
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 ϵ -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:

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

63 where u, γ 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 η is the step size. $\xi_k^{\mathbf{x}}$ are normal distributed in \mathbb{R}^d satisfying
 65 that :

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

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 Δ . 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 σ^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 Δ and con-
 91 centration parameters (λ^*, μ^*) but hides multiplicative terms that depend polynomially on the other
 92 parameters (e.g., dimension d , friction γ , inverse mass u and gradients variance σ^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 ϵ and step size η , 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 μ^* 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 ϵ and step size η , 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 ϵ term and a $\log(\epsilon^{-1})$ term. The
 113 difference is that for the SGHMCLP-F algorithm, the quantization error Δ affects the multiplicative
 114 constant of the $\log(\epsilon^{-1})$ term. Without Δ , one can choose a small ϵ and a larger batch size (i.e., a
 115 smaller σ^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. $\tilde{\mathcal{O}}\left(\epsilon + (\log(\epsilon^{-1})\Delta)^{1/2}\right)$) ver-
 119 sus $\tilde{\mathcal{O}}\left(\epsilon + \log(\epsilon^{-1})\Delta^{1/2}\right)$) distance for non-log-concave target distribution within fewer iterations
 120 (i.e., $\tilde{\mathcal{O}}\left(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1})\right)$ versus $\tilde{\mathcal{O}}\left(\epsilon^{-4}\lambda^{*-1}\log^5(\epsilon^{-1})\right)$).

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. ϵ and achieves better upper bound
 129 $\tilde{\mathcal{O}}\left(\epsilon^{-2}\log^{3/2}(\epsilon^{-1})\Delta^{1/2}\right)$ versus $\tilde{\mathcal{O}}\left(\epsilon^{-4}\log^5(\epsilon^{-1})\Delta^{1/2}\right)$.

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 $\text{Var}_{\mathbf{v}}^{hmc} = u(1 - e^{-2\gamma\eta})$
 138 and $\text{Var}_{\mathbf{x}}^{hmc} = u\gamma^{-2}(2\gamma\eta + 4e^{-\gamma\eta} - e^{-2\gamma\eta} - 3)$, the VC SGHMCLP-L can be done as:

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

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, dissativity 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 ϵ and step size η , the K -th iteration of the VC*
 143 *SGHMCLP-L update (4), i.e., \mathbf{x}_K , satisfies*

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

144 for some K satisfying

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

145 Comparing with Theorem 2, the variance corrected quantization can improve the upper bound
 146 w.r.t. ϵ from $\tilde{\mathcal{O}}\left(\epsilon^{-2}\log^{3/2}(\epsilon^{-1})\Delta^{1/2}\right)$ to $\tilde{\mathcal{O}}\left(\epsilon^{-1}\log(\epsilon^{-1})\Delta^{1/2}\right)$. 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. ϵ and achieves better upper bounds ($\tilde{\mathcal{O}}\left(\epsilon + \log(\epsilon^{-1})\epsilon^{-1}\Delta^{1/2}\right)$ versus
 150 $\tilde{\mathcal{O}}\left(\epsilon + \log^3(\epsilon^{-1})\epsilon^{-2}\Delta^{1/2}\right)$).

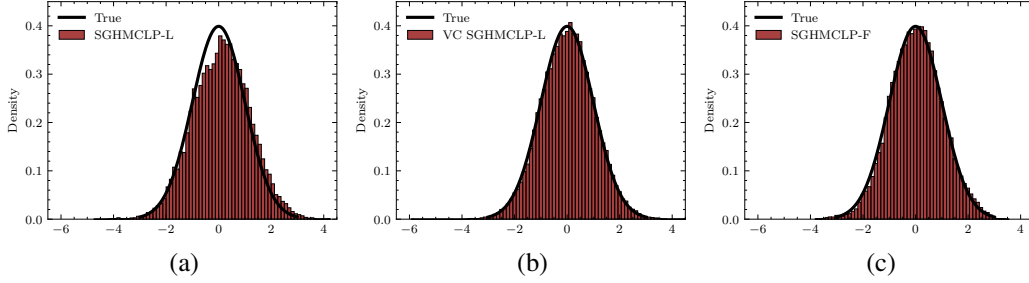


Figure 1: Low-precision SGHMC on Gaussian distribution. (a): SGHMCLP-L. (b): VC SGHMCLP-L. (c): SGHMCLP-F.

