  $M$  such that

$$\| \nabla U(x) - \nabla U(y) \| \leq M \| x - y \|; \text{ for any } x, y \in \mathbb{R}^d;$$

267  
 268 **Assumption 2 (Strongly Log-Convex)** The energy function  $U$  is  $m$ -strongly log-convex, i.e., there exists a positive constant  $m$  such that,

$$U(y) \leq U(x) + \langle \nabla U(x), y - x \rangle + \frac{m_1}{2} \| y - x \|^2; \text{ for any } x, y \in \mathbb{R}^d;$$

270  
 271 **Assumption 3 (Dissaptiveness)** There exist constants  $m_2, b > 0$ , such that the following holds

$$\langle \nabla U(x), x \rangle \geq m_2 \| x \|^2 - b; \text{ for any } x \in \mathbb{R}^d.$$

272  
 273 **Assumption 4 (Bounded Variance)** There exists a constant  $\ell > 0$ , such that the following holds

$$\mathbb{E} \| \nabla U(x) - \nabla U(x) \|^2 \leq \ell^2; \text{ for any } x \in \mathbb{R}^d;$$

274  
 275 Beyond the above assumptions, we further define  $m = M = m_1$  and  $\mu = M = m_2$  as the condition  
 276 number for strongly log-concave and non-log-concave target distribution respectively, and denote the  
 277 global minimum of  $U(x)$  as  $x^*$ . Assumption 3 is the standard assumption [Raginsky et al., 2017,  
 278 Zou et al., 2019, Gao et al., 2022] in the analysis of sampling from non-log-concave distributions and  
 279 is essential to guarantee the convergence of underdamped Langevin dynamics. Now we introduce  
 280 the of SGHMCLP-F for strongly log-concave and non-log-concave target distribution in Theorem 4  
 281 and 5 respectively.

282 **Theorem 4.** Suppose Assumptions 1, 2 and 4 hold and the minimum satisfies  $\ell^2 < D^2$ . Fur-  
 283 thermore, let  $p$  denote the target distribution  $\pi$  and  $v$ . Given any sufficiently small, if we set  
 284 the step size to be

$$\epsilon = \min \left( \frac{1}{4792325(d=m_1 + D^2)}; \frac{\mu^2}{1440 \mu^2 (M^2 + 1) \frac{2d}{4} + \mu^2} \right);$$

285 then after  $K$  steps starting with initial points  $x_0 = v_0 = 0$ , the output  $(x_K; v_K)$  of the SGHMCLP-  
 286 F in (3) satisfies

$$W_2(p(x_K; v_K); \pi) = \mathcal{O}(\epsilon);$$

287 for some  $K$  satisfying

$$K \geq \frac{1}{\epsilon} \log \frac{1}{\epsilon} \frac{36}{m_1} + D^2 \frac{1}{\epsilon} A = \mathcal{O} \left( \frac{1}{\epsilon^2} \log \frac{1}{\epsilon} \right);$$

288 Theorem 5. Suppose Assumptions 1, 3 and 4 hold. Furthermore, let  $p$  denote the target distribution  
 289 of  $x$  and  $v$ . Given initialization  $x_0 = v_0 = 0$  and  $\sigma^2 = 4Mu$ , for any sufficiently small  $\epsilon$ , if we set  
 290 the step size to be  $\epsilon = \mathcal{O}\left(\frac{\sigma^2}{\log(1/\epsilon)}\right)$  and also satisfy

$$\min \left( \frac{\sigma^2}{4(8Mu + u + 22\sigma^2)}; \frac{4u^2}{4Mu + 3\sigma^2}; \frac{6bu}{(4Mu + 3\sigma^2)d}; \frac{1}{8}; \frac{m_2}{12(21u + \sigma^2)M^2}; \frac{8(\sigma^2 + 2u)}{(20u + \sigma^2)} \right);$$

291 then, the  $K$ -th iteration of the SGHMCLP-F update (3), i.e.,  $x_K$  and  $v_K$ , satisfies

$$W_2(p(x_K; v_K); p) = \mathcal{O}\left(\frac{\sigma^2}{\epsilon} + A \log \frac{1}{\epsilon}\right);$$

292 for some  $K$  satisfying

$$K = \mathcal{O}\left(\frac{1}{\epsilon^2} \log^2 \frac{1}{\epsilon}\right);$$

293 where constants are defined as  $A = \max\left\{\frac{p}{2d + \sigma^2}, \frac{D}{2d + \sigma^2}\right\}$ ; and  $\epsilon$  is a constant w.r.t.  
 294 dimension  $d$ , denoting the concentration rate of the underdamped Langevin dynamics [Zou et al.,  
 295 2019].

296 Theorem 1 in Zhang et al. [2022] implies that for strongly log-concave target distribution, the  
 297 low-precision SGLD with full-precision gradient accumulators can achieve  $\epsilon$ -accuracy within  
 298  $\mathcal{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$  iterations.

299 Thus, the theorem of SGHMCLP-F does not showcase any advantage over SGLDLP-F. This is not  
 300 surprising, since the quantization applied to the gradients in the full-precision gradient accumulator  
 301 algorithm is equivalent to adding extra noise to the stochastic gradients. As theoretically shown by  
 302 Cheng et al. [2018] for strongly-log-concave target distribution, HMC doesn't exhibit any advantage  
 303 over the unadjusted Langevin algorithm when stochastic gradients are used.

304 However, as shown in the Theorem 5, for non-log-concave distributions, the low-precision SGHMC  
 305 displays faster convergence speed and a better dependence on the quantization compared to  
 306 SGLD. Besides the discussion in Theorem 1, we can discuss the upper bound due to the fact  
 307 that  $\log(x) \leq x^{1-\epsilon}$ , one can tune the choice of  $\epsilon$  and  $\sigma$ , and achieve  $\mathcal{O}\left(\frac{1}{\epsilon^{(1+2\epsilon)}}\right)$  2-Wasserstein  
 308 bound for non-log-concave target distribution. Furthermore, based on Theorem 10, after carefully  
 309 choosing the stepsize, the 2-Wasserstein distance of the SGLDLP-F algorithm can be further  
 310 bounded by  $\mathcal{O}\left(\frac{1}{\epsilon^{(2+2\epsilon)}}\right)$  which is worse than the bound  $\mathcal{O}\left(\frac{1}{\epsilon^{(1+2\epsilon)}}\right)$  obtained by SGHMC.  
 311 Next, we introduce the convergence analysis of SGHMCLP-L for strongly log-concave and non-  
 312 log-concave target distribution in Theorem 6 and 7 respectively.

313 Theorem 6. Let Assumption 1, 2 and 4 hold and the minimum satisfies  $\sigma^2 < D^2$ . Furthermore,  
 314 let  $p$  denote the target distribution of  $x$ . Given any sufficiently small  $\epsilon$ , if we set the step size  
 315 to be

$$\epsilon = \min \left( \frac{1}{6635525 \frac{d}{m_1} + D^2}; \frac{1}{2880 \frac{1}{u} \frac{d}{4} + \sigma^2} \right);$$

316 then after  $K$  steps starting with initial points  $x_0 = v_0 = 0$ , the output  $(x_K; v_K)$  of the SGHMCLP-L  
 317 in (4) satisfies

$$W_2(p(x_K; v_K); p) = \mathcal{O}\left(\frac{\sigma^2}{\epsilon} + A \log \frac{1}{\epsilon}\right); \quad (11)$$

318 for some  $K$  satisfying

$$K = \mathcal{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right) A = \mathcal{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right);$$

319 Compared with Theorem 2 in Zhang et al. [2022], We cannot show the advantages of low-precision  
 320 SGHMC over SGLD for strongly log-concave target distribution. However, for non-log-concave tar-  
 321 get distribution, we show SGHMCLP-L can achieve lower distance in smaller iterations. Next, we  
 322 present the convergence theorem of SGHMCLP-L for non-log-concave target distribution. Besides  
 323 the discussion in Theorem 2, by the same argument in Theorem 1's discussion after carefully choos-  
 324 ing the stepsize, the 2-Wasserstein distance of SGHMCLP-L to non-log-concave target distribution  
 325 can be further bounded as  $\mathcal{O}(e^{-(3+6\epsilon)})$ , and the distance of the sample obtained by SGLDLP-L  
 326 can be bounded as  $\mathcal{O}(e^{-10(1+\epsilon)})$ . Thus the low-precision SGHMC is more robust to the quan-  
 327 tization error than SGLD. Next, we present the convergence analysis of VC SGHMCLP-L in (8).  
 328 We begin with the formal definition of the variance-corrected quantization function. Instead of  
 329 adding real value Gaussian noise and quantizing the weights, we can design a categorical sampler  
 330 that samples from the space  $\mathcal{S}$ ;  $\mathbb{P}$  with the desired expectation and variance as

$$\text{Cat}(\mathbb{P}; v) = \begin{cases} \geq \frac{v+2}{2^2} & ; \quad w: \frac{v+2}{2^2} \\ > \frac{v+2}{2^2} & ; \quad w: \frac{v+2}{2^2} \\ > 0; & \text{otherwise} \end{cases} \quad (12)$$

331 Based on the sampler 12, we design the variance correction quantization function in the algo-  
 332 rithm 1.

333 Theorem 7. Let Assumptions 1, 3 and 4 hold. If  $4Mu$  and we set the step size to be

334  $\mathcal{O}\left(\frac{2}{\log(1+\epsilon)}\right)$ , also satisfied

$$\min \left\{ \frac{s}{4(8Mu + u + 22^2)}; \frac{4u^2}{4Mu + 3^2}; \frac{6bu}{(4Mu + 3^2)d}; \frac{1}{8}; \frac{m_2}{12(21u + )M^2}; \frac{8(\epsilon^2 + 2u)}{(20u + )} \right\};$$

335 let  $p$  denote the target distribution  $\mathbb{P}(x; v)$  then after  $K$  steps starting at the initial point  $x_0 =$   
 336  $v_0 = 0$  the output  $(x_K; v_K)$  of SGHMCLP-L in 4 satisfies

$$W_2(p(x_K; v_K); p) = \mathcal{O} \left( \frac{s}{\max\{ \epsilon^2; \log \frac{1}{\epsilon} + \frac{\log^{3-2\epsilon} 1}{2} p \} } \right); \quad (13)$$

337 for some  $K$  satisfying

$$K = \mathcal{O} \left( \frac{1}{2} \log^2 \frac{1}{\epsilon} \right);$$

338 Theorem 8. Let Assumption 1, 2 and 4 hold and the minimum satisfies  $\epsilon^2 < D^2$ . Furthermore,  
 339 let  $p$  denote the target distribution of  $x$  and  $v$ . Given any sufficiently small, if we set the stepsize  
 340 to be

$$= \min \left\{ \frac{8}{6635525 \frac{d}{m_1} + D^2}; \frac{9}{90u^2 2d + 360u^2 2^2}; \frac{1}{1} \right\};$$

341 after  $K$  steps starting from the initial point  $x_0 = v_0 = 0$  the output  $(x_K; v_K)$  of the VC SGHMCLP-  
 342 L in algorithm 2 satisfies

$$W_2(p(x_K; v_K); p) = \mathcal{O} \left( \frac{1}{\epsilon} \right); \quad (14)$$

343 for some  $K$  satisfying

$$K = \frac{1}{\epsilon} \log \left( \frac{36 \frac{d}{m_1} + D^2}{\epsilon} \right) = \mathcal{O} \left( \frac{1}{\epsilon^2} \log \frac{1}{\epsilon} \right);$$

344 Theorem 8 shows that the variance corrected quantization function can solve the overdispersion  
 345 problem we observe for the VC SGHMCLP-L algorithm for strongly log-concave distribution.  
 346 The  $W_2$  distance between the sample distribution and target distribution can be arbitrarily close  
 347 to  $\mathcal{O}(\epsilon)$ . Compared to the Theorem 3 in Zhang et al. [2022], the VC SGHMCLP-L doesn't  
 348 showcase its advantage over VC SGLDLP-L for strongly log-concave distribution, however for  
 349 non-log-concave target distribution we show VC SGHMCLP-L can achieve lower Wasserstein  
 350 distance in smaller iterations. Next, we provide the convergence analysis of the VC SGHMCLP-L  
 351 for non-log-concave distribution.

---

**Algorithm 1 Variance-Corrected Quantization Function  $Q^{VC}$ .**


---

input:  $(x, v, \sigma)$   $Q^{VC}$  returns a variable with mean  $x$  and variance  $v$   
 $v_0 = 4\sigma^2$   $v_0$  is the largest possible variance that stochastic rounding can cause  
 if  $v > v_0$  then  $v = v_0$  add a small Gaussian noise and sample from the discrete grid to make up the remaining variance  
 $x = x + \sqrt{v - v_0} N(0; I_d)$   
 $r = x - Q^d(x)$   
 for all  $i$  do  
   sample  $c_i$  from  $\text{Ca}(r_i; v_0)$  as in (12)  
 end for  
 $Q^d(x) + \text{sign}(r) \cdot c$   
 else  $v \leq v_0$  sample from the discrete grid to achieve the target variance  
 $r = Q^s(x)$   
 for all  $i$  do  
    $v_s = 1 - \frac{|r_i|}{d}$   $v_s = r_i^2 + \frac{|r_i|}{d} (|r_i| + \text{sign}(r_i))^2$   
   if  $v > v_s$  then  
     sample  $c_i$  from  $\text{Ca}(0; v - v_s)$  as in (12)  
      $r_i = Q^s(x)_i + c_i$   
   else  
      $r_i = Q^s(x)_i$   
   end if  
 end for  
 end if  
 clip if outside representable range  
 return

---

352 Theorem 9. Let Assumption 1, 3 and 4 hold. If  $\frac{1}{\epsilon} \geq 4Mu$  and we set the step size to be  
 353  $\epsilon = \frac{2}{\log(1/\epsilon)}$ , also satisfied

$$\min \left( \frac{s}{4(8Mu + u + 22\epsilon^2)}; \frac{4u^2}{4Mu + 3\epsilon^2}; \frac{6bu}{(4Mu + 3\epsilon^2)d}; \frac{1}{8}; \frac{m_2}{12(21u + \epsilon^2)M^2}; \frac{8(\epsilon^2 + 2u)}{(20u + \epsilon^2)} \right)$$

354 We further assume that  $Q_G(r - U(x)) \leq G^2$ , let  $p$  be the target distribution of  $x$  then after  
 355  $K$  steps starting at the initial point  $x_0 = v_0 = 0$  the output  $(x_K)$  of the VC SGHMCLP-L in  
 356 algorithm 2 satisfies

$$W_2(p(x_K); p) = \mathcal{O} \left( \frac{s}{\max\{2; \log \frac{1}{\epsilon}\} + \frac{\log \frac{1}{p}}{\epsilon}} \right); \quad (15)$$

357 for some  $K$  satisfying

$$K = \mathcal{O} \left( \frac{1}{2\epsilon^2} \log^2 \frac{1}{\epsilon} \right);$$

## 358 B Stochastic Gradient Langevin Dynamics Result

359 In order to sample from the target distribution, Langevin dynamics-based samplers, such as over-  
 360 damped Langevin MCMC and underdamped Langevin MCMC methods, are widely used when  
 361 the evaluation of  $U(x)$  is expensive due to a large sample size. The continuous-time overdamped  
 362 Langevin MCMC can be represented by the following stochastic differential equation(SDE):

$$dx_t = -r \cdot U(x_t) + \sqrt{P} dB_t; \quad (16)$$

363 where  $B_t$  represents the standard Brownian motion. Under some mild conditions, it can  
 364 be proved that the invariant distribution of (16) converges to the target distribution  $p(x) \propto \exp(-U(x))$ . To

Table 3: Theoretical results of the achieved Wasserstein distance and the required gradient complexity for both log-concave (italic) non-log-concave (bold) target distributions, where  $\epsilon$  is any sufficiently small constant,  $\delta$  is the quantization error, and  $\beta$  and  $\beta'$  denote the concentration rate of underdamped and overdamped Langevin dynamics respectively.

	Gradient Complexity	Achieved Wasserstein
Full-precision gradient accumulators		
SGLD/SGHMC (Theorem 4)	$\mathcal{O}(\log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon})$
SGLD (Theorem 10)	$\mathcal{O}(\frac{1}{\epsilon^4} \log^5 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \log \frac{1}{\epsilon})$
SGHMC (Theorem 5)	$\mathcal{O}(\frac{1}{\epsilon^2} \log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \frac{1}{\log(\frac{1}{\epsilon})})$
Low-precision gradient accumulators		
SGLD/SGHMC (Theorem 6)	$\mathcal{O}(\log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \frac{1}{\epsilon})$
VC SGLD/VC SGHMC (Theorem 8)	$\mathcal{O}(\log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \frac{1}{\epsilon})$
SGLD (Theorem 11)	$\mathcal{O}(\frac{1}{\epsilon^4} \log^5 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \log^5 \frac{1}{\epsilon})$
VC SGLD (Theorem 12)	$\mathcal{O}(\frac{1}{\epsilon^4} \log^3 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \log^3 \frac{1}{\epsilon})$
SGHMC (Theorem 7)	$\mathcal{O}(\frac{1}{\epsilon^2} \log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \log^{3=2} \frac{1}{\epsilon})$
VC SGHMC (Theorem 9)	$\mathcal{O}(\frac{1}{\epsilon^2} \log^2 \frac{1}{\epsilon})$	$\mathcal{O}(\frac{\delta}{\epsilon} + \log \frac{1}{\epsilon})$

365 reduce the computational cost of evaluating  $\nabla U(x)$ , Welling and Teh [2011] proposed the Stochastic  
366 Gradient Langevin Dynamics (SGLD) and updates the weights using stochastic gradients:

$$x_{k+1} = x_k - \eta \nabla U(x_k) + \sqrt{\frac{\eta}{2}} \bar{\epsilon}_{k+1}; \quad (17)$$

367 where  $\eta$  is the stepsize,  $\bar{\epsilon}_{k+1}$  is a standard Gaussian noise, and  $\nabla U(x_k)$  is an unbiased estimation  
368 of  $\nabla U(x_k)$ . Despite the additional noise induced by stochastic gradient estimations, SGLD can still  
369 converge to the target distribution.