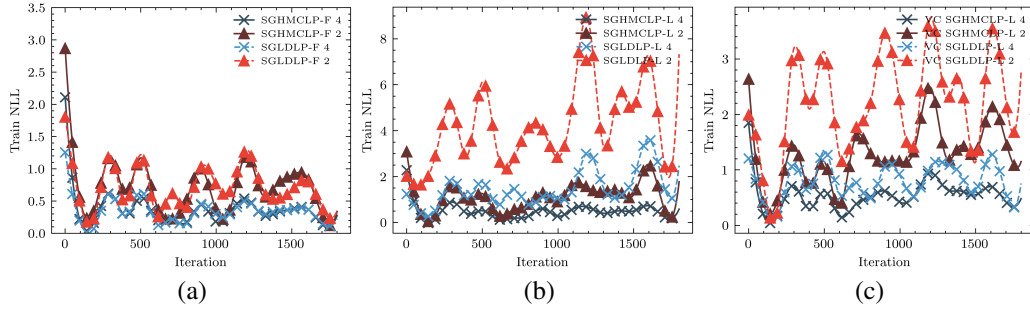


Figure 2: Training NLL of low-precision SGHMC and SGLD on logistic model with MNIST in terms of different numbers of fractional bits. (a): Methods with full-precision gradient accumulators. (b): Methods with low-precision gradient accumulators. (c): Variance corrected quantization.

151 Interestingly, the naïve SGHMCLP-L has similar dependence on the quantization error Δ 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, *qtorch* [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 fit 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 Δ . 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.

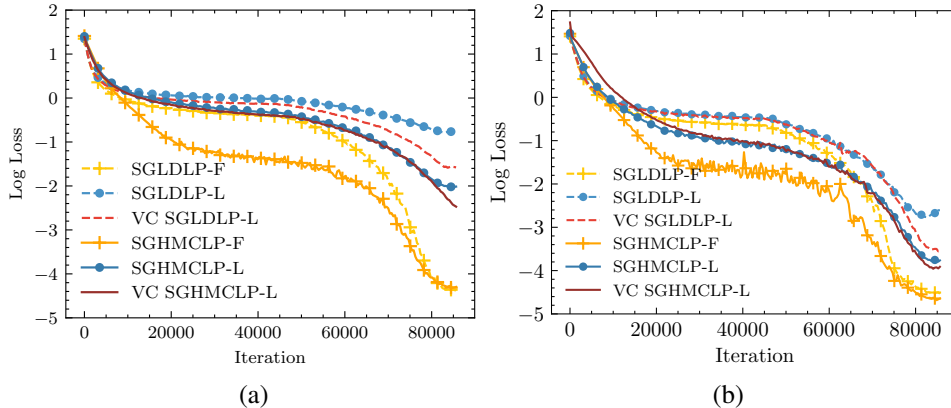


Figure 3: Log of training NLL of low-precision SGHMC and SGLD on ResNet-18 with CIFAR100 and constant step sizes. (a): 8-bit Fixed Point. (b): 8-bit Block Float Point.