370 The low-precision SGLD with full-precision gradient accumulators (SGLDLP-F) only quantizes  
371 weights before computing the gradient. The update rule can be defined as:

$$x_{k+1} = x_k - \eta \nabla U(Q_W(x_k)) + \sqrt{\frac{\eta}{2}} \bar{\epsilon}_{k+1}; \quad (18)$$

372 Zhang et al. [2022] shows that the SGLDLP-F outperforms its counterpart low-precision SGD with  
373 full-gradient accumulators (SGDLP-F). The computation costs can be further reduced using low-  
374 precision gradient accumulators by only keeping low-precision weights. Low-precision SGLD with  
375 low-precision gradient accumulators (SGLDLP-L) can be defined as the following:

$$x_{k+1} = Q_W(x_k) - \eta \nabla U(x_k) + \sqrt{\frac{\eta}{2}} \bar{\epsilon}_{k+1}; \quad (19)$$

376 Zhang et al. [2022] studied the convergence property of both SGLDLP-F and SGLDLP-L under  
377 strongly-log-concave distributions, and showed that a small stepsize deteriorates the performance of  
378 SGLDLP-L. To mitigate this problem, Zhang et al. [2022] proposed a variance-corrected quantiza-  
379 tion function.

380 Theorem 10. Suppose Assumptions 1, 3 and 4 hold. Let  $\beta$  have the same definition in Theorem 5,  
381 and  $\beta'$  be the concentration number of (16). After  $K$  steps starting with initial point  $x_0 = 0$ , if we  
382 set the stepsize to be  $\eta = \frac{\delta}{\log(\frac{1}{\epsilon})^4}$ . The output  $x_K$  of SGLDLP-F in (18) satisfies

$$W_2(p(x_K); p) \leq \frac{\delta}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\epsilon}; \quad (20)$$

383 provided

$$K = \mathcal{O}\left(\frac{1}{\delta} \log^5 \frac{1}{\epsilon}\right);$$

---

**Algorithm 2 Variance-Corrected Low-Precision SGHMC (VC SGLDLP-L).**


---

given: Step size  $\eta$ , friction  $\gamma$ , inverse mass  $m$ , number of training iterations  $K$ , gradient quantizer  $Q_G$ , quantization gap  $u$  and upper bound of low-precision representation  $U$ . Let  $\text{Var}_v^{\text{hmc}} = u(1 - e^{-2})$  and  $\text{Var}_x^{\text{hmc}} = u^2(2 + 4e^{-\gamma} - e^{-2} - 3)$  and  $S_v = 1$ . **Initialize the scaling parameters.**  
 for  $k = 1 : K$  do  
   rescale  $v_k = v_k \cdot S_v$  **Restore the velocity before update**  
   update  $(v_{k+1}) = v_k e^{-\gamma} + \eta(1 - e^{-\gamma})Q_G(r \nabla(x_k))$   
   update  $(x_{k+1}) = x_k + \eta(1 - e^{-\gamma})v_k + u^2(1 + e^{-\gamma})Q_G(r \nabla(x_k))$   
   update  $S_v = \frac{k \cdot (v_{k+1})_{k_1}}{U}$  **Update the Scaling**  
   update  $v_{k+1} = Q^{\text{vc}}(v_{k+1}); \text{Var}_v^{\text{hmc}} = S_v^2$ ;  
   update  $x_{k+1} = Q^{\text{vc}}(x_{k+1}); \text{Var}_x^{\text{hmc}} = S_v^2$ ;  
 end for  
 output: samples  $x_k$

---

384 Theorem 10 shows that the low-precision SGLD with full-precision gradient accumulators can converge to the non-log-concave target distribution provided a small gradient variance and quantization error. Next, we present the SGLDLP-L's result.

387 Theorem 11. Let Assumptions 1, 3 and 4 hold. If we set the step size to be  $\Theta\left(\frac{1}{\log(1/\epsilon)}\right)^4$ , after  $K$  steps starting at the initial point  $x_0 = 0$  the output  $x_K$  of the SGLDLP-L in (19) satisfies

$$W_2(p(x_K); p) = \Theta\left(\frac{1}{\log(1/\epsilon)}\right)^4 + \frac{p}{\max f^2}; \text{glog} \frac{1}{\epsilon} + \frac{\log^5 \frac{1}{\epsilon} p}{4}; \quad (21)$$

389 provided

$$K = \Theta\left(\frac{1}{\log(1/\epsilon)}\right)^4 \log^5 \frac{1}{\epsilon} :$$

390 The VC SGLDLP-L can be done as:

$$x_{k+1} = Q^{\text{vc}}(x_k + \eta Q_G(r \nabla(x_k))); \quad (22)$$

391 Theorem 12. Let Assumption 1, 3 and 4 hold. If we set the step size to be  $\Theta\left(\frac{1}{\log^4(1/\epsilon)}\right)$ , after  $K$  steps from the initial point  $x_0 = 0$  the output  $x_K$  of VC SGLDLP-L in (22) satisfies

$$W_2(p(x_K); p) = \Theta\left(\frac{1}{\log^4(1/\epsilon)}\right)^4 + \frac{p}{\max f^2}; \text{glog} \frac{1}{\epsilon} + \frac{\log^3 \frac{1}{\epsilon} p}{2}; \quad (23)$$

393 provided

$$K = \Theta\left(\frac{1}{\log^4(1/\epsilon)}\right)^4 \log^5 \frac{1}{\epsilon} :$$

## 394 C Technical Detail

395 In this section, we disclose more details of empirical experiments. When implementing low-precision SGHMC on classification task in the CIFAR-10 and CIFAR-100 dataset, we observed that the momentum terms tend to gather in a small range around zero in which case the low-precision representations of end up in gathering only few points, thus the momentum information is seriously lost and cause in performance degradation. In order to tackle this problem and fully utilize all the low-precision representations, we borrow the idea of rescaling from the bit centering trick and adopted to the low-precision SGHMC method. The detailed algorithm is listed in Algorithm 2.

402 In Algorithm 2, we introduce the bit centering trick [De Sa et al., 2018] to enhance the variance corrected quantization function. Bit centering trick is a technique to increase the accuracy low-precision training algorithm by recentering and rescaling representable bits making low-precision

405 numbers closer to its real full-precision counterpart. We borrow the idea of rescaling to enhance  
 406 the variance-corrected quantization function. Based on the discussion in previous paragraph, when  
 407 the desired variance is small the variance corrected quantization has a high chance to match the  
 408 variance. By scaling up the weights, additional to increasing the accuracy of low-precision repre-  
 409 sentation also increase the desired variance resulting in a lower chance of fail in variance corrected  
 410 quantization.

## 411 D Proof of Main Theorems

### 412 D.1 Proof of Theorem 4

413 Section 3.1 introduces low-precision HMC with full-precision gradient accumulators (SGHMCLP-  
 414 F) as:

$$\begin{aligned} v_{k+1} &= v_k e^{-\eta u} (1 - e^{-\eta u}) Q_G(r \psi(Q_W(x_k))) + \frac{v}{k} \\ v_{k+1} &= x_k + \eta (1 - e^{-\eta u}) v_k + u^{-2} (\eta + e^{-\eta u}) Q_G(r \psi(Q_W(x_k))) + \frac{x}{k}; \end{aligned}$$

415 In this section, we prove the convergence of SGHMCLP-F in terms of Wasserstein distance for  
 416 strongly-log-concave target distribution via coupling argument. To simplify the notation we define  
 417 the quantized stochastic gradients as:

$$g(x) := Q_G(r \psi(Q_W(x))) \quad (24)$$

$$=: r U(x) + \xi \quad (25)$$

418 Lemma 13. For any  $x \in \mathbb{R}^d$ , the random noise of the low-precision gradients defined in (25)  
 419 satisfies:

$$\begin{aligned} kE[k^2] &= M^2 \frac{2d}{4} \\ E[k^2] &= (M^2 + 1) \frac{2d}{4} + \sigma^2. \end{aligned}$$

420

421 We follow the proof in Cheng et al. [2018]. Denote  $\mathcal{B}(\mathbb{R}^d)$  the Borel  $\sigma$ -field of  $\mathbb{R}^d$ . Given  
 422 probability measures  $\mu$  and  $\nu$  on  $(\mathbb{R}^d; \mathcal{B}(\mathbb{R}^d))$ , we define a transference plan between  $\mu$  and  $\nu$  as  
 423 a probability measure  $\alpha$  on  $(\mathbb{R}^d \times \mathbb{R}^d; \mathcal{B}(\mathbb{R}^d \times \mathbb{R}^d))$  such that for all sets  $A \in \mathcal{B}(\mathbb{R}^d)$ ,  $(A \times \mathbb{R}^d) = \alpha$   
 424 and  $(\mathbb{R}^d \times A) = \alpha$ . We denote  $\Pi(\mu, \nu)$  as the set of all transference plans. A pair of random  
 425 variables  $(x; y)$  is called a coupling if there exists a  $\alpha \in \Pi(\mu, \nu)$  such that  $(x; y)$  is distributed  
 426 according to  $\alpha$ . (With some abuse of notation, we will also refer to  $\alpha$  as the coupling.)

427 In order to calculate the Wasserstein distance from the proposed sample  $(x_k; v_k)$  and the target  
 428 distribution sample  $(x; v)$ , we define sample  $\alpha_k = (x_k; x_k + v_k)$  and the target distribution  
 429 sample  $\alpha = (x; x + v)$ . Let  $p_k = (x_k; v_k)$  and  $\mathfrak{b}$  be the operator that maps from  $p_k$  to  $p_{k+1}$   
 430 i.e.

$$p_{k+1} = \mathfrak{b} p_k:$$

431 The solution  $(x_t; v_t)$  of the continuous underdamped Langevin dynamics with exact gradient satis-  
 432 es the following equations:

$$\begin{aligned} v_t &= v_0 e^{-\eta t} - \eta \int_0^t e^{-\eta(t-s)} r U(x_s) ds + \sqrt{\frac{\eta}{2u}} \int_0^t e^{-\eta(t-s)} dB_s; \quad (26) \\ x_t &= x_0 + \int_0^t v_s ds; \end{aligned}$$

433 Let  $\mathfrak{a}$  denote the operator that maps to the solution of continuous underdamped Langevin dy-  
 434 namics in (26) after time step  $t$ . Notice the solution  $(v_t; x_t)$  of the discrete underdamped Langevin  
 435 dynamics with an exact gradient can be written as

$$\begin{aligned} v_t &= v_0 e^{-\eta t} - \eta \int_0^t e^{-\eta(t-s)} r U(x_0) ds + \sqrt{\frac{\eta}{2u}} \int_0^t e^{-\eta(t-s)} dB_s; \quad (27) \\ x_t &= x_0 + \int_0^t v_s ds; \end{aligned}$$

436 We can also define a similar operator for the discrete underdamped Langevin dynamics solution  
 437  $p_t = (x_t; v_t)$ , let  $e_t$  be the operator that maps  $p_0$  to  $p_t$ . Furthermore the SGHMCLP-F can be  
 438 written as:

$$v_t = v_0 e^{-\gamma t} + \int_0^t e^{-\gamma(t-s)} g(x_0) ds + \sqrt{\frac{2\mu}{\gamma}} \int_0^t e^{-\gamma(t-s)} dB_s; \quad (28)$$

$$x_t = x_0 + \int_0^t v_s ds;$$

439 Given  $g(x_0) = \gamma U(x_0) + \mu$  and  $x_0 = x_0$ , we know:

$$v_t = v_0 e^{-\gamma t} + \int_0^t e^{-\gamma(t-s)} g(x_0) ds \quad (29)$$

$$x_t = x_0 + \int_0^t \int_0^s e^{-\gamma(s-r)} g(x_0) ds dr :$$

440 Lemma 14. Let  $q_0$  be some initial distribution and  $e_t$  and  $e$  be the operator we defined above for  
 441 discrete Langevin dynamics with exact full-precision gradients and low-precision gradients respec-  
 442 tively. If the step size  $\epsilon > 0$ , then the Wasserstein distance satisfies

$$W_2^2(e^{-\epsilon} q_0; q) \leq W_2^2(q_0; q) + \frac{\epsilon^2}{5} \frac{L^2}{2\mu} + 5\epsilon^2 \frac{d}{4} (M^2 + 1) + \epsilon^2 :$$

443 The lemma 14 says that if starting from the same distribution after one step of low-precision update  
 444 the Wasserstein distance from the target distribution is bounded by the distance after one step of  
 445 exact gradients plus  $\epsilon^2 \frac{L^2}{5}$ . Furthermore from the corollary 7 in Cheng et al. [2018] we know  
 446 that for any  $\epsilon > 0$ ;  $K \geq 1$ :

$$W_2^2(e^{-\epsilon} q; q) \leq e^{-2\epsilon} W_2^2(q; q); \quad (30)$$

447 where  $\epsilon = M\epsilon_1$  is the condition number. Let  $\epsilon_K$  denote the discretization error bound from Theorem 9 and Lemma 8 (sandwich inequality) in Cheng et al.  
 448 [2018], we get

$$W_2(e^{-\epsilon} q; e^{-\epsilon} q) \leq 2W_2(p_1; e^{-\epsilon} p_1) \leq \frac{\epsilon_K}{5} :$$

450 By triangle inequality:

$$W_2(e^{-\epsilon} q; q) \leq W_2(e^{-\epsilon} q; e^{-\epsilon} q) + W_2(e^{-\epsilon} q; q) \leq \frac{\epsilon_K}{5} + e^{-2\epsilon} W_2(q; q):$$

451 Combine this with the result in Lemma 14 we have,

$$W_2^2(e^{-\epsilon} q; q) \leq e^{-2\epsilon} W_2^2(q; q) + \frac{\epsilon^2}{5} \frac{L^2}{2\mu} + 5\epsilon^2 \frac{d}{4} (M^2 + 1) + \epsilon^2 :$$

452 By invoking the Lemma 7 in Dalalyan and Karagulyan [2019] we can bound the Wasserstein  
 453 distance by:

$$W_2(q_K; q) \leq e^{-K\epsilon} W_2(q_0; q) + \frac{\epsilon^2 \frac{L^2}{5} + \frac{\epsilon M^2}{2} \frac{d}{4}}{1 - e^{-2\epsilon}} + \frac{5\epsilon^2 \frac{d}{4} (M^2 + 1) + \epsilon^2}{1 - e^{-2\epsilon}} :$$

454 Finally by sandwich inequality we have:

$$W_2(p_K; p) \leq 4e^{-K\epsilon} W_2(p_0; p) + 4 \frac{\epsilon^2 \frac{L^2}{5} + \frac{\epsilon M^2}{2} \frac{d}{4}}{1 - e^{-2\epsilon}} + \frac{20\epsilon^2 \frac{d}{4} (M^2 + 1) + \epsilon^2}{1 - e^{-2\epsilon}} :$$



455 Now we let the first term less than  $\frac{1}{3}$ , from the lemma 13 in [Cheng et al., 2018] we know that  
 456  $W_2(\rho_K; p) \leq 3 \frac{d}{m_1} + D^2$ . So we can choose  $K$  as the following,

$$K \geq \frac{2}{\epsilon} \log \frac{1}{3\epsilon} \frac{d}{m_1} + D^2 :$$

457 Next, we choose a stepsize  $p = \frac{1}{479232=5(d=m_1+D^2)}$  to ensure the second term is controlled below

458  $= 3 + \frac{16}{2} \frac{uM}{5} \frac{p^{5d}}{2}$ . Since  $e^{-2} = \frac{1}{4}$  and definition of  $E_K$ ,

$$\begin{aligned} & 4 \frac{2 \frac{q}{2} \frac{8E_K}{5} + \frac{uM}{2} \frac{p^{5d}}{2}}{1 - e^{-2}} = 4 \frac{2 \frac{q}{2} \frac{8E_K}{5} + \frac{uM}{2} \frac{p^{5d}}{2}}{=4} = 16 \frac{1}{1} \left( \frac{8E_K}{5} + \frac{uM}{2} \frac{p^{5d}}{2} \right) \\ & = 3 + \frac{16}{2} \frac{uM}{5} \frac{p^{5d}}{2} : \end{aligned}$$

459 Finally by choosing the stepsize satisfied that,

$$\frac{M \frac{p^{5d}}{2}}{120u (M^2 + 1) \frac{2d}{4} + 2} ;$$

460 the third term can be bounded as:

$$\begin{aligned} & \frac{20u^2 (M^2 + 1) \frac{2d}{4} + 2}{2 \frac{q}{2} \frac{8E_K}{5} + \frac{uM}{2} \frac{p^{5d}}{2} + 1 - e^{-2}} = \frac{20u^2 (M^2 + 1) \frac{2d}{4} + 2}{5u^2 (M^2 + 1) \frac{2d}{4} + 2} \\ & \frac{20u^2 (M^2 + 1) \frac{2d}{4} + 2}{\frac{uM}{2} \frac{p^{5d}}{2}} = 40u \frac{(M^2 + 1) \frac{2d}{4} + 2}{M \frac{p^{5d}}{2}} = 3: \end{aligned}$$

461 This complete the proof.

## 462 D.2 Proof of Theorem 5

463 In this section we analyze the Wasserstein distance between the sample in (3) and the  
 464 target distribution, given the target distribution satisfies the assumption 1 and 3. We follow the  
 465 proof in Raginsky et al. [2017]. To analyze the Wasserstein distance, we first calculate the distance  
 466 between solutions of low-precision discrete underdamped Langevin dynamics and solutions of the ideal  
 467 continuous underdamped Langevin dynamics, also the distance between solutions of the ideal  
 468 continuous underdamped Langevin dynamics and the target distribution.

469 Again let  $\rho_k = (x_k; v_k)$  denote the low-precision sample from (3) at iteration  $k$ , let  $\rho_t = (x_t; v_t)$   
 470 denote the sample from the ideal continuous underdamped Langevin dynamics in 26 iterations  
 471 the Wasserstein distance between  $\rho_k$  and the target distribution  $p$  can be bounded as:

$$W_2(\rho_k; p) \leq W_2(\rho_k; \rho_t) + W_2(\rho_t; p) :$$

472 We first bound  $W_2(\rho_k; \rho_t)$  by invoking the weighted CKP inequality Bolley and Villani [2005],

$$W_2^2(\rho_k; \rho_t) \leq \frac{q}{D_{KL}(\rho_k || \rho_t)} + \frac{q}{4 D_{KL}(\rho_k || \rho_t)} ;$$

473 where  $\frac{q}{2} = 2 \inf_{\gamma > 0} \frac{1}{\gamma} = 3 - 2 + \log E_{\rho_k} [\exp(-\gamma(x_k^2 + v_k^2))]^{-1}$ . We define a Lyapunov  
 474 function for every  $(x; v) \in \mathbb{R}^d \times \mathbb{R}^d$

$$E(x; v) = \gamma(x^2 + v^2) + 2v = \gamma(x^2 + 8u(U(x)) + U(x)) = 2 :$$

475 Note that  $\gamma(x^2 + v^2) + 2v = \gamma(x^2 + 8u(U(x)) + U(x)) = 2$ , we can have:

$$E(x; v) = \gamma(x^2 + v^2) + 2v = \gamma(x^2 + 8u(U(x)) + U(x)) = 2 :$$

476 Given assumptions 2 and 3 hold and apply Lemma B.4 in Zou et al. [2019], we can get

$$\frac{2}{S} \inf_{0 < \frac{m_2}{128u}; \frac{m_2}{32}g} \frac{1}{S} \frac{3}{2} + 2 E(X_0; V_0) + \frac{32M u (4d + 2b + m_2 k x^2)}{2m_2}$$

$$2 E(X_0; V_0) + \frac{32M u (4d + 2b + m_2 k x^2) + 16(12um_2 + 3^{-2})}{2m_2} := \quad :$$

477 It remains to bound the divergence between the distribution  $\mathbb{P}_K$  and  $\mathbb{P}_K$ . We first define a continuous  
 478 interpolation of the low-precision sample  $(x_k; v_k)$ ,

$$dv_t = v_t dt - uG_t dt + \frac{p}{2} \overline{u} dB_t \quad (31)$$

$$dx_t = v_t dt; \quad (32)$$

479 where  $G_t = \sum_{k=0}^{\lfloor t/h \rfloor} \mathfrak{g}(x_k) 1_{t \in [k, (k+1)h)}$ . Integrating this equation from time 0 to  $t$ , we can get

$$v_t = v_0 + \int_0^t v_s ds - \int_0^t uG_s ds + \frac{p}{2} \int_0^t \overline{u} dB_s$$

$$x_t = x_0 + \int_0^t v_s ds:$$

480 Notice that when  $t = k$ , the solution of (31) has the same distribution with the low-precision  
 481 sample  $(x_k; v_k)$ . Now by Girsanov formula we can compute the Radon-Nikodym derivative  
 482 with respect to  $\mathbb{P}_K$  as follow:

$$\frac{d\mathbb{P}_K}{d\mathbb{P}_K} = \exp \left( \int_0^t \frac{r}{2} \int_0^s U(x_s) G_s ds - \frac{u}{4} \int_0^t \int_0^s \text{tr} U(x_s) G_s ds \right) :$$

483 It follows that

$$D_{KL}(\mathbb{P}_K \parallel \mathbb{P}_K) = E_{\mathbb{P}_K} \log \frac{d\mathbb{P}_K}{d\mathbb{P}_K} \quad (33)$$

$$= \frac{u}{4} E \int_0^k \text{tr} U(x_s) G_s ds$$

$$= \frac{u}{4} \sum_{k=0}^{\lfloor T/h \rfloor} E \int_k^{(k+1)h} \text{tr} U(x_s) G_s ds$$

$$= \frac{u}{4} \sum_{k=0}^{\lfloor T/h \rfloor} E \int_k^{(k+1)h} \text{tr} U(x_s) \mathfrak{g}(x_k) ds:$$

484 Furthermore, in the  $k$ -th interval, we have

$$E \int_k^{(k+1)h} \text{tr} U(x_s) \mathfrak{g}(x_k) ds \leq 2E \int_k^{(k+1)h} \text{tr} U(x_k) ds + 2E \int_k^{(k+1)h} \text{tr} U(x_k) \mathfrak{g}(x_k) ds : \quad (34)$$

485 We now bound the first term in the RHS of the (34). By the smooth Assumption 1, we have

$$E \int_k^{(k+1)h} \text{tr} U(x_s) \mathfrak{g}(x_k) ds \leq M^2 E \int_k^{(k+1)h} \|x_s - x_k\|^2 ds :$$

486 Notice that

$$x_s = x_k + \int_k^s v_r dr$$

$$= x_k + \int_k^s v_k e^{-(r-k)} + \int_k^s \int_k^r v_k e^{-(r-z)} \mathfrak{g}(x_k) dz + \frac{p}{2} \int_k^s \int_k^r \overline{u} e^{-(r-z)} dB_z dr:$$

487 This further implies that:

$$\begin{aligned}
 \mathbb{E} \|x_k\|^2 &= \mathbb{E} \left[ \sum_{k=0}^s v_k e^{-(r_k)} u \sum_{k=0}^r e^{-(r_z)} g(x_k) dz + \rho \frac{Z_r}{2u} e^{-(r_z)} dB_z \right]^2 \\
 &= 3 \sum_{k=0}^s v_k e^{-(k-r)} dr + 3 \sum_{k=0}^s \sum_{k=0}^r u g(x_k) e^{-(z-r)} dz dr + 6 \rho \sum_{k=0}^s \sum_{k=0}^r e^{-(r_z)} dB_z dr \\
 &= 3 \sum_{k=0}^s v_k k^2 + 3 u^2 \sum_{k=0}^s k g(x_k) k^2 + 3 \frac{u}{2} \sum_{k=0}^s (s-k) + 4 e^{-(s-k)} e^{2(s-k)} \sum_{k=0}^s d \\
 &= 3 \sum_{k=0}^s v_k k^2 + u^2 \sum_{k=0}^s k g(x_k) k^2 + 2 du ; \tag{35}
 \end{aligned}$$

488 where we use inequality  $x e^{-x} \leq 1 - x + x^2/2$  for  $x > 0$  and  $k \leq s - (k+1)$  to get the  
 489 last inequality. Given this analysis we can bound the first term in the RHS of (34)

$$\mathbb{E} \sum_{k=0}^h \|U(x_s) - U(x_k)\|^2 \leq 3M^2 \sum_{k=0}^s \mathbb{E} v_k k^2 + u^2 \sum_{k=0}^s k g(x_k) k^2 + 2 du ;$$

490 By lemma 13, the second term in the RHS of (34) can be bounded as:

$$\mathbb{E} \sum_{k=0}^h \|U(x_k) - g(x_k)\|^2 \leq (M^2 + 1) \frac{2d}{4} + \sum_{k=0}^s ;$$

491 We need to introduce a lemma to bound  $\sup_k v_k k^2$ ,  $\sup_k v_k k^2$  and  $\sup_k k g(x_k) k^2$ .

492 Lemma 15. Under Assumptions 1 and 3, if we set the stepsize statis ed the following condition:

$$\min \left\{ \frac{s}{4(8Mu + u + 22)} ; \frac{4u^2}{4Mu + 3} ; \frac{6bu}{(4Mu + 3)d} ; \frac{1}{8} ; \frac{m_2}{12(21u + )M^2} ; \frac{8(\sum_{k=0}^s + 2u)}{(20u + )} \right\} ;$$

493 then for all  $\rho > 0$  the  $\mathbb{E} \|x_k\|^2$ ,  $\mathbb{E} \|v_k\|^2$  and  $\mathbb{E} \|k g(x_k)\|^2$  can be bounded as

$$\begin{aligned}
 \mathbb{E} \|x_k\|^2 &\leq \bar{E} + C_0 (M^2 + 1) \frac{2d}{4} + \sum_{k=0}^s \\
 \mathbb{E} \|v_k\|^2 &\leq 2\bar{E} + 2C_0 (M^2 + 1) \frac{2d}{4} + \sum_{k=0}^s \\
 \mathbb{E} \|k g(x_k)\|^2 &\leq 2 (M^2 + 1) \frac{2d}{4} + \sum_{k=0}^s + 4M^2 \bar{E} + 4G^2
 \end{aligned}$$

494 where  $\bar{E}$  and  $C_0$  are de ned as:

$$\begin{aligned}
 \bar{E} &= \mathbb{E} [E(x_0; v_0)] + \frac{24(21u + )uM}{m_2^3} G^2 + \frac{96(d+b)uM}{m_2^2} ; \quad G = \mathbb{E} \|U(0)\| \\
 C_0 &= \frac{96u \sum_{k=0}^s + 2u}{m_2^4} ;
 \end{aligned}$$

495 We now ready to bound  $E \int_{\mathbb{R}^d} U(x_s) \mathfrak{g}(x_k) k^2 dx$  as:

$$\begin{aligned}
 & E \int_{\mathbb{R}^d} U(x_s) \mathfrak{g}(x_k) k^2 dx = 2E \int_{\mathbb{R}^d} U(x_s) \int_{\mathbb{R}^d} U(x_k) k^2 dx + 2E \int_{\mathbb{R}^d} U(x_k) \mathfrak{g}(x_k) k^2 dx \\
 & \leq 6M^2 \int_{\mathbb{R}^d} E \int_{\mathbb{R}^d} k^2 dx + u^2 \int_{\mathbb{R}^d} E \int_{\mathbb{R}^d} \mathfrak{g}(x_k) k^2 dx + 2 \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx \\
 & \leq 6M^2 \int_{\mathbb{R}^d} (2 + 4M^2 u^2) \bar{E} dx + (2C_0 + 2 + 2u^2) \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx + 4u^2 \int_{\mathbb{R}^d} 2G^2 + 2 du \\
 & + 2 \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx \\
 & \leq 6M^2 \int_{\mathbb{R}^d} (2 + 4M^2 u^2) \bar{E} dx + 4u^2 \int_{\mathbb{R}^d} 2G^2 + 2 du \\
 & + 6M^2 \int_{\mathbb{R}^d} (2C_0 + 2 + 2u^2) dx + 2 \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx :
 \end{aligned}$$

496 Thus the divergence can be bounded as:

$$\begin{aligned}
 D_{KL}(\rho_K \| \rho_K) & \leq \frac{3u}{2} M^2 \int_{\mathbb{R}^d} (2 + 4M^2 u^2) \bar{E} dx + 4u^2 \int_{\mathbb{R}^d} 2G^2 + 2 du \\
 & + \frac{u}{4} \int_{\mathbb{R}^d} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx + \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx :
 \end{aligned}$$

497 By the weighted CKP inequality and given  $\theta = 1$ ,

$$\begin{aligned}
 W_2(\rho_K; \rho_K) & \leq \frac{q}{D_{KL}(\rho_K \| \rho_K)} + \frac{q}{4} \frac{q}{D_{KL}(\rho_K \| \rho_K)} \\
 & \leq \mathfrak{C}_0 \rho_K + \mathfrak{C}_1 \rho_K :
 \end{aligned}$$

498 where the constants  $\mathfrak{C}_0$ ,  $\mathfrak{C}_1$  and  $\mathfrak{A}$  are defined as:

$$\begin{aligned}
 \mathfrak{C}_0 & = \frac{\int_{\mathbb{R}^d} \frac{3u}{2} M^2 (2 + 4M^2 u^2) \bar{E} dx + 4u^2 \int_{\mathbb{R}^d} 2G^2 + 2 du + \int_{\mathbb{R}^d} \frac{3u}{2} M^2 (2 + 4M^2 u^2) \bar{E} dx + 4u^2 \int_{\mathbb{R}^d} 2G^2 + 2 du}{\int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx + \int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx} \\
 \mathfrak{C}_1 & = \frac{\int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx + \int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx}{\int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx + \int_{\mathbb{R}^d} \frac{u}{4} (6M^2 (2C_0 + 2 + 2u^2) + 2) dx} \\
 \mathfrak{A} & = \max \left\{ \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx ; \int_{\mathbb{R}^d} (M^2 + 1) \frac{2^d}{4} dx \right\} :
 \end{aligned}$$

499 Finally by the Lemma A.2 in Zou et al. [2019], we can have

$$W_2(\rho_K; \rho) \leq \rho e^{-K} ;$$

500 where  $\rho = e^{-\Theta(d)}$  denotes the concentration rate of the underdamped Langevin dynamics and  $\Theta$  is a constant of order  $\Theta(1)$ . Combining this inequality with previous analysis we can prove:

$$W_2(\rho_K; \rho) \leq \mathfrak{C}_0 \rho + \mathfrak{C}_1 \rho e^{-K} + \rho e^{-K} : \quad (36)$$

502 In order to bound the Wasserstein distance, we need to set

$$\mathfrak{C}_0 \rho e^{-K} = \frac{\rho}{2} \quad \text{and} \quad \rho e^{-K} = \frac{\rho}{2} : \quad (37)$$

503 Solving the equation (37), we can have

$$K = \frac{\log \frac{\rho}{2}}{\rho} \quad \text{and} \quad \rho = \frac{\rho}{4^{-2} \mathfrak{C}_0^2 K} :$$

504 Combining these two we can have

$$\rho = \frac{\rho}{4^{-2} \mathfrak{C}_0^2 \log \frac{\rho}{2}} \quad \text{and} \quad K = \frac{4^{-2} \mathfrak{C}_0^2 \log^2 \frac{\rho}{2}}{\rho} :$$

505 Plugging in (36) completes the proof.

506 D.3 Proof of Theorem 10

507 In this section we generalize the convergence analysis of LPSGLDLP-F in Zhang et al. [2022] to  
 508 non-log-concave target distribution. We prove a more general version of theorem 10 following the  
 509 same proof outlines in Raginsky et al. [2017]. We further introduce an assumption about the initial  
 510 distribution  $p_0$ .

511 Assumption 5. The probability  $p_0$  of the initial hypothesis  $x_0$  has a bounded and strictly positive  
 512 density and satisfies the following:

$$I_0 := \log \int_{\mathbb{R}^d} e^{kx \cdot k^2} p_0(x) dx < 1 :$$

513 Note that for initial distribution  $x_0 = 0$ , the value  $I_0 = 0$  is bounded and the assumption is  
 514 satisfied. Recall the Overdamped Langevin dynamics is

$$dx_t = -r \nabla U(x_t) dt + \sqrt{\frac{p}{2}} dB_t : \quad (38)$$

515 We further define the value of the energy function and the gradient at  $p_0$  in the following:

$$|U(0)| = G_0; \quad \|\nabla U(0)\| = G_1 :$$

516 In order to analyze the convergence of SGLD for non-log-concave distribution, we need to introduce  
 517 extra assumptions.