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 \pm 0.10	4.52 \pm 0.07	4.78 \pm 0.08	5.19 \pm 0.09	5.07 \pm 0.04	5.08 \pm 0.08	4.75 \pm 0.21	4.58 \pm 0.07	4.93 \pm 0.09
CIFAR-100	22.34 \pm 0.22	22.40 \pm 0.04	22.37 \pm 0.04	23.71 \pm 0.18	23.36 \pm 0.10	23.54 \pm 0.10	22.86 \pm 0.14	22.70 \pm 0.22	22.39 \pm 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 SGLDLP-F and the
 183 naive SGHMCLP-L significantly outperforms naive SGLDLP-L and SGLDLP-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 naive
 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 \pm 0.22	7.81 \pm 0.07	7.03 \pm 0.23	6.63 \pm 0.01	6.60 \pm 0.06	5.86 \pm 0.18	5.75 \pm 0.05	5.51 \pm 0.01	5.38 \pm 0.06	5.15 \pm 0.08
CIFAR-100	28.42 \pm 0.35	27.15 \pm 0.35	26.73 \pm 0.12	26.57 \pm 0.10	26.43 \pm 0.19	26.75 \pm 0.11	26.11 \pm 0.38	25.14 \pm 0.11	25.29 \pm 0.03	24.45 \pm 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*, 31, 2018.
- 203 F. Bolley and C. Villani. Weighted csiszár-kullback-pinsker inequalities and applications to trans-
204 portation inequalities. In *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. In *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. In *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 gradient. *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. Ré. 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. In *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. In *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. In *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, Ł. 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*, 31, 2018.
- 245 M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Pro-*
246 *ceedings of the 28th international conference on machine learning (ICML-11)*, 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*, 32, 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:

$$\begin{aligned} d\mathbf{v}_t &= -\gamma\mathbf{v}_t dt - u\nabla U(\mathbf{x}_t)dt + \sqrt{2\gamma u}d\mathbf{B}_t \\ d\mathbf{x}_t &= \mathbf{v}_t dt. \end{aligned} \quad (9)$$

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

$$Q^s(\theta) = \begin{cases} \Delta \lfloor \frac{\theta}{\Delta} \rfloor, & \text{w.p. } \lfloor \frac{\theta}{\Delta} \rfloor - \frac{\theta}{\Delta} \\ \Delta \lceil \frac{\theta}{\Delta} \rceil, & \text{w.p. } 1 - (\lceil \frac{\theta}{\Delta} \rceil - \frac{\theta}{\Delta}). \end{cases} \quad (10)$$

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*

$$\|\nabla U(\mathbf{x}) - \nabla U(\mathbf{y})\|^2 \leq M^2 \|\mathbf{x} - \mathbf{y}\|^2, \quad \text{for any } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

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,*

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

270

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

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

272

273 **Assumption 4 (Bounded Variance).** *There exists a constant $\sigma^2 > 0$, such that the following holds*

$$\mathbb{E} \left\| \nabla \tilde{U}(\mathbf{x}) - \nabla U(\mathbf{x}) \right\|^2 \leq \sigma^2, \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

274

275 Beyond the above assumptions, we further define $\kappa_1 = M/m_1$ and $\kappa_2 = M/m_2$ as the condition
 276 number for strongly log-concave and non-log-concave target distribution respectively, and denote the
 277 global minimum of $U(\mathbf{x})$ as \mathbf{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 $\|\mathbf{x}^*\|^2 < \mathcal{D}^2$. Fur-*
 283 *thermore, let p^* denote the target distribution of \mathbf{x} and \mathbf{v} . Given any sufficiently small ϵ , if we set*
 284 *the step size to be*

$$\eta = \min \left\{ \frac{\epsilon \kappa_1^{-1}}{\sqrt{479232/5(d/m_1 + \mathcal{D}^2)}}, \frac{\epsilon^2}{1440\kappa_1 u^2 [(M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2]} \right\},$$

285 *then after K steps starting with initial points $\mathbf{x}_0 = \mathbf{v}_0 = 0$, the output $(\mathbf{x}_K, \mathbf{v}_K)$ of the SGHMCLP-*
 286 *F in (3) satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) \leq \tilde{\mathcal{O}}(\epsilon + \Delta),$$

287 *for some K satisfying*

$$K \leq \frac{\kappa_1}{\eta} \log \left(\frac{36 \left(\frac{d}{m_1} + \mathcal{D}^2 \right)}{\epsilon} \right) = \tilde{\mathcal{O}}(\epsilon^{-2} \log(\epsilon^{-1}) \Delta^2).$$