518 Then the solution of the Langevin dynamics should satisfy

$$x_t = x_0 + \int_0^t -r \nabla U(x_s) ds + \sqrt{\frac{p}{2}} \int_0^t dB_s : \quad (39)$$

519 To analyze the LPSGLDLP-F in (18), we define a continuous interpolation of the low-precision  
 520 sample as:

$$\hat{x}_t = \hat{x}_0 + \int_0^t G_s ds + \sqrt{\frac{p}{2}} \int_0^t dB_s ; \quad (40)$$

521 where  $G_s = \int_{\mathbb{R}^d} g(\hat{x}_k) 1_{s \in [k, (k+1))} ds$ . The Wasserstein distance can be bounded as

$$W_2(p_K; p) \leq W_2(p_K; \hat{p}_K) + W_2(\hat{p}_K; p) ;$$

522 where the first term of the RHS can be bounded via the weighted CKP inequality

$$W_2(p_K; \hat{p}_K) \leq C_{p_K} \left( \frac{1}{D_{KL}(p_K || \hat{p}_K)} + \frac{D_{KL}(p_K || \hat{p}_K)}{2} \right)^{\frac{1}{q}} ;$$

523 where the constant  $C_{p_K} = 2 \inf_{\gamma > 0} \left( \frac{1}{\gamma} + \log \int_{\mathbb{R}^d} e^{\gamma \|k\|^2} p_K(dk) \right)$ . By Lemma 4 in Raginsky  
 524 et al. [2017] and assuming  $\gamma > 1$ , we can write:

$$W_2^2(p_K; \hat{p}_K) \leq (12 + 8(\gamma_0 + 2b + 2d)K) D_{KL}(p_K || \hat{p}_K) + \frac{q}{D_{KL}(p_K || \hat{p}_K)} ;$$

525 Now we bound the term  $D_{KL}(p_K || \hat{p}_K)$ . The Radon-Nikodym derivative of  $\hat{p}_K$  w.r.t  $p_K$  is  
 526 the following

$$\frac{d\hat{p}_K}{dp_K} = \exp \left( \int_0^T \left( \frac{1}{2} \int_0^t (r \nabla U(x_s) \cdot G_s) dB_s - \frac{1}{4} \int_0^t \|\nabla U(x_s) \cdot G_s\|^2 ds \right) \right) ;$$

527 Thus, we have:

$$\begin{aligned}
 D_{KL}(p_k || \hat{p}_k) &= E_{p_k} \log \frac{dp_k}{d\hat{p}_k} \\
 &= \frac{1}{4} \sum_{k=0}^Z E_{kr} \int_{U(x_s)}^h G_s k^2 ds \\
 &= \frac{1}{4} \sum_{k=0}^Z E_{kr} \int_{U(x_s)}^h g(x_k) k^2 ds \\
 &= \frac{1}{2} \sum_{k=0}^Z E_{kr} \int_{U(x_s)}^h r U(x_k) k^2 ds \\
 &+ \frac{1}{2} \sum_{k=0}^Z E_{kr} \int_{U(x_k)}^h g(x_k) k^2 ds \\
 &+ \frac{M^2}{2} \sum_{k=0}^Z E_{kr} \int_{U(x_s)}^h k x_s x_k k^2 ds \\
 &+ \frac{1}{2} \sum_{k=0}^Z E_{kr} \int_{U(x_k)}^h g(x_k) k^2 ds; \tag{41}
 \end{aligned}$$

528 We now bound the first term in the RHS of the equation 41, from the update rule in 40 we know:

$$\begin{aligned}
 x_s - x_k &= (s - k)g(x_k) + \frac{p}{2}(B_s - B_k) \\
 &= (s - k)r U(x_k) + (s - k)(r U(x_k) - g(x_k)) + \frac{p}{2}(B_s - B_k);
 \end{aligned}$$

529 thus,

$$\begin{aligned}
 E_{kr} \int_{U(x_s)}^h k x_s x_k k^2 ds &\leq 3^2 E_{kr} \int_{U(x_k)}^h k^2 ds + 3^2 E_{kr} \int_{U(x_k)}^h g(x_k) k^2 ds + 6d \\
 &\leq 3^2 (M E[kx_k k] + G)^2 + 3^2 (M^2 + 1) \frac{2d}{4} + 2^2 + 6d; \tag{42}
 \end{aligned}$$

530 Similarly, we need a uniform bound  $E_{kr} \int_{U(x_k)}^h k x_k k^2 ds$ .

531 Lemma 16. Under assumptions 1, 3 and 4, if we set the stepsize  $2e^{-0; 1 \wedge \frac{m_2}{2M^2}}$ , then for all

532  $k \geq 0$ , the  $E_{kr} \int_{U(x_k)}^h k x_k k^2 ds$  can be bounded as

$$E_{kr} \int_{U(x_k)}^h k x_k k^2 ds \leq E + \frac{2(M^2 + 1)2d}{4m_2};$$

533 provided  $E = E_{kr} \int_{U(x_0)}^h k x_0 k^2 ds + \frac{M}{m_2} (2b + 2G^2 + 2d)$ :

534 Using this bound, we can further bound  $E_{kr} \int_{U(x_s)}^h k x_s x_k k^2 ds$  as:

$$\begin{aligned}
 E_{kr} \int_{U(x_s)}^h k x_s x_k k^2 ds &\leq 6^2 M^2 E + \frac{2(M^2 + 1)2d}{m_2} + 6^2 G^2 + 3^2 (M^2 + 1) \frac{2d}{4} + 2^2 + 6d \\
 &\leq 6^2 M^2 E + 6^2 G^2 + 6d + \frac{12^2 M^2 (M^2 + 1)}{m_2} + 3(M^2 + 1) \frac{2^2 d}{4} + 3^2 2^2; \\
 &=: \bar{E} + C \frac{2d}{4} + 3^2 2^2
 \end{aligned}$$

535 where the constant  $\bar{E}$  and  $C$  are defined as:

$$\begin{aligned}
 \bar{E} &= 6M^2 E + 6G^2 + 6d \\
 C &= \frac{12^2 M^2 (M^2 + 1)}{m_2} + 3(M^2 + 1);
 \end{aligned}$$

536 Thus the divergence can be bounded as:

$$\begin{aligned}
 D_{KL}(\rho_K || \rho_K) &= \frac{M^2}{2} \bar{E} + C \frac{2d}{4} + 3^2 K^2 + \frac{1}{2} (M^2 + 1) \frac{2d}{4} + 2^2 K \\
 &= \frac{M^2}{2} \bar{E} K^2 + \frac{M^2}{2} C^2 + \frac{1}{2} (M^2 + 1) \frac{2d}{4} K + \frac{3M^2 + 1}{2} 2^2 K \\
 &= \frac{M^2}{2} \bar{E} K^2 + \frac{M^2}{2} C + \frac{1}{2} (M^2 + 1) \frac{2d}{4} K + \frac{3M^2 + 1}{2} 2^2 K \\
 &=: C_0 K^2 + C_1 \frac{2d}{4} K + C_2 2^2 K :
 \end{aligned}$$

537 We are ready to bound the Wasserstein distance,

$$\begin{aligned}
 W_2^2(\rho_K ; \rho_K) &= (12 + 8(\rho_0 + 2b + 2d)) (C_0 + \frac{p}{C_0})^{p-} + C_1 + \frac{p}{C_1} A + C_2 + \frac{p}{C_2} B (K)^2 \\
 &=: \mathfrak{C}_0^{2p-} + \mathfrak{C}_1^2 A + \mathfrak{C}_2^2 B (K)^2 ;
 \end{aligned}$$

538 where the constants are defined as:

$$\begin{aligned}
 A &= \max \left( \frac{2d}{4}, \frac{r}{\frac{2d}{4}} \right) \\
 B &= \max \left( n^2, \frac{p-0}{2} \right) \\
 \mathfrak{C}_0^2 &= (12 + 8(\rho_0 + 2b + 2d)) C_0 + \frac{p}{C_0} \\
 \mathfrak{C}_1^2 &= (12 + 8(\rho_0 + 2b + 2d)) C_1 + \frac{p}{C_1} \\
 \mathfrak{C}_2^2 &= (12 + 8(\rho_0 + 2b + 2d)) C_2 + \frac{p}{C_2} :
 \end{aligned}$$

539 From Proposition 9 in the paper Raginsky et al. [2017], we know that

$$\begin{aligned}
 W_2(\rho_K ; p) &= \frac{s}{2C_{LS}} \log k p_0 k_1 + \frac{d}{2} \log \frac{3}{m} + \frac{M}{3} \rho_0 + B \frac{p-0}{0} + G_0 + \frac{b}{2} \log 3 e^{K=C_{LS}} \\
 &=: \mathfrak{C}_3 e^{K=C_{LS}}
 \end{aligned}$$

540 Finally, we can have

$$W_2(\rho_K ; p) = \mathfrak{C}_0^{1=4} + \mathfrak{C}_1^p A + \mathfrak{C}_2^p B K + \mathfrak{C}_3 e^{K=C_{LS}} : \quad (43)$$

541 In order to bound the Wasserstein distance, we need to set

$$\mathfrak{C}_0 K^{5=4} = \frac{1}{2} \quad \text{and} \quad \mathfrak{C}_3 e^{K=C_{LS}} = \frac{1}{2} : \quad (44)$$

542 Solving the (37), we can have

$$K = C_{LS} \log \frac{2\mathfrak{C}_3}{\mathfrak{C}_0} \quad \text{and} \quad = \frac{4}{16\mathfrak{C}_0^4 (K)^4} :$$

543 Combining these two we can have

$$= \frac{4}{16\mathfrak{C}_0^4 C_{LS}^4 \log^4 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}} \quad \text{and} \quad K = \frac{16\mathfrak{C}_0^4 C_{LS}^5 \log^5 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}}{4} :$$

544 Plugging  $K$  and into (43) completes the proof.

545 D.4 Proof of Theorem 6

546 Recall the SGHMCLP-L the update rule:

$$\begin{aligned} v_{k+1} &= Q_W v_k e^{-u} (1 - e^{-u}) Q_G(r \mathcal{U}(x_k)) + \frac{v}{k} \\ x_{k+1} &= Q_W x_k + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-u}) Q_G(r \mathcal{U}(x_k)) + \frac{x}{k} : \end{aligned}$$

547 If we let  $\frac{x}{k}$  and  $\frac{v}{k}$  denote the quantization error,

$$\begin{aligned} \frac{x}{k} &= Q_W v_k e^{-u} (1 - e^{-u}) Q_G(r \mathcal{U}(x_s)) + \frac{v}{k} - v_k e^{-u} (1 - e^{-u}) Q_G(r \mathcal{U}(x_s)) + \frac{v}{k} \\ \frac{v}{k} &= Q_W x_s + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-u}) Q_G(r \mathcal{U}(x_s)) + \frac{x}{k} \\ &\quad x_s + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-u}) Q_G(r \mathcal{U}(x_s)) + \frac{x}{k} ; \end{aligned}$$

548 we can rewrite the update rule as:

$$\begin{aligned} v_{k+1} &= v_k e^{-u} (1 - e^{-u}) Q_G(r \mathcal{U}(x_s)) + \frac{v}{k} + \frac{v}{k} \\ x_{k+1} &= x_k + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-u}) Q_G(r \mathcal{U}(x_k)) + \frac{x}{k} + \frac{x}{k} : \quad (45) \end{aligned}$$

549 Similarly, we can define a continuous interpolation of (45) for  $(0; \infty]$ .

$$\begin{aligned} v_t &= v_0 e^{-t} u \int_0^t e^{-(t-s)} (r \mathcal{U}(x_0) + \frac{v}{s}) ds + \frac{v}{2u} \int_0^t e^{-(t-s)} dB_s + \int_0^t v(s) ds \\ x_t &= x_0 + \int_0^t v(s) ds + \int_0^t x(s) ds; \quad (46) \end{aligned}$$

550 where the  $v = Q_G(r \mathcal{U}(x_0)) - r \mathcal{U}(x_0)$  the function  $v(s), x(s)$  are defined as:

$$\begin{aligned} v(s) &= \frac{v}{k} = 1_{s \in (k; (k+1))} \\ x(s) &= \frac{x}{k} = 1_{s \in (k; (k+1))} : \end{aligned}$$

551 If we let  $\rho_0 = (x_0; v_0)$  be the initial sample and  $\rho_t = (x_t; v_t)$  be the sample that satisfies the  
 552 previous equations, we can define an operator  $\hat{\rho}_t$  that maps  $\rho_0$  to  $\rho_t$  i.e.,  $\rho_t = \hat{\rho}_t \rho_0$ . Notice that  
 553 since  $\hat{\rho}_t$  is the continuous interpolation of (4), then  $\rho_k = \hat{\rho}_k \rho_0 = (x_k; v_k)$ . Similarly, we define  
 554  $\mathbf{q}_k = (x_k; v_k + x_k) = (x_k; \frac{x}{k})$  as a tool to analyze the convergence of

555 We are now ready to compute the Wasserstein distance between  $\rho_0$  and  $\rho_t$ . Let  $\gamma_1$  be all of the  
 556 couplings between  $\rho_0$  and  $\rho_t$ , and  $\gamma_2$  be all of the couplings between  $\rho_0$  and  $\rho_t$ . Let  $r_1$  be the  
 557 optimal coupling between  $\rho_0$  and  $\rho_t$ . By taking the difference between (46) and (27),

$$\frac{x}{k} = \frac{x}{k} + u \int_0^R \int_0^R e^{-(s-r)} ds dr + \int_0^R \int_0^R e^{-(s-r)} ds dr + \int_0^R \int_0^R e^{-(s-r)} ds dr + \int_0^R \int_0^R e^{-(s-r)} ds dr + \int_0^R \int_0^R e^{-(s-r)} ds dr :$$



558 Let us now analyze the Wasserstein distance between  $\hat{q}_n$  and  $q$ ,

$$\begin{aligned}
 W_2^2(\hat{q}_n; q) &= E_{r_1} \left[ \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R R_r e^{-(s-r)} ds dr + \int_0^R R_r e^{-(s-r)} ds dr + \int_0^R R_r e^{-(s-r)} ds dr + \int_0^R R_r e^{-(s-r)} ds dr \right] \\
 &= E_{r_1} \left[ \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr \right] \\
 &= W_2^2(q; q) + 4u^2 \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr + \int_0^R \int_0^{R_r} e^{-(s-r)} ds dr \\
 &= W_2^2(q; q) + 4u^2 \left( \frac{4}{4} + \frac{2d}{4} + \frac{2}{4} + 2u^2 E_k \frac{h}{k} k^2 + E_k \frac{h}{k} k^2 \right) \\
 &= W_2^2(q; q) + 5u^2 \left( \frac{2d}{4} + \frac{2}{4} + 2u^2 (A + B) \right);
 \end{aligned}$$

559 where the constants  $A, B$  are the uniform bounds of  $E[k \frac{h}{k} k]$  and  $E[k \frac{v}{k} k]$  respectively. Furthermore  
 560 from the corollary 7 in Cheng et al. [2018] we know that for any  $\epsilon > 0$ ,  $\delta > 0$ ;  $K$ :

$$W_2^2(\hat{q}_n; q) \leq e^{-\delta} W_2^2(q; q); \quad (47)$$

561 where  $\delta = M = m_1$  is the condition number. From the discretization error bound from theorem 9  
 562 and lemma 8 (sandwich inequality) in Cheng et al. [2018], we get

$$W_2(\hat{q}_n; e \hat{q}_n) \leq 2W_2(\hat{p}_n; e \hat{p}_n) \leq \frac{8E_k}{5}.$$

563 By triangle inequality:

$$W_2(e \hat{q}_n; q) \leq W_2(\hat{q}_n; e \hat{q}_n) + W_2(\hat{q}_n; q) \leq \frac{8E_k}{5} + e^{-\delta} W_2(q; q);$$

564 further implies the following inequality:

$$W_2^2(\hat{q}_n; q) \leq e^{-\delta} W_2^2(q; q) + \frac{8E_k}{5} + 5u^2 \left( \frac{2d}{4} + \frac{2}{4} + 2u^2 (A + B) \right);$$

565 By invoking the Lemma 7 in Dalalyan and Karagulyan [2019] we can bound the Wasserstein dis-  
 566 tance by:

$$\begin{aligned}
 W_2(q_K; q) &\leq e^{-K} W_2(q; q) + \frac{q \frac{8E_k}{5}}{1 - e^{-2}} \\
 &+ \frac{q \frac{8E_k}{5} + p \frac{8E_k}{5}}{1 - e^{-2}} \frac{q \frac{8E_k}{5}}{5u^2 \left( \frac{2d}{4} + \frac{2}{4} + 2u^2 (A + B) \right)};
 \end{aligned}$$

567 Finally by sandwich inequality we have:

$$\begin{aligned}
 W_2(p_K; p) &\leq 4e^{-K} W_2(q; q) + \frac{4 \frac{8E_k}{5}}{1 - e^{-2}} \\
 &+ \frac{20u^2 \left( \frac{2d}{4} + \frac{2}{4} + 8u^2 (A + B) \right)}{1 - e^{-2}} \frac{q \frac{8E_k}{5}}{5u^2 \left( \frac{2d}{4} + \frac{2}{4} + 2u^2 (A + B) \right)};
 \end{aligned} \quad (48)$$

568 And in this case, we know that  $E[k_x^k]$  and  $E[k_y^k]$  can be bounded by  $\frac{2^d}{4}$ . Finally, we can have:

$$W_2(p_K; p) \leq 4e^{K-2} W_2(q_0; q) + \frac{q \frac{8E_K}{5}}{1 - e^{-2}} + \frac{20u^{2-2} \frac{2^d}{4} + 2 + 4u^{2-2}d}{5u^{2-2} \frac{2^d}{4} + 2 + u^{2-2}d}:$$

569 Now we let the first term less than 3, from the lemma 13 in [Cheng et al., 2018] we know that  
 570  $W_2(q_0; q) \leq 3 \frac{d}{m_1} + D^2$ . So we can choose  $K$  as the following,

$$K \geq \frac{2}{1} \log 36 \frac{d}{m_1} + D^2:$$

571 Next, we choose a stepsize  $p = \frac{1}{479232=5(d=m_1+D^2)}$  to ensure the second term is controlled below

572  $\leq 3$ . Since  $1 - e^{-2} = 4$  and definition of  $E_K$ ,

$$4 \frac{q \frac{8E_K}{5}}{1 - e^{-2}} = 4 \frac{q \frac{8E_K}{5}}{4} = 16 \frac{q \frac{8E_K}{5}}{5} \leq 3:$$

573 Finally by choosing the stepsize satisfied that,

$$\frac{2}{2880} u \frac{2^d}{4} + 2;$$

574 the third term can be bounded as:

$$\begin{aligned} & \frac{q \frac{20u^{2-2} (M^2+1) \frac{2^d}{4} + 2 + 4u^{2-2}d}{5} + p \frac{q \frac{20u^{2-2} (M^2+1) \frac{2^d}{4} + 2 + 4u^{2-2}d}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2}}{1 - e^{-2}} \\ & \leq \frac{p \frac{20u^{2-2} (M^2+1) \frac{2^d}{4} + 2 + 4u^{2-2}d}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2}}{4} + \frac{20u^{2-2} (M^2+1) \frac{2^d}{4} + 2 + 4u^{2-2}d}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2} \\ & \leq \frac{20u^{2-2} (M^2+1) \frac{2^d}{4} + 2 + 4u^{2-2}d}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2} + \frac{8u^{2-2}d^p}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2} \\ & = 3 + \frac{8u^{2-2}d^p}{5u^{2-2} (M^2+1) \frac{2^d}{4} + 2}: \end{aligned}$$

575 This complete the proof.

### 576 D.5 Proof of Theorem 7

577 In this section, we analyze the convergence of SGHMCLP-L when the target distribution is non-log-  
 578 concave. Recall the continuous interpolation of the SGHMCLP-L,

$$\begin{aligned} v_t &= v_0 + \int_0^t v_s ds + \int_0^t G_s ds + p \int_0^t \frac{Z_t}{2u} e^{-(t-s)} dB_s + \int_0^t v(s) ds \\ x_t &= x_0 + \int_0^t v_s ds + \int_0^t x(s) ds; \end{aligned}$$

579 where  $G_s = \prod_{k=0}^P Q_G(r, U(x_k^0)) 1_{s \in (k; (k+1))}$ . And we define an intermediate process by  $v_t =$

580  $v_t + x(t):$

$$v_t^0 = v_0^0 + \int_0^t (v_s^0 - x(s)) ds + \int_0^t G_s ds + \rho \int_0^t \frac{Z_t}{2u} e^{-(t-s)} dB_s + \int_0^t (v(s) + \frac{1}{T} x(t)) ds$$

$$x_t^0 = x_0^0 + \int_0^t v_s^0 ds: \tag{49}$$

581 By integrating the underdamped Langevin dynamic (9), we can have:

$$v_t = v_0 + \int_0^t (v_s - x(s)) ds + \int_0^t r U(x_s) ds + \rho \int_0^t \frac{Z_t}{2u} e^{-(t-s)} dB_s$$

$$x_t = x_0 + \int_0^t v_s ds: \tag{50}$$

582 Notice that the process  $x^0$  has the same distribution with  $x$ , thus in the following analysis we study  
583 the convergence of the intermediate process  $(x_k^0; v_k^0)$ . By taking the difference of equation

584 (49) with (50) and the Girsanov formula, we can derive the Radon-Nikodym derivative of  $\mathbb{P}_K^0$  w.r.t  
585  $\mathbb{P}_K^0$ :

$$\frac{d\mathbb{P}_K}{d\mathbb{P}_K^0} = \exp \left( \int_0^T \frac{Z_T}{2u} (x(s) + v(s) + \frac{1}{T} x(T) + r U(x_s) - G_s) dB_s - \frac{1}{4} \int_0^T \frac{Z_T}{u} k^2 (x(s) + v(s) + \frac{1}{T} x(T) + r U(x_s) - G_s)^2 ds \right)$$

586 Thus the divergence can be bounded as:

$$D_{KL}(\mathbb{P}_K || \mathbb{P}_K^0) = E_{\mathbb{P}_K} \log \frac{d\mathbb{P}_K}{d\mathbb{P}_K^0}$$

$$= \frac{u}{4} \int_0^T E \left( x(s) + v(s) + \frac{1}{T} x(T) + r U(x_s) - G_s \right)^2 ds$$

$$= \frac{u}{4T} E \int_0^T k^2 x(s)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (v(s) + x(s) + r U(x_s) - G_s)^2 ds$$

$$= \frac{u}{4T} E \int_0^T k^2 x(s)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (v(s))^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 x(s)^2 ds$$

$$+ \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (r U(x_s) - G_s)^2 ds$$

$$+ \frac{u}{4T} E \int_0^T k^2 x(s)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (v(s) = k^2)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 x(s) = k^2 ds$$

$$+ \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (r U(x_s) - Q_G(r U(x_k)))^2 ds$$

$$+ \frac{u}{4T} E \int_0^T k^2 x(s)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (v(s) = k^2)^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 x(s) = k^2 ds$$

$$+ \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (r U(x_s) - r U(x_k))^2 ds + \frac{u}{4} \sum_{k=0}^{k+1} E \int_0^T k^2 (r U(x_k) - Q_G(r U(x_k)))^2 ds: \tag{51}$$

587 By assumption 1, we know that:

$$E_{h,k} \leq M^2 E_{h,k}^0$$

588 From the same analysis in (35), we can derive:

$$E_{h,k} \leq 3M^2 E_{h,k}^0 + u^2 E_{h,k}^0 + 2du$$

589 Now we need to derive a uniform bound of  $E_{h,k}^0$  and  $E_{h,k}^0$ .

590 Lemma 17. Let Assumptions 3 and 1 hold. If we set the step size to the following condition

$$\min \left\{ \frac{4u^2}{4(8Mu + u^2 + 22d^2)}; \frac{6bu}{4Mu + 3d^2}; \frac{m_2}{6(22u + d^2)M^2} \right\};$$

591 then for all  $k > 0$   $E_{h,k}^0$  and  $E_{h,k}^0$  can be bounded as follow:

$$E_{h,k}^0 \leq E + C^2 d; \quad E_{h,k}^0 \leq 2E + 2C^2 d;$$

592 where constant  $E$  and  $C$  are defined as:

$$E = E[E(x_0; v_0)] + \frac{54}{m_2^4} (4u + d^2) u^2 + \frac{12(22u + d^2)uM^3}{m_2^3} G^2 + \frac{96(d + b)uM}{m_2^2}$$

$$C = \frac{27}{2m_2^4} (4u + d^2) u;$$

593

594 Thus,

$$\begin{aligned} E_{h,k} &\leq 3M^2 E_{h,k}^0 + u^2 E_{h,k}^0 + \frac{2d}{4} + 2M^2 E_{h,k}^0 + 2G^2 + 2du \\ &\leq 3M^2 (2E + 2C^2 d) + u^2 (2E + 2C^2 d) + \frac{2d}{4} + 2M^2 (2E + 2C^2 d) + 2G^2 + 2du \\ &\leq 3M^2 (2E + 2C^2 d) + 2u^2 M^2 (E + C^2 d) + u^2 (2E + 2C^2 d) + 2u^2 G^2 + 2du \end{aligned}$$

595 Now we can go back to the divergence of  $\rho_k$ ,

$$D_{KL}(\rho_k \| \rho_k)$$

$$\begin{aligned} &\frac{u}{4T^2} E_{h,k}^0 + \frac{u}{4} \sum_{k=0}^{Z^{(k+1)}} E_{h,k}^0 ds + \frac{u}{4} \sum_{k=0}^{Z^{(k+1)}} E_{h,k}^0 ds \\ &+ \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + \frac{2d}{4} + 2u^2 M^2 C^2 d + u^2 (2E + 2C^2 d) + 2u^2 G^2 + 2du + \frac{u}{4} K \left( \frac{2d}{4} + 2 \right) \\ &+ \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + \frac{2d}{4} + 2u^2 M^2 C^2 d + u^2 (2E + 2C^2 d) + 2u^2 G^2 + 2du + \frac{u}{4} K \left( \frac{2d}{4} + 2 \right) \\ &+ \frac{u}{16T^2} + \frac{uK}{8} \frac{2d}{4} \\ &+ \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 (2E + 2C^2 d) + 2u^2 G^2 + 2du + \frac{u}{4} K^2 \\ &+ \frac{u}{4} 3M^2 K^3 C^2 + 2u^2 M^2 C^2 + \frac{uK}{16} + \frac{u}{16T^2} + \frac{uK}{8} 2d \\ &=: C_0 K^3 + C_1 K^2 + C_2 K^2; \end{aligned}$$

596 where the constant  $C_0$ ,  $C_1$  and  $C_2$  are defined as:

$$C_0 = \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 (2E + 2C^2 d) + 2u^2 G^2 + 2du$$

$$C_1 = \frac{u}{4}$$

$$C_2 = \frac{u}{4} 3M^2 K^3 C^2 + 2u^2 M^2 C^2 + \frac{u}{16} + \frac{u}{16T^2} + \frac{u}{8} d;$$

597 By the weighted CKP inequality and given 1,

$$W_2(p_K; p_K) = \frac{q}{D_{KL}(p_K || p_K)} + \frac{q}{4 D_{KL}(p_K || p_K)} \left( \mathfrak{C}_0^{p-} + \mathfrak{C}_1 \mathfrak{A}^{p-} \overline{K} + \mathfrak{C}_2^{p-} \overline{K} \right); \quad (52)$$

598 where the constants are defined as:

$$\begin{aligned} \mathfrak{C}_0 &= \frac{p}{C_0} + \frac{p_4}{C_0} \\ \mathfrak{C}_1 &= \frac{p}{C_1} + \frac{p_4}{C_1} \\ \mathfrak{C}_2 &= \frac{p}{C_2} + \frac{p_4}{C_2} \\ \mathfrak{A} &= \max_{\cdot} ; p- : \end{aligned}$$

599 From the same analysis in (36), we can have:

$$W_2(p_K; p) = \mathfrak{C}_0^{p-} + \mathfrak{C}_1 \mathfrak{A}^{p-} \overline{K} + \mathfrak{C}_2^{p-} \overline{K} + \rho e^{-K}; \quad (53)$$

600 In order to bound the Wasserstein distance, we need to set

$$\mathfrak{C}_0^{p-} \overline{K}^2 = \frac{\rho}{2} \quad \text{and} \quad \rho e^{-K} = \frac{\rho}{2}; \quad (54)$$

601 Solving the equation (54), we can have

$$K = \frac{\log \frac{\rho}{2}}{\rho} \quad \text{and} \quad \overline{K} = \frac{2}{4^{-2} \mathfrak{C}_0^2 K};$$

602 Combining these two we can have

$$\overline{K} = \frac{2}{4^{-2} \mathfrak{C}_0^2 \log \frac{\rho}{2}} \quad \text{and} \quad K = \frac{4^{-2} \mathfrak{C}_0^2 \log^2 \frac{\rho}{2}}{2 \left( \frac{\rho}{2} \right)^2};$$

603 Plugging in (53) completes the proof.

## 604 D.6 Proof of Theorem 11

605 In this section we generalize the convergence analysis of SGLDLP-L in Zhang et al. [2022] to non-  
606 log-concave target distribution. Following the same proof outlines in Raginsky et al. [2017]. Recall  
607 the LPSGLDLP-L update rule 19 is the following,

$$\begin{aligned} x_{k+1} &= Q_W(x_k, r \mathcal{U}(x_k) + \frac{p}{2} \overline{z}_{k+1}) \\ &=: x_k, r \mathcal{U}(x_k) + \frac{p}{2} \overline{z}_{k+1} + \kappa; \end{aligned}$$

608 where  $\kappa$  is defined as:

$$\kappa = Q_W(x_k, r \mathcal{U}(x_k) + \frac{p}{2} \overline{z}_{k+1}) - x_k, r \mathcal{U}(x_k) + \frac{p}{2} \overline{z}_{k+1};$$

609 Thus, we can define a continuous interpolation of the SGLDLP-L as:

$$x_t = x_0 + \int_0^t G_s ds + \frac{p}{2} \int_0^t dB(s) + \int_0^t \kappa(s) ds;$$

610 where  $G_s = \frac{p}{k=0} Q_G(r \mathcal{U}(x_k)) 1_{s \in (k; (k+1))}$  and  $\kappa(s) = \frac{p}{k=0} \kappa = 1_{s \in (k; (k+1))}$ . By taking the

611 difference of the interpolation with the discrete estimation of Langevin process in equation 39, we can  
612 derive the Radon-Nikodym derivative of  $\mathbb{P}_K$  w.r.t  $p_K$  as:

$$\frac{d\mathbb{P}_K}{dp_K} = \exp \left( \frac{1}{2} \int_0^T (r \mathcal{U}(x_s) - G_s) dB(s) - \frac{1}{4} \int_0^T \kappa(r \mathcal{U}(x_s) - G_s)^2 ds \right);$$

613 Thus, the divergence can be computed as:

$$\begin{aligned}
 D_{KL}(p_K || p_K) &= \frac{1}{4} \int_0^Z E_{kr} U(x_s) G_s(s) k^2 ds \\
 &= \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{kr} U(x_s) Q_G(r U(x_k)) k^2 ds \\
 &= \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{kr} U(x_s) Q_G(r U(x_k))^2 ds + \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{k_k} k^2 ds \\
 &= \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{kr} U(x_s) r U(x_k) k^2 ds + \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{kr} U(x_k) Q_G(r U(x_k))^2 ds \\
 &\quad + \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{k_k} k^2 ds \\
 &= \frac{M^2}{4} \int_{k=0}^{Z^{(k+1)}} E_{kx_s} x_k k^2 ds + \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{kr} U(x_k) Q_G(r U(x_k))^2 ds \\
 &\quad + \frac{1}{4} \int_{k=0}^{Z^{(k+1)}} E_{k_k} k^2 ds: \tag{55}
 \end{aligned}$$

614 From the same analysis in (35), we know that

$$\begin{aligned}
 E_{kx_s} x_k k^2 &= 3^2 E_{kr} U(x_k) k^2 + 3^2 E_{kr} U(x_k) Q_G(r U(x_k))^2 + 6d \\
 &= 3^2 M E_{kx_k} k^2 + G^2 + 3^2 \frac{2d}{4} + 2^2 + 6d:
 \end{aligned}$$

615 Again, we need to derive a uniform bound of  $E_{kx_k} k^2$ ,

$$\begin{aligned}
 E_{kx_{k+1}} k^2 &= E_{x_k} Q_G(r U(x_k))^2 + 2E_{k_{k+1}} k^2 + E_{k_k} k^2 \\
 &= E_{x_k} r U(x_k) + r U(x_k) Q_G(r U(x_k))^2 + 2d + E_{k_k} k^2 \\
 &= E_{x_k} r U(x_k) + r U(x_k) Q_G(r U(x_k))^2 + E_{k_k} k^2 + 2d \\
 &= E_{kx_k} r U(x_k) k^2 + 2E_{kr} U(x_k) Q_G(r U(x_k))^2 + E_{k_k} k^2 + 2d:
 \end{aligned}$$

616 By plugging in the inequality we derived before:

$$E_{kx_k} r U(x_k) k^2 \leq (1 - 2m_2 + 2^2 M^2) E_{kx_k} k^2 + 2b + 2^2 G^2:$$

617 we can have:

$$E_{kx_{k+1}} k^2 \leq (1 - 2m_2 + 2^2 M^2) E_{kx_k} k^2 + 2b + 2^2 G^2 + \frac{2^2 d}{4} + 2^2 + E_{k_k} k^2 + 2d:$$

618 Thus for any  $\alpha \in (0, 1)$  and  $1 - 2m_2 + 2M^2 > 0$ , we can bound  $E_k x_k^i$  for any  
 619  $k > 0$  as:

$$E_k x_k^i \leq E_k x_0^i + \frac{1}{2(m_2 - M^2)} (2b + 2G^2 + \frac{2d}{4} + \alpha + 2d) + \frac{E_k x_k^i}{2(m_2 - M^2)}$$

$$E_k x_0^i + \frac{1}{m_2} (2b + 2G^2 + \frac{2d}{4} + \alpha + 2d) + \frac{E_k x_k^i}{m_2}$$

$$E_k + \frac{2d}{4m_2} + \frac{E_k x_k^i}{m_2};$$

620 where the constant  $\bar{E}$  is defined as:

$$\bar{E} = E_k x_0^i + \frac{1}{m_2} (2b + 2G^2 + \alpha + 2d);$$

621 Thus, we can have,

$$E_k x_s^i \leq \bar{E} + \frac{2d}{4m_2} + \frac{E_k x_k^i}{m_2} A + 6G^2 + 3\alpha + \frac{2d}{4} + \alpha + 6d$$

$$\bar{E} + 3\alpha + \frac{6 + 3m_2}{4m_2} \alpha + \frac{6E_k x_k^i}{m_2};$$

622 Plugging this into the equation (55), we can have,

$$D_{KL}(p_K; \hat{p}_K) \leq \frac{M\bar{E}}{4} K^2 + \frac{3M}{4} K^3 + \frac{(6 + 3m_2)M}{16m_2} \alpha K^3 + \frac{6ME_k x_k^i}{4m_2} K^2 + \frac{1}{4} \left( \frac{2d}{4} + \alpha \right) K + \frac{KE}{4}$$

$$\frac{M\bar{E}}{4} K^2 + \frac{3M + 1}{4} \alpha K^2 + \frac{((6 + 3m_2)M + m_2)d}{16m_2} \alpha K^2 + \frac{6M}{4m_2} + \frac{1}{4} KE_k x_k^i;$$

623 By the fact that  $E_k x_k^i \leq \frac{2d}{4}$ , we can further bound the divergence as:

$$D_{KL}(p_K; \hat{p}_K) \leq \frac{M\bar{E}}{4} K^2 + \frac{3M + 1}{4} \alpha K^2 + \frac{((12 + 3m_2)M + m_2)d}{16m_2} \alpha + \frac{d}{16} \alpha K^2$$

$$=: C_0 K^2 + C_1 \alpha K^2 + C_2 \alpha K;$$

624 where the constants are defined as:

$$C_0 = \frac{M\bar{E}}{4}$$

$$C_1 = \frac{3M + 1}{4}$$

$$C_2 = \frac{((12 + 3m_2)M + m_2)d}{16m_2} + \frac{d}{16};$$

625 We are ready to bound the Wasserstein distance,

$$W_2^2(p_K; \hat{p}_K) \leq (12 + 8(\alpha + 2b + 2d)) C_0 + \alpha C_0 + C_1 + \alpha C_1 A (K^2)^2 + C_2 + \alpha C_2 K^2$$

$$=: \hat{C}_0^{2p} + \hat{C}_1^2 A (K^2)^2 + \hat{C}_2^2 K^2;$$

626 where the constants are defined as:

$$A = \max_{\alpha \in (0, 1)} \alpha^{2p-2};$$

$$\hat{C}_0^2 = (12 + 8(\alpha + 2b + 2d)) C_0 + \alpha C_0$$

$$\hat{C}_1^2 = (12 + 8(\alpha + 2b + 2d)) C_1 + \alpha C_1$$

$$\hat{C}_2^2 = (12 + 8(\alpha + 2b + 2d)) C_2 + \alpha C_2;$$

627 From Proposition 9 in the paper Raginsky et al. [2017], we know that

$$W_2(\mathbb{P}_K; p) \leq \frac{S}{2C_{LS}} \log k p_0 k_1 + \frac{d}{2} \log \frac{3}{m} + \frac{M_0}{3} + B_0^p + G_0 + \frac{b}{2} \log 3 e^{K=C_{LS}}$$

$$=: \mathfrak{C}_3 e^{K=C_{LS}}$$

628 Finally, we can have

$$W_2(\mathbb{P}_K; p) \leq \mathfrak{C}_0^{1-4} + \mathfrak{C}_1^p \bar{A} K + \mathfrak{C}_2^{p-p} \bar{K}^2 + \mathfrak{C}_3 e^{K=C_{LS}} : (56)$$