288 **Theorem 5.** Suppose Assumptions 1, 3 and 4 hold. Furthermore, let p^* denote the target distribution
 289 of \mathbf{x} and \mathbf{v} . Given initialization $\mathbf{x}_0 = \mathbf{v}_0 = 0$ and $\gamma^2 \leq 4Mu$, for any sufficiently small ϵ , if we set
 290 the step size to be $\eta = \tilde{\mathcal{O}}\left(\frac{\mu^* \epsilon^2}{\log(1/\epsilon)}\right)$ and also satisfy

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

291 then, the K -th iteration of the SGHMCLP-F update (3), i.e., \mathbf{x}_K and \mathbf{v}_K , satisfies

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

292 for some K satisfying

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

293 where constants are defined as: $\tilde{A} = \max\{\sqrt{\Delta^2 d + \sigma^2}, \sqrt[4]{\Delta^2 d + \sigma^2}\}$, and μ^* 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 ϵ accuracy within
 298 $\tilde{\mathcal{O}}(\epsilon^{-2} \log(\epsilon^{-1}) \Delta^2)$ 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 error Δ compared to
 306 SGLD. Besides the discussion in Theorem 1, we can discuss the upper w.r.t. to Δ , due to the fact
 307 that $\log(x) \leq x^{1/e}$, one can tune the choice of ϵ and η , and achieve a $\tilde{\mathcal{O}}(\Delta^{e/(1+2e)})$ 2-Wasserstein
 308 bound for non-log-concave target distribution. Furthermore, based on Theorem 10, after carefully
 309 choosing the stepsize η , the 2-Wasserstein distance of the SGLDLP-P algorithm can be further
 310 bounded by $\tilde{\mathcal{O}}(\Delta^{e/(2+2e)})$ which is worse than the bound $\tilde{\mathcal{O}}(\Delta^{e/(1+2e)})$ 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 $\|\mathbf{x}^*\|^2 < \mathcal{D}^2$. Furthermore,
 314 let p^* denote the target distribution of \mathbf{v} and \mathbf{x} . Given any sufficiently small ϵ , if we set the step size
 315 η to be

$$\eta = \min \left\{ \frac{\epsilon \kappa_1^{-1}}{\sqrt{663552/5 \left(\frac{d}{m_1} + \mathcal{D}^2\right)}}, \frac{\epsilon^2}{2880 \kappa_1 u \left(\frac{\Delta^2 d}{4} + \sigma^2\right)} \right\},$$

316 then after K steps starting with initial points $\mathbf{x}_0 = \mathbf{v}_0 = 0$, the output $(\mathbf{x}_K, \mathbf{v}_K)$ of the SGHMCLP-L
 317 in (4) satisfies

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) = \tilde{\mathcal{O}}\left(\epsilon + \frac{\Delta}{\epsilon}\right), \quad (11)$$

318 for some K satisfying

$$K \leq \frac{\kappa_1}{\eta} \log\left(\frac{36 \left(\frac{d}{m_1} + \mathcal{D}^2\right)}{\epsilon}\right) = \tilde{\mathcal{O}}(\epsilon^{-2} \log(\epsilon^{-1}) \Delta^2).$$

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 η , the 2-Wasserstein distance of SGHMCLP-L to non-log-concave target distribution
325 can be further bounded as $\tilde{O}(\Delta^{e/(3+6e)})$, and the distance of the sample obtained by SGLDLP-L
326 can be bounded as $\tilde{O}(\Delta^{e/10(1+e)})$. 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 Q^{vc} . Instead of
329 adding real value Gaussian noise and quantizing the weights, we can design a categorical sampler
330 that samples from the space $\{\Delta, -\Delta, 0\}$ with the desired expectation μ and variance v as

$$\text{Cat}(\mu, v) = \begin{cases} \Delta, & w.p. \frac{v+\mu^2+\mu\Delta}{2\Delta^2} \\ -\Delta, & w.p. \frac{v+\mu^2-\mu\Delta}{2\Delta^2} \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

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

333 **Theorem 7.** *Let Assumptions 1, 3 and 4 hold. If $\gamma^2 \leq 4Mu$ and we set the step size to be $\eta =$
334 $\tilde{O}\left(\frac{\mu^* \epsilon^2}{\log(1/\epsilon)}\right)$, also satisfied*