629 In order to bound the  $\mathfrak{C}$ -Wasserstein distance, we need to set

$$\mathfrak{C}_0 K^{5-4} = \frac{1}{2} \text{ and } \mathfrak{C}_3 e^{K=C_{LS}} = \frac{1}{2} : (57)$$

630 Solving the (57), we can have

$$K = C_{LS} \log \frac{2\mathfrak{C}_3}{\mathfrak{C}_0} \text{ and } \frac{16\mathfrak{C}_0^4}{(K)^4} :$$

631 Combining these two we can have

$$\frac{16\mathfrak{C}_0^4 C_{LS}^4 \log^4 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}}{16\mathfrak{C}_0^4 C_{LS}^4 \log^4 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}} \text{ and } K = \frac{16\mathfrak{C}_0^4 C_{LS}^5 \log^5 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}}{4} :$$

632 Plugging  $K$  and into (56) completes the proof.

### 633 D.7 Proof of Theorem 8

634 In this section, we analyze the convergence of VC SGHMCLP-L, recall the VC SGHMCLP-L up-  
635 date rule is the following,

$$v_{k+1} = Q^{VC} v_k e^{-u} (1 - e^{-u}) Q_G(r \psi(x_k)); \text{Var}_v;$$

$$x_{k+1} = Q^{VC} x_k + (1 - e^{-u}) v_k + u^{-2} + e^{-1} Q_G(r \psi(x_k)); \text{Var}_x; : (58)$$

636 If we let  $\tilde{x}_k$  and  $\tilde{v}_k$  denote the quantization error,

$$\tilde{v}_k = Q^{VC} v_k e^{-u} (1 - e^{-u}) Q_G(r \psi(x_k)); \text{Var}_v; \quad v_k e^{-u} (1 - e^{-u}) Q_G(r \psi(x_k)) + \tilde{v}_k$$

$$\tilde{x}_k = Q^{VC} x_k + (1 - e^{-u}) v_k + u^{-2} + e^{-1} Q_G(r \psi(x_k)); \text{Var}_x;$$

$$x_k + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-1}) Q_G(r \psi(x_k)) + \tilde{x}_k ;$$

637 we can rewrite the update rule as:

$$v_{k+1} = v_k e^{-u} (1 - e^{-u}) Q_G(r \psi(x_k)) + \tilde{v}_k + \tilde{v}_k$$

$$x_{k+1} = x_k + (1 - e^{-u}) v_k + u^{-2} (1 + e^{-1}) Q_G(r \psi(x_k)) + \tilde{x}_k + \tilde{x}_k :$$

638 Next, we first derive a uniform bound on  $\mathbb{E} \|\tilde{x}_k\|_2^2$ . In this section and the following section, we  
639 further assume the norm of quantized stochastic gradients are bounded.

640 Assumption 6. For any  $x \in \mathbb{R}^d$ , there exists a constant  $G$  and the quantized stochastic gradients at  
641  $x$  satisfies the following

$$\mathbb{E} \|Q_G(r \psi(x))\|_2^2 \leq G^2 :$$



642 By the definition of the variance corrected quantization function, when  $\text{Var}_v > 0 = \frac{2}{4}$ , if

643 we let  $k$  denote  $v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}$ ,

$$E[k^2] = E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 + \frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2$$

644 Let

$$b = Q_G(r \vartheta(x_k)) + \frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2$$

645 then

$$E[k^2] = E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 + \frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2 + b \text{sign}(r)c^2$$

$$= E\left[\frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2 + b \text{sign}(r)c^2\right]$$

$$= E\left[\frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2 + b \text{sign}(r)c^2\right] + E[kb + \text{sign}(r)c^2]$$

$$= \frac{2\text{Var}_v d}{4ud} + \frac{p}{4ud} \text{sign}(r)c^2 \tag{59}$$

646 When  $\text{Var}_v < \frac{2}{4}$ ,

$$E[k^2] = E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 + \frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2$$

$$= E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 + E\left[\frac{p}{\text{Var}_v} Q_G(r \vartheta(x_k))^2\right]$$

$$\max\{2E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2, 2\text{Var}_v\} + \frac{p}{4ud} Q_G(r \vartheta(x_k))^2 \tag{60}$$

647 Using the bound equation (6) in Li and De Sa [2019] gives us,

$$E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 \leq \frac{p}{d} E[kv_k k] + E[Q_G(r \vartheta(x_k))]^2$$

648 Now we need to derive a uniform bound  $E[kv_k k]$ , by the update rule, we know that,

$$E[kv_{k+1} k^2] = E[v_k e^{-u^{-1}(1-e^{-Q_G(r \vartheta(x_k))})}]^2 + \frac{p}{4ud} Q_G(r \vartheta(x_k))^2 + 2ud + E[k^2]$$

$$(1 + \frac{p}{4ud}) E[kv_k k^2] + \frac{2}{4ud} + 1 + u^2 E[Q_G(r \vartheta(x_k))]^2 + 2ud + E[k^2]$$

$$(1 + \frac{p}{4ud}) E[kv_k k^2] + 3u^2 = G^2 + 2ud + E[k^2]$$

649 When  $E[k_k^i] \leq 2V ar_v d < 4ud$ , the inequality can be further written as:

$$E[k_{k+1}^i] (1 - 2) E[k_k^i] + 3u^2 = G^2 + 6ud$$

$$E[k_0^i] + \frac{6u^2 G^2}{2} + \frac{12ud}{2}$$

$$E[k_0^i] + \frac{6u^2 G^2}{2} + 12ud:$$

650 If  $E[k_k^i] \leq 2E[v_k e^{-u(1-e)}] Q_G(r \mathcal{U}(x_k)) - Q^s v_k e^{-u(1-e)} Q_G(r \mathcal{U}(x_k))^2$ ,

651 the inequality can be written as:

$$E[k_{k+1}^i] (1 - 2) E[k_k^i] + 3u^2 = G^2 + 2ud + 2 \frac{1 - e^{-p_d}}{p_d} E[k_k] + E^h Q_G(r \mathcal{U}(x_k))$$

$$(1 - 2) E[k_k^i] + 3u^2 = G^2 + 2ud + 2 \frac{1 - e^{-p_d}}{p_d} E[k_k^i] + G$$

$$p \frac{1 - e^{-p_d}}{1 - 2} E[k_k^i] + p \frac{1 - e^{-p_d}}{1 - 2} + 3u^2 = G^2 + 2ud + 2 \frac{1 - e^{-p_d}}{p_d} G:$$

652 Thus,

$$E[k_k] \frac{1 - e^{-p_d}}{p_d} + \frac{1 - e^{-p_d}}{1 - 2} p \frac{1 - e^{-p_d}}{p_d} + \frac{3u^2 = G^2 + 2ud + 2 \frac{1 - e^{-p_d}}{p_d} G}{p \frac{1 - e^{-p_d}}{1 - 2} + \frac{1 - e^{-p_d}}{1 - 2} 3u^2 = G^2 + 2ud + 2 \frac{1 - e^{-p_d}}{p_d} G}$$

$$\frac{1 - e^{-p_d}}{p_d} + \frac{1 - e^{-p_d}}{1 - 2} + \frac{q}{6u^2 = 2G^2 + 4ud + 4} p \frac{1 - e^{-p_d}}{p_d}$$

$$\frac{1 - e^{-p_d}}{p_d} + \frac{q}{6u^2 = 2G^2 + 4ud + 4} p \frac{1 - e^{-p_d}}{p_d}:$$

653 Finally, we can have:

$$E[k_k] \max \left( \frac{1 - e^{-p_d}}{p_d} + \frac{q}{6u^2 = 2G^2 + 4ud + 4} p \frac{1 - e^{-p_d}}{p_d}; \right.$$

$$\left. \frac{1 - e^{-p_d}}{p_d} + \frac{6u^2 G^2}{2} + \frac{12ud}{2} \right) =: A^0.$$

654 Thus, we can have,

$$E[v_k e^{-u(1-e)}] Q_G(r \mathcal{U}(x_k)) - Q^s v_k e^{-u(1-e)} Q_G(r \mathcal{U}(x_k))^2$$

$$p_d(A^0 + G);$$

655 and we can bound the  $E[k_k^i]$  as,

$$E[k_k^i] \max^n \left( p_d(A^0 + G); 4ud \right)$$

$$= \max^n \left( p_d(A^0 + G); 4ud \right)$$

$$=: A: \tag{61}$$

656 Now we bound the  $E_k^h x_k^i k^2$ . When  $V_{ar_x} \geq 0$ , as the same analysis in (59) we can show,

$$E_k^h x_k^i k^2 \leq 2V_{ar_x} d + 4ud^2;$$

657 If  $V_{ar_x} < 0$ , and let  $x_k = x_k + \frac{1}{2}(1 - e^{-\frac{1}{2}u})v_k + u^{-2}(1 + e^{-\frac{1}{2}u})Q_G(r \mathcal{U}(x_k))$ , by  
658 the same analysis in (60) we can have:

$$E_k^h x_k^i k^2 \leq \max_n \{ 2E_k^h x_k^i Q^s(x_k) k^2; 2V_{ar_x} d \};$$

659 Again using the bound equation (6) in Li and De Sa [2019] gives us,

$$\begin{aligned} E_k^h x_k^i Q^s(x_k) k^2 &\leq E_k^h \left[ \frac{1}{2} (1 - e^{-\frac{1}{2}u}) v_k + u^{-2} (1 + e^{-\frac{1}{2}u}) Q_G(r \mathcal{U}(x_k)) \right]^2 k^2 \\ &\leq E[kv_k k] + \frac{u^{-2}}{2} E[Q_G(r \mathcal{U}(x_k))]^2 k^2 \\ &\leq p_{dE}^- [kv_k k] + \frac{u^{-2}}{2} p_{dE}^- E[Q_G(r \mathcal{U}(x_k))]^2 k^2 \\ &\leq p_{dA}^- + \frac{u^{-2}}{2} p_{dG}^-; \end{aligned}$$

660 Thus, we can have,

$$\begin{aligned} E_k^h x_k^i k^2 &\leq \max_n \{ p_{dA}^- + u^{-2} p_{dG}^-; 4ud^2 \} \\ &\leq \max_n \{ p_{dA}^- + u^{-2} p_{dG}^-; 4ud^2 \} \\ &=: B; \end{aligned} \tag{62}$$

661 Then follow the same analysis of (48), we can show

$$\begin{aligned} W_2(\rho_K; p) &\leq 4e^{-K} = 2^{-1} W_2(q_0; q) + \frac{4^{-2} q \frac{8E_K}{5}}{1 - e^{-2^{-1}}} \\ &\quad + \frac{20u^{2^{-2}} \frac{2^{-2}d}{4} + 2^{-2} + 8u^2 (A + B)}{2^{-2} \frac{8E_K}{5} + p \frac{1}{1 - e^{-1}} q \frac{2^{-2}d}{4} + 2^{-2} + 2u^2 (A + B)}; \end{aligned}$$

662 Now we let the first term less than  $\epsilon=3$ , from the Lemma 13 in [Cheng et al., 2018] we know that

663  $W_2(q_0; q) \leq 3 \frac{d}{m_1} + D^2$ . So we can choose  $K$  as the following,

$$K \geq \frac{2^{-1}}{\log 36} \left( \frac{d}{m_1} + D^2 \right);$$

664 Next, we choose a stepsize  $p = \frac{1}{479232=5(d=m_1 + D^2)}$  to ensure the second term is controlled below

665  $\epsilon=3$ . Since  $1 - e^{-2^{-1}} = \frac{1}{4}$  and definition of  $E_K$ ,

$$4 \frac{q \frac{8E_K}{5}}{1 - e^{-2^{-1}}} = 4 \frac{q \frac{8E_K}{5}}{1 - e^{-2^{-1}}} = 16 \frac{q \frac{8E_K}{5}}{1 - e^{-2^{-1}}} = 3;$$

666 Finally by choosing the stepsize satisfied that,

$$\frac{2}{2880} u \frac{2^{-2}d}{4} + 2^{-2};$$

667 the third term can be bounded as:

$$\begin{aligned}
 & \frac{20u^2 \left( \frac{2d}{4} + 2 + 8u^2 (A+B) \right)}{2 \left( \frac{8E_k}{5} + 1 \right) e^{-1} \left( 5u^2 \left( \frac{2d}{4} + 2 + 2u^2 (A+B) \right) \right)} \\
 & \frac{20u^2 \left( \frac{2d}{4} + 2 + 8u^2 (A+B) \right)}{1 \left( e^{-1} \right) \left( 5u^2 \left( \frac{2d}{4} + 2 + 2u^2 (A+B) \right) \right)} \quad \frac{20u^2 \left( \frac{2d}{4} + 2 + 8u^2 (A+B) \right)}{=4 \left( 5u^2 \left( \frac{2d}{4} + 2 + 2u^2 (A+B) \right) \right)} \\
 & \frac{4 \cdot 20u^2 \left( \frac{2d}{4} + 2 + 8u^2 (A+B) \right)}{=3+8 \left( 2 \cdot 1u^2 (A+B) \right)}:
 \end{aligned}$$

668 This complete the proof.

### 669 D.8 Proof of Theorem 9

670 Similarly, from the analysis in (61), we know that

$$E \left[ k \sum_k^i k^2 \right] \leq A; \tag{63}$$

671 where  $A = \max_n \left( \frac{1}{d} (A^0 + G) \right); 4ud^2$ . By the analysis in (59), we know that  $\text{Var}_x^{\text{hmc}} < \frac{2}{4}$ ,  
 672 we can have

$$E \left[ k \sum_k^i k^2 \right] \leq 4ud^2 \tag{64}$$

673 by (62), if  $\text{Var}_x^{\text{hmc}} < \frac{2}{4}$ ,

$$E \left[ k \sum_k^i k^2 \right] \leq B; \tag{65}$$

674 where  $B = \max_n \left( \frac{1}{2} \left( \frac{1}{d} A^0 + u \left( \frac{1}{d} G \right) \right); 4ud^2 \right)$ . Thus, we can define the following:

$$E \left[ k \sum_k^i k^2 \right] = B; \tag{66}$$

675 where  $B$  is defined as:

$$B = \begin{cases} 4ud; & \text{if } \text{Var}_x^{\text{hmc}} < \frac{2}{4} \\ B; & \text{else} \end{cases}$$

676 Combining the bound  $\int_{\mathbb{R}^n} \frac{h}{k} \chi_{k^2}^i$ ,  $\int_{\mathbb{R}^n} \frac{h}{k} \chi_{k^2}^i$  with (51), we can show,

$$\begin{aligned}
 & D_{KL}(p_K; j; p_K) \\
 & \frac{u}{4T^2} \int_{\mathbb{R}^n} \frac{h}{k} \chi_{k^2}^i + \frac{u}{4} \sum_{k=0}^Z \int_{\mathbb{R}^n} \frac{h}{k} \chi_{k^2}^i ds + \frac{u}{4} \sum_{k=0}^Z \int_{\mathbb{R}^n} \frac{h}{k} \chi_{k^2}^i ds \\
 & + \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 + 2u^2 M^2 C^2 d + u^2 + 2u^2 G^2 + 2du + \frac{u}{4} K \left( \frac{2d}{4} + \right)^2 \\
 & + \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 + 2u^2 M^2 C^2 d + u^2 + 2u^2 G^2 + 2du + \frac{u}{4} K \left( \frac{2d}{4} + \right)^2 \\
 & + \frac{uB}{4T} + \frac{uKA}{4} + \frac{uKB}{4} \\
 & + \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 + 2u^2 M^2 C^2 d + u^2 + 2u^2 G^2 + 2du + \frac{u}{4} K \left( \frac{2d}{4} + \right)^2 \\
 & + \frac{uKA}{4} + \frac{uKB}{2} \\
 & + \frac{u}{4} 3M^2 K^3 + 2u^2 M^2 E + u^2 + 2u^2 G^2 + 2du + \frac{u}{4} K^2 + \frac{u}{16} K^2 d + \frac{uKA}{4} + \frac{uKB}{2} \\
 & =: C_0 K^3 + C_1 K^2 + C_2 K^2 + C_3 KA + C_4 KB;
 \end{aligned}$$

677 where the constants are defined as

$$\begin{aligned}
 C_0 &= \frac{u}{4} 3M^2 + 2u^2 M^2 E + u^2 + 2u^2 G^2 + 2du \\
 C_1 &= \frac{u}{4} \\
 C_2 &= \frac{u}{16} d \\
 C_3 &= \frac{u}{4} \\
 C_4 &= \frac{u}{2};
 \end{aligned}$$

678 By the weighted CKP inequality and given (1),

$$\begin{aligned}
 W_2(p_K; p_K) &= \frac{q}{D_{KL}(p_K; j; p_K)} + \frac{q}{4 D_{KL}(p_K; j; p_K)} \\
 &= \mathfrak{C}_0^{p-} + \mathfrak{C}_1 \overline{K} + \mathfrak{C}_2^{p-} \overline{K} + \mathfrak{C}_3 \overline{KA} + \mathfrak{C}_4 \overline{KB};
 \end{aligned}$$

679 where the constants are defined as:

$$\begin{aligned}
 \mathfrak{C}_0 &= -\frac{p}{C_0} + \frac{p}{4} \frac{1}{C_0} \\
 \mathfrak{C}_1 &= -\frac{p}{C_1} + \frac{p}{4} \frac{1}{C_1} \\
 \mathfrak{C}_2 &= -\frac{p}{C_2} + \frac{p}{4} \frac{1}{C_2} \\
 \mathfrak{C}_3 &= -\frac{p}{C_3} + \frac{p}{4} \frac{1}{C_3} \\
 \mathfrak{C}_4 &= -\frac{p}{C_4} + \frac{p}{4} \frac{1}{C_4} \\
 \mathfrak{A}^2 &= \max_{n \geq 2} \left\{ \frac{p}{2} \right\};
 \end{aligned}$$

680 From the same analysis of (36), we can have:

$$W_2(p_K; p) = \mathfrak{C}_0^{p-} + \mathfrak{C}_1 \overline{K} + \mathfrak{C}_2^{p-} \overline{K} + \mathfrak{C}_3 \overline{KA} + \mathfrak{C}_4 \overline{KB} + \theta e^{-K} : (67)$$

681 In order to bound the Wasserstein distance, we need to set

$$\epsilon_0^p \overline{K}^2 = \frac{1}{2} \quad \text{and} \quad \epsilon_0 e^{-K} = \frac{1}{2}; \quad (68)$$

682 Solving the equation (68), we can have

$$K = \frac{\log \frac{2}{\epsilon_0}}{\epsilon_0} \quad \text{and} \quad \overline{K} = \frac{2}{4^{-2} \epsilon_0^2 K};$$

683 Combining these two we can have

$$\overline{K} = \frac{2}{4^{-2} \epsilon_0^2 \log \frac{2}{\epsilon_0}} \quad \text{and} \quad K = \frac{4^{-2} \epsilon_0^2 \log^2 \frac{2}{\epsilon_0}}{2 \left( \frac{2}{4^{-2} \epsilon_0^2 \log \frac{2}{\epsilon_0}} \right)^2};$$

684 Plugging in (67) completes the proof.

## 685 D.9 Proof of Theorem 12

686 Recall that the update of VC SGLDLP-L is

$$\begin{aligned} x_{k+1} &= Q^{vc} x_k - Q_G(r \psi(x_k)); 2; \\ &= x_k - Q_G(r \psi(x_k)) + \epsilon_0^p \overline{K} x_k + \epsilon_0^p \overline{K} x_k; \end{aligned}$$

687 where  $\overline{K}$  is defined as

$$\overline{K} = Q^{vc} x_k - Q_G(r \psi(x_k)); 2; \quad x_k - Q_G(r \psi(x_k)) + \epsilon_0^p \overline{K} x_k;$$

688 From analysis in Zhang et al. [2022], we know that

$$\begin{aligned} E_k \left[ \max_{k \leq i \leq k^2} \langle G; \psi \rangle \right] &\leq 5d \\ &=: A; \end{aligned}$$

689 Combining the analysis in section D.6, we can show,

$$\begin{aligned} D_{KL}(\rho_K \| \rho_K) &\leq \frac{M \overline{E}}{4} K^2 + \frac{3M+1}{4} K^2 + \frac{((6+3m_2)M+m_2)d}{16m_2} K^2 + \frac{6M}{4m_2} + \frac{1}{4} K E_k \left[ \max_{k \leq i \leq k^2} \langle G; \psi \rangle \right] \\ &\leq \frac{M \overline{E}}{4} K^2 + \frac{3M+1}{4} K^2 + \frac{((6+3m_2)M+m_2)d}{16m_2} K^2 + \frac{6M}{4m_2} + \frac{1}{4} KA \\ &=: C_0 K^2 + C_1 K^2 + C_2 K^2 + C_3 KA; \end{aligned}$$

690 where the constants  $C_0, C_1, C_2$  and  $C_3$  are defined as:

$$\begin{aligned} C_0 &= \frac{M \overline{E}}{4} \\ C_1 &= \frac{3M+1}{4} \\ C_2 &= \frac{((6+3m_2)M+m_2)d}{16m_2} \\ C_3 &= \frac{6M+m_2}{m_2} \end{aligned}$$

691 We are ready to bound the Wasserstein distance,

$$\begin{aligned} W_2^2(\rho_K; \rho_K) &\leq (12+8(\epsilon_0+2b+2d)) \left( C_0 + \epsilon_0^p \overline{C_0} + C_1 + \epsilon_0^p \overline{C_1} \right) (K)^2 + C_2 + \epsilon_0^p \overline{C_2} (K)^2 \\ &\quad + C_3 + \epsilon_0^p \overline{C_3} AK^2 \\ &=: \epsilon_0^2 + \epsilon_1^2 A + \epsilon_2^2 (K)^2 + \epsilon_3^2 AK^2; \end{aligned}$$

692 where the constants are defined as:

$$\begin{aligned}
 A &= \max\{n^{-2}, p^{-\frac{1}{2}}\} \\
 A &= \max\{A, \frac{1}{A}\} \\
 \mathfrak{C}_0^2 &= (12 + 8(a_0 + 2b + 2d)) C_0 + \frac{p}{C_0} \\
 \mathfrak{C}_1^2 &= (12 + 8(a_0 + 2b + 2d)) C_1 + \frac{p}{C_1} \\
 \mathfrak{C}_2^2 &= (12 + 8(a_0 + 2b + 2d)) C_2 + \frac{p}{C_2} \\
 \mathfrak{C}_3^2 &= (12 + 8(a_0 + 2b + 2d)) C_3 + \frac{p}{C_3} :
 \end{aligned}$$

693 From Proposition 9 in the paper Raginsky et al. [2017], we know that

$$\begin{aligned}
 W_2(\mathbb{P}_K; p) &\leq 2C_{LS} \log kp_0 k_1 + \frac{d}{2} \log \frac{3}{m} + \frac{M}{3} + B \frac{p}{C_0} + G_0 + \frac{b}{2} \log 3 e^{K=C_{LS}} \\
 &=: \mathfrak{C}_4 e^{K=C_{LS}}
 \end{aligned}$$

694 Finally, we can have

$$W_2(\mathbb{P}_K; p) \leq \mathfrak{C}_0 p^{-\frac{1}{2}} + \mathfrak{C}_1 \frac{p}{A} + \mathfrak{C}_2 p^{-\frac{1}{2}} K + \mathfrak{C}_3 \frac{p}{A} \frac{1}{K^2} + \mathfrak{C}_4 e^{K=C_{LS}} : \quad (69)$$

695 In order to bound the 2-Wasserstein distance, we need to set

$$\mathfrak{C}_0 K^{-\frac{5}{4}} = \frac{1}{2} \quad \text{and} \quad \mathfrak{C}_3 e^{K=C_{LS}} = \frac{1}{2} : \quad (70)$$

696 Solving the (70), we can have

$$K = C_{LS} \log \frac{2\mathfrak{C}_3}{\mathfrak{C}_0} \quad \text{and} \quad K = \frac{4}{16\mathfrak{C}_0^4 (K)^4} :$$

697 Combining these two we can have

$$= \frac{4}{16\mathfrak{C}_0^4 C_{LS}^4 \log^4 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}} \quad \text{and} \quad K = \frac{16\mathfrak{C}_0^4 C_{LS}^5 \log^5 \frac{2\mathfrak{C}_3}{\mathfrak{C}_0}}{4} :$$

698 Plugging  $K$  and  $\mathfrak{C}_4$  into (69) completes the proof.

## 699 E Technical Proofs

### 700 E.1 Proof of Lemma 13

701 Proof. By the definition of  $\mathfrak{C}_i$  in (25)

$$\begin{aligned}
 \mathbb{E} \|k\|^2 &= \mathbb{E} \|g(x) - \mathbb{E} r U(x)\|^2 \\
 &= \mathbb{E} \left\| \mathbb{E}_h r U(Q_w(x)) - \mathbb{E} r U(x) \right\|_i^2 \\
 &= \mathbb{E} \left\| \mathbb{E}_h (k r U(Q_w(x)) - r U(x)) \right\|_i^2 \\
 &= M^2 \mathbb{E} \|k Q_w(x) - r U(x)\|^2 \\
 &= M \frac{2d}{4} :
 \end{aligned}$$

702 We also know that from the definition that

$$\begin{aligned}
 E \|k\|^2 &= E \|g(x) - U(x)\|^2 \\
 &= E \|Q_G(r - \theta(Q_W(x))) - \theta(Q_W(x)) + r - \theta(Q_W(x)) - U(Q_W(x)) + r - U(Q_W(x)) - U(x)\|^2 \\
 &= E \|Q_G(r - \theta(Q_W(x))) - \theta(Q_W(x))\|^2 + E \|r - \theta(Q_W(x)) - U(Q_W(x))\|^2 + E \|r - U(Q_W(x)) - U(x)\|^2 \\
 &= \frac{2d}{4} + 2 + M^2 E \|Q_W(x) - x\|^2 \\
 &= (M^2 + 1) \frac{2d}{4} + 2;
 \end{aligned}$$

703 where in the first inequality, we apply Assumptions 1 and 4.

704

□

705 E.2 Proof of Lemma 14

706 Proof. Let  $\pi_1$  be the set of all couplings between  $q_0$  and  $q$  and  $\pi_2$  be the set of all couplings  
 707 between  $q_0$  and  $q$ . Let  $r_1$  be the optimal coupling between  $q_0$  and  $q$ , i.e.

$$E(\cdot; \cdot)_{r_1}[\|k\|^2] = W_2^2(q_0; q):$$

708 Let  $x_!; x_!$   $r_1$ . We define the random variable  $x_!$  as

$$x_! = x_! + u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr + \int_0^R e^{-s} ds :$$

709 By equation (29),  $x_!; x_!$  define a valid coupling between  $q_0$  and  $q$ . Now we can analyze  
 710 the Wasserstein distance between  $q_0$  and  $q$ .

$$\begin{aligned}
 W_2^2(q_0; q) &= E_{r_1} \left[ \|x_! + u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr + \int_0^R e^{-s} ds - x_!\|^2 \right] \\
 &= E_{r_1} \left[ \|x_! - x_!\|^2 + 2u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr E \left[ \|x_! - x_!\| \int_0^R e^{-s} ds \right] \right. \\
 &\quad \left. + E_{r_1} \left[ u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr \int_0^R e^{-s} ds \right] \right] \\
 &= W_2^2(q_0; q) + 2u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr E \left[ \|x_! - x_!\| \int_0^R e^{-s} ds \right] \\
 &= W_2^2(q_0; q) + 2u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr \left( \frac{2d}{4} + 2 \right) E_{r_1} \|k\|^2 \\
 &= W_2^2(q_0; q) + 2u \int_0^R \int_0^{R_r} e^{-(s+r)} ds dr (M^2 + 1) \frac{2d}{4} + 2 :
 \end{aligned} \tag{71}$$

711

□

712 E.3 Proof of Lemma 15

713 Proof. In order to get the upper bound  $k_k$  and  $v_k$ , we bound the Lyapunov function  
 714  $E(x_k; v_k)$ . By the smooth Assumption 1, we know

$$U(x_{k+1}) - U(x) = U(x_k) + h r U(x_k); x_{k+1} - x_k + M^2 \|x_{k+1} - x_k\|^2 - U(x) :$$

715 Recall the definition of the Lyapunov function

$$E(x_{k+1}; v_{k+1}) = \|x_{k+1}\|^2 + \|v_{k+1}\|^2 = \|k\|^2 + 8u(U(x_{k+1}) - U(x)) = 2:$$



716 For the first two terms we have

$$\begin{aligned} kx_{k+1} + 2v_{k+1} &= kx_k + 2v_k + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) \\ &= kx_k + 2v_k + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) \\ &= kx_k + 2v_k + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) + \frac{h}{2} \left( \frac{1}{k^2} - \frac{1}{(k+1)^2} \right) \end{aligned}$$

717 This implies the following:

$$\begin{aligned} E[E(x_{k+1}; v_{k+1})] &= E[E(x_k; v_k)] + 4E[hx_k; x_{k+1} - x_k] + \frac{4}{2}E[hv_k; v_{k+1} - v_k] + \frac{4}{2}E[hv_k; x_{k+1} - x_k] \\ &+ \frac{8}{2}E[hv_k; v_{k+1} - v_k] + \frac{8u}{2}E[hv_k; r U(x_k); x_{k+1} - x_k] + M=2kx_{k+1} - x_k k^2 \\ &+ E[kx_{k+1} - x_k k^2] + E[kv_{k+1} - v_k k^2] : \end{aligned} \tag{72}$$

718 By the update rule in (3), we know that

$$\begin{aligned} E[hx_k; x_{k+1} - x_k] &= \frac{1}{2} \frac{e^{-u}}{1+e^{-u}} E[hx_k; v_k] + \frac{u(1+e^{-u})}{2} E[hx_k; g(x_k)]; \\ E[hx_k; v_{k+1} - v_k] &= (1 - e^{-u}) E[hx_k; v_k] - \frac{u(1 - e^{-u})}{2} E[hx_k; g(x_k)]; \\ E[hv_k; x_{k+1} - x_k] &= \frac{1}{2} \frac{e^{-u}}{1+e^{-u}} E[hv_k k^2] + \frac{u(1+e^{-u})}{2} E[hv_k; g(x_k)]; \\ E[hv_k; v_{k+1} - v_k] &= (1 - e^{-u}) E[hv_k k^2] - \frac{u(1 - e^{-u})}{2} E[hv_k; g(x_k)]; \end{aligned}$$

719 Plug into the (72) yields:

$$\begin{aligned} E[E(x_{k+1}; v_{k+1})] &= E[E(x_k; v_k)] + \frac{4u(2 - 2e^{-u})}{2} E[hx_k; g(x_k)] - \frac{4(1 - e^{-u})}{2} E[hv_k k^2] \\ &+ \frac{4u(1+e^{-u})}{3} E[hv_k; g(x_k)] + \frac{8u(1 - e^{-u})}{3} E[hv_k; r U(x_k); g(x_k)] \\ &+ \frac{8u^2(1+e^{-u})}{4} E[hv_k; g(x_k)] + \frac{4Mu}{2} + 3 E[kx_{k+1} - x_k k^2] \\ &+ \frac{8}{2} E[hv_{k+1} - v_k k^2] : \end{aligned} \tag{73}$$

720 By Assumption 3, we know that  $h \leq \frac{1}{2} k^2$ . We then assume  $1 = (8)$  and use the inequality  $x - e^{-x} \leq 1 - x^2$  for any  $x \geq 0$ , it follows that

$$\begin{aligned} &\frac{4u(2 - 2e^{-u})}{2} E[hx_k; g(x_k)] \\ &= \frac{4u(2 - 2e^{-u})}{2} (E[hx_k; r U(x_k)] + E[hx_k; g(x_k) - r U(x_k)]) \\ &\leq \frac{4u(2 - 2e^{-u})}{2} \left( \frac{h}{2} E[kx_k k^2] + b + \frac{4u(2 - 2e^{-u})}{2} \left( \frac{1}{8} E[kx_k k^2] + 2E[kg(x_k) - r U(x_k)] k^2 \right) \right) \\ &\leq \frac{3m_2 u}{2} E[kx_k k^2] + \frac{4u b}{2} + \frac{8u}{2} E[kg(x_k) - r U(x_k)] k^2 ; \end{aligned}$$

722 where the first inequality is because of the Young's inequality and Assumption 1 and the last in-  
723 equality is based on the inequality that  $(x - e^{-x})^2 \leq 2 - 2e^{-x}$ . Again by Young's  
724 inequality and the update rule in (3) we have:

$$\begin{aligned} E[kx_{k+1} - x_k k^2] &\leq 2 E[hv_k k^2] + u^2 = 2E[kg(x_k)] k^2 + E[kx_k k^2] \\ E[kv_{k+1} - v_k k^2] &\leq 2 E[hv_k k^2] + 2u^2 = 2E[kg(x_k)] k^2 + E[kv_k k^2] : \end{aligned}$$

725 It is easy to verify the fact that  $E[k_k^2] \leq 2ud$  and  $E[k_k^2] \leq 2ud^2$ . Thus,

$$E[E(x_{k+1}; v_{k+1})] = E[E(x_k; v_k)] + \frac{3um^2}{2} E[k_k^2] + \frac{3(1-e^{-u})}{2} (8Mu + u^2 + 22u^2) E[k_k^2] + \frac{36u^2 + 2u^2 + 4Mu + 3}{2} E[kg(x_k)k^2] + \frac{2u^2}{2} E[kr U(x_k)k^2] + \frac{8u(2+2u)}{3} E[kr U(x_k)k^2] + \frac{(8Mu + 6^2)ud^2 + 4(4d+b)u}{2} :$$

726 If we set  $\min \left( \frac{s}{4(8Mu + u^2 + 22u^2)}; \frac{4u^2}{4Mu + 3^2}; \frac{6bu}{(4Mu + 3^2)d} \right)$ ;

727 we can obtain the following,

$$E[E(x_{k+1}; v_{k+1})] = E[E(x_k; v_k)] + \frac{3um^2}{2} E[k_k^2] + \frac{2}{2} E[k_k^2] + \frac{(20u + )u^2}{2} E[kg(x_k)k^2] + \frac{2u^2}{2} E[kr U(x_k)k^2] + \frac{8u(2+2u)}{3} E[kr U(x_k)k^2] + \frac{16(d+b)u}{2} : \quad (74)$$

728 Furthermore we can bound  $E[kg(x_k)k^2]$  by the following analysis:

$$E[kg(x_k)k^2] \leq 2E[kg(x_k)kr U(x_k)k^2] + 2E[kr U(x_k)k^2] \leq 2(M^2 + 1) \frac{2d}{4} + u^2 + 4M^2 E[k_k^2] + 4G^2; \quad (75)$$

729 where  $G^2$  is the bound of the gradient  $\| \nabla_{\theta} \text{kr U}(0) \|^2 = G^2$ . Thus we can have:

$$E[E(x_{k+1}; v_{k+1})] = E[E(x_k; v_k)] + \frac{3um^2}{2} E[k_k^2] + \frac{2}{2} E[k_k^2] + \frac{(21u + )4M^2u^2}{2} E[k_k^2] + \frac{2(20u + )u^2 + 8u(2+2u)}{2} (M^2 + 1) \frac{2d}{4} + u^2 + \frac{(21u + )4u^2}{2} G^2 + \frac{16(d+b)u}{2} :$$

730 If we set the stepsize

$$\min \left( \frac{m^2}{12(21u + )M^2}; \frac{8(2+2u)}{(20u + )} \right) ;$$

731 then we have:

$$E[E(x_{k+1}; v_{k+1})] = E[E(x_k; v_k)] + \frac{8um^2}{3} E[k_k^2] + \frac{2}{2} E[k_k^2] + \frac{16u(2+2u)}{3} (M^2 + 1) \frac{2d}{4} + u^2 + \frac{(21u + )4u^2}{2} G^2 + \frac{16(d+b)u}{2} :$$

732 Furthermore by Young's inequality and Assumption 1, we can bound the Lyapunov function by the following:

$$E(x; v) \leq 2kxk^2 + \frac{12}{2} + \frac{2uM}{2} - 3kxk^2 + 6kxk^2 :$$

734 Then if  $2 \leq 4Mu$ , we have

$$E(x; v) = \frac{16uM}{2} kxk^2 + \frac{12}{2} kvk^2 + \frac{12uM}{2} kxk^2: \quad (76)$$

735 Thus,

$$E[E(x_{k+1}; v_{k+1})] = 1 - \frac{m_2}{6M} E[E(x_k; v_k)] + \frac{16u}{3} \frac{2+2u}{3} (M^2+1) \frac{2d}{4} + 2 + \frac{(21u+4)4u^2}{2} G^2 + \frac{16(d+b)u}{2}:$$

736 Finally we show that

$$\begin{aligned} \sup_k E[E(x_k; v_k)] &= E[E(x_0; v_0)] + \frac{6M}{m_2} \frac{16u}{3} \frac{2+2u}{3} (M^2+1) \frac{2d}{4} + 2 \\ &+ \frac{6M}{m_2} \frac{(21u+4)4u^2}{2} G^2 + \frac{6M}{m_2} \frac{16(d+b)u}{2} \\ &= E[E(x_0; v_0)] + \frac{96u}{m_2^4} \frac{2+2u}{3} (M^2+1) \frac{2d}{4} + 2 + \frac{24(21u+4)uM}{m_2^3} G^2 + \frac{96(d+b)uM}{m_2^2} \\ &= \bar{E} + C_0 (M^2+1) \frac{2d}{4} + 2; \end{aligned} \quad (77)$$

737 where  $\bar{E} = E[E(x_0; v_0)] + \frac{24(21u+4)uM}{m_2^3} G^2 + \frac{96(d+b)uM}{m_2^2}$  and  $C_0 = \frac{96u(2+2u)}{m_2^4}$ . Moreover by the  
738 definition of Laypunov function, we know  $E(x; v) = \max\{kxk^2; 2kv = k^2g\}$ . This further implies  
739 that

$$\begin{aligned} E[kx_k k^2] &\leq \bar{E} + C_0 (M^2+1) \frac{2d}{4} + 2 \\ E[kv_k k^2] &\leq 2\bar{E} + 2C_0 (M^2+1) \frac{2d}{4} + 2: \end{aligned}$$

740 Combining with equation (75) we can bound  $E[kg(x_k)k^2]$  as:

$$E[kg(x_k)k^2] \leq 2(M^2+1) \frac{2d}{4} + 2 + 4M^2\bar{E} + 4G^2: \quad (78)$$

741

□

#### 742 E.4 Proof of Lemma 16

743 Proof. By the update rule in (18), we have:

$$\begin{aligned} E[kx_{k+1}k^2] &= E[kx_k g(x_k)k^2 + \frac{p}{8} E[hx_k g(x_k); k_{k+1}i] + 2 E[k_{k+1}k^2] \\ &= E[kx_k g(x_k)k^2 + 2d \\ &= E[kx_k (r U(x_k) - (g(x_k) - r U(Q_W(x_k)))) (r U(Q_W(x_k)) - r U(x_k))k^2 + 2d \\ &= E[kx_k (r U(x_k) - (r U(Q_W(x_k)) - r U(x_k)))k^2 + 2E[kg(x_k) - r U(Q_W(x_k))]k^2 + 2d \\ &= (E[kx_k (r U(x_k) - r U(Q_W(x_k)))k^2 + E[kr U(Q_W(x_k)) - r U(x_k)]k^2) + \frac{2d}{4} + 2d: \end{aligned}$$

744 We know the fact that:

$$\begin{aligned} E[kx_k (r U(x_k) - r U(Q_W(x_k)))k^2] &= E[kx_k k^2] - 2 E[hx_k; r U(x_k)i] + 2 E[kr U(x_k)k^2] \\ &= E[kx_k k^2] + 2b - m_2 E[kx_k k^2] + 2M^2 E[kx_k k^2] + G^2 \\ &= (1 - 2m_2 + 2M^2) E[kx_k k^2] + 2b + 2G^2: \end{aligned}$$

745 For any  $2 > 0; 1 \wedge \frac{m_2}{2M^2}$ , if  $0 < 1 - 2m_2 + 2^2M^2 < 1$  and set  $c = \frac{m_2 + 2M^2}{1 - 2m_2 + 2^2M^2}$ , then we  
 746 have:

$$\begin{aligned} E kx_{k+1} k^2 &= (1+c) E kx_k k^2 + \frac{1}{c} E kr U(x_k) k^2 + \frac{2d}{4} + 2d \\ &= \frac{1}{m_2 + 2M^2} E kx_k k^2 + \frac{1}{m_2 + 2M^2} \frac{m_2 + 2M^2}{4} \frac{2d}{4} + \frac{1}{1 - 2m_2 + 2^2M^2} (2b + 2^2G^2) \\ &\quad + \frac{2d}{4} + 2d: \end{aligned}$$

747 For any  $\gamma > 0$  we can bound the recursive equations as:

$$\begin{aligned} E kx_k k^2 &= E kx_0 k^2 + \frac{1}{2(m_2 + M^2)^2} \frac{m_2 + 2M^2}{4} \frac{2d}{4} + \frac{1}{(1 - 2m_2 + 2^2M^2)(m_2 + M^2)} (2b + 2^2G^2) \\ &\quad + \frac{1}{(m_2 + M^2)} \frac{2d}{4} + 2d \\ &= E kx_0 k^2 + \frac{1}{(m_2 + M^2)^2} \frac{m_2 + 2M^2}{4} \frac{2d}{4} + \frac{1}{(1 - 2m_2 + 2^2M^2)(m_2 + M^2)} (2b + 2^2G^2) \\ &\quad + \frac{1}{m_2 + M^2} \frac{2d}{4} + 2d \\ &= E kx_0 k^2 + \frac{2M^2}{m_2} \frac{2d}{4} + \frac{2}{m_2} (2b + 2^2G^2) + \frac{2}{m_2} \frac{2d}{4} + 2d: \end{aligned}$$

748 Now if we let  $E = E kx_0 k^2 + \frac{M}{m_2} (2b + 2^2G^2) + 2d$ , then we can write:

$$E kx_k k^2 = E + \frac{2M^2 + 1}{m_2} \frac{2d}{4}:$$

749

□

## 750 E.5 Proof of Lemma 17

751 Proof. From the same analysis in (74), if we set

$$\min \left( \frac{s}{4(8Mu + u^2 + 22^2)}; \frac{4u^2}{4Mu + 3^2}; \frac{6bu}{(4Mu + 3^2)d} \right);$$

752 we can obtain the following,

$$\begin{aligned} E[E(x_{k+1}; v_{k+1})] &= E[E(x_k; v_k)] + \frac{3um_2}{2} E kx_k k^2 + \frac{2}{2} E kv_k k^2 + \frac{(20u + )u^2}{2} E Q_G(r U(x_k))^2 \\ &\quad + \frac{2u^2}{2} E kr U(x_k) k^2 + \frac{8u^2 + 2u}{3} E r U(x_k) Q_G(r U(x_k))^2 + \frac{16(d+b)u}{4}: \end{aligned} \tag{79}$$

753 By assumption 1, we can bound  $E Q_G(r U(x_k))^2$  by the following,

$$\begin{aligned} E kQ_G(r U(x_k)) k^2 &= E Q_G(r U(x_k))^2 + E Q_G(r U(x_k))^2 + 2E kr U(x_k) r U(0) k^2 + 2E kr U(0) k^2 \\ &\quad + \frac{2d}{4} + 2 + 2M^2 E kx_k k^2 + 2G^2: \end{aligned}$$

754 Plugging this bound into equation 79, we can have:

$$\begin{aligned}
\mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] + \frac{2(20u + \gamma)u\eta^2 M^2}{\gamma^2} \mathbb{E} [\|\mathbf{x}_k\|^2] \\
&\quad + \frac{(20u + \gamma)u\eta^2}{\gamma^2} \left( \frac{\Delta^2 d}{4} + \sigma^2 + 2G^2 \right) + \frac{2u^2\eta^2}{\gamma^2} \left( 2M^2 \mathbb{E} [\|\mathbf{x}_k\|^2] + 2G^2 \right) \\
&\quad + \frac{8u\eta(\gamma^2 + 2u)}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) + \frac{16(d+b)u\eta}{\gamma} \\
&\leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} \mathbb{E} [\|\mathbf{x}_k\|^2] \\
&\quad + \frac{(20u + \gamma)\gamma u\eta^2 + 8(\gamma^2 + 2u)u\eta}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} G^2 + \frac{16(d+b)u\eta}{\gamma} \\
&\leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} \mathbb{E} [\|\mathbf{x}_k\|^2] \\
&\quad + \frac{(36u + 9\gamma^2)u\eta}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} G^2 + \frac{16(d+b)u\eta}{\gamma}.
\end{aligned}$$

755 If we set the step size  $\eta \leq \frac{m_2}{6(22u + \gamma)M^2}$ , we can have:

$$\begin{aligned}
\mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{8um_2\eta}{3\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] \\
&\quad + \frac{(36u + 9\gamma^2)u\eta}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} G^2 + \frac{16(d+b)u\eta}{\gamma}.
\end{aligned}$$

756 Again from the same analysis in (76), if  $\gamma^2 \leq 4Mu$ , we have

$$\mathcal{E}(x, v) \leq \frac{16uM}{\gamma^2} \|x\|^2 + \frac{12}{\gamma^2} \|v\|^2 + \frac{12uM}{\gamma^2} \|x^*\|^2.$$

757 Thus,

$$\begin{aligned}
\mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \left( 1 - \frac{\gamma m_2 \eta}{6M} \right) \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] + \frac{(36u + 9\gamma^2)u\eta}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
&\quad + \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} G^2 + \frac{16(d+b)u\eta}{\gamma}.
\end{aligned}$$

758 Finally, we show that for any  $k > 0$ ,

$$\begin{aligned}
\mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] &\leq \mathbb{E} [\mathcal{E}(x_0, v_0)] + \frac{6M}{\gamma m_2 \eta} \frac{(36u + 9\gamma^2)u\eta}{\gamma^3} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
&\quad + \frac{6M}{\gamma m_2 \eta} \frac{2(22u + \gamma)u\eta^2 M^2}{\gamma^2} G^2 + \frac{6M}{\gamma m_2 \eta} \frac{16(d+b)u\eta}{\gamma} \\
&\leq \mathbb{E} [\mathcal{E}(x_0, v_0)] + \frac{54(4u + \gamma^2)u}{m_2 \gamma^4} \left( \frac{\Delta^2 d}{4} + \sigma^2 \right) + \frac{12(22u + \gamma)uM^3}{m_2 \gamma^3} G^2 + \frac{96(d+b)uM}{m_2 \gamma^2} \\
&=: \mathcal{E} + C\Delta^2 d.
\end{aligned}$$

759 Finally by the fact that  $\mathbb{E} [\|\mathbf{x}_k\|^2] \leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)]$  and  $\mathbb{E} [\|\mathbf{v}_k\|^2] \leq \gamma^2 \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] / 2$  we can  
760 get our claim in Lemma 17.

761

□

## 762 F Additional Experiment Results

763 In this section, we provide additional experiment results.

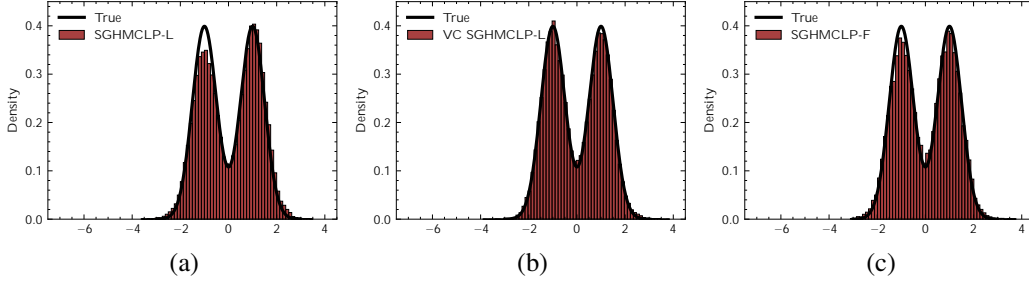


Figure 5: Low-precision SGHMC with stepsize equal to 0.01 on a Gaussian mixture distribution. (a): SGHMCLP-L. (b): VC SGHMCLP-L. (c): SGHMCLP-F.

764 **F.1 Sampling from Gaussian mixture Distribution**

765 We first demonstrate the performance of Low-precision SGHMC for fitting a strongly log-concave  
 766 distribution. In this case, we use the standard Gaussian distribution as the representative of the  
 767 strongly log-concave distribution. The simulation result is shown in Figure 1. As in the Figure 1  
 768 and 5 displayed, the sample obtained from naïve SGHMCLP-L has a larger variance than the target  
 769 distribution. This verifies the results we prove in Theorem 6 and 7. This is because in addition to the  
 770 Gaussian noise the naïve quantizer in order to be unbiased introduces an extra noise which increases  
 771 the variance of the sample. The variance corrected quantizer solves this problem by quantizing the  
 772 mean of each sample and letting the variance of the quantizer equal to the variance  $\text{Var}_x^{hmc}$  de-  
 773 fined by the Hamiltonian dynamics 9. The variance-corrected SGHMC with low-precision gradient  
 774 accumulators (VC SGHMCLP-L) doesn't suffer from the larger variance problem as the variance  
 775 corrected quantization matches the variance defined in (2).

776 We also study in which case the variance corrected  
 777 quantization function is advantageous  
 778 over the naïve stochastic quantization func-  
 779 tion. We test the 2-Wasserstein distance of VC  
 780 SGHMCLP-L and SGHMCLP-L over different  
 781 variances. The result is shown in Figure 4. We  
 782 found that when the variance  $\text{Var}_x^{hmc}$  is close  
 783 to the largest quantization variance  $\Delta^2/4$ , the  
 784 variance corrected quantization function shows  
 785 the largest advantage over the naïve quantiza-  
 786 tion. When the variance  $\text{Var}_x^{hmc}$  is less than  
 787  $\Delta^2/4$  the correction has a chance to fail and  
 788 when it is 100 times the quantization variance,  
 789 the advantage of variance corrected quantiza-  
 790 tion shows less advantage. One possible reason  
 791 is the quantization noise eliminated by variance  
 792 corrected quantization function is not critical  
 793 compared with the intrinsic variance needed.

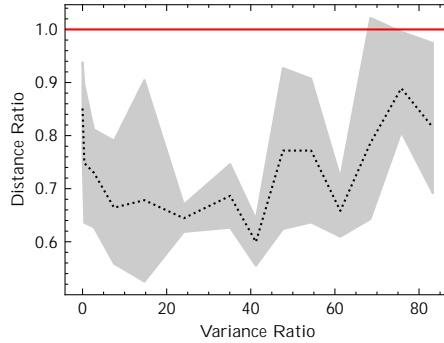


Figure 4: Wasserstein Distance Ratio of VC SGHMCLP-L & SGHMCLP-L (Smaller is better). The dashed line is the 2-Wasserstein distance to the target distribution ratio between the sample obtained by VC SGHMCLP-L and SGHMCLP-L.

794 **F.2 Multi-layer perception**

795 We present the low-precision SGHMC with MLP on the MNIST dataset in Figure 6. We observe  
 796 similar results as the low-precision SGHMC with the logistic model.

797 **F.3 CIFAR-10 & CIFAR-100**

798 In this section, we present some additional results for experiments on computer vision tasks in  
 799 CIFAR datasets.