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

335 let p^* denote the target distribution of (\mathbf{x}, \mathbf{v}) then after K steps starting at the initial point $\mathbf{x}_0 =$
336 $\mathbf{v}_0 = 0$ the output $(\mathbf{x}_K, \mathbf{v}_K)$ of SGHMCLP-L in 4 satisfies

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) = \tilde{O} \left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\} \log\left(\frac{1}{\epsilon}\right) + \frac{\log^{3/2}\left(\frac{1}{\epsilon}\right)}{\epsilon^2} \sqrt{\Delta}} \right), \quad (13)$$

337 for some K satisfying

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

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

$$\eta = \min \left\{ \frac{\epsilon^2}{663552/5 \left(\frac{d}{m_1} + \mathcal{D}^2 \right) \kappa_1^2}, \frac{\epsilon^2}{90u^2 \Delta^2 d \kappa_1 + 360u^2 \sigma^2 \kappa_1} \right\}$$

341 after K steps starting from the initial point $\mathbf{x}_0 = \mathbf{v}_0 = 0$ the output $(\mathbf{x}_K, \mathbf{v}_K)$ of the VC SGHMCLP-
342 L in algorithm 2 satisfies

$$\mathcal{W}_2(p(\mathbf{x}_K, \mathbf{v}_K), p^*) = \tilde{O} \left(\epsilon + \sqrt{\Delta} \right), \quad (14)$$

343 for some K satisfying

$$K \leq \frac{\kappa_1}{\eta} \log \left(\frac{36 \left(\frac{d}{m_1} + \mathcal{D}^2 \right)}{\epsilon} \right) = \tilde{O} \left(\epsilon^{-2} \log(\epsilon^{-1}) \Delta^2 \right).$$

344 Theorem 8 shows that the variance corrected quantization function can solve the overdispersion
345 problem we observe for the naive SGHMCLP-L algorithm for strongly log-concave distribution.
346 The \mathcal{W}_2 distance between the sample distribution and target distribution can be arbitrarily close
347 to $\tilde{O}(\sqrt{\Delta})$. 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 distribtuion, however for
349 non-log-concave target distribution we show VC SGHMCLP-L can achieve lower 2-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: (μ, v, Δ) $\{Q^{vc}$ returns a variable with mean μ and variance $v\}$
 $v_0 \leftarrow \Delta^2/4$ $\{\Delta^2/4$ is the largest possible variance that stochastic rounding can cause $\}$
if $v > v_0$ **then** $\{\text{add a small Gaussian noise and sample from the discrete grid to make up the remaining variance}\}$
 $x \leftarrow \mu + \sqrt{v - v_0}\xi$, where $\xi \sim \mathcal{N}(0, I_d)$
 $r \leftarrow x - Q^d(x)$
for all i **do**
 sample c_i from $\text{Cat}(|r_i|, v_0)$ as in (12)
end for
 $\theta \leftarrow Q^d(x) + \text{sign}(r) \odot c$
else $\{\text{sample from the discrete grid to achieve the target variance}\}$
 $r \leftarrow \mu - Q^s(\mu)$
for all i **do**
 $v_s \leftarrow \left(1 - \frac{|r_i|}{\Delta}\right) \cdot r_i^2 + \frac{|r_i|}{\Delta} \cdot (-r_i + \text{sign}(r_i)\Delta)^2$
 if $v > v_s$ **then**
 sample c_i from $\text{Cat}(0, v - v_s)$ as in (12)
 $\theta_i \leftarrow Q^s(\mu)_i + c_i$
 else
 $\theta_i \leftarrow Q^s(\mu)_i$
 end if
end for
end if
clip θ if outside representable range
return θ

352 **Theorem 9.** Let Assumption 1, 3 and 4 hold. If $\gamma^2 \leq 4Mu$ and we set the step size to be $\eta =$
353 $\tilde{\mathcal{O}}\left(\frac{\mu^* \epsilon^2}{\log(1/\epsilon)}\right)$, also satisfied

$$\eta \leq \min \left\{ \frac{\gamma}{4(8Mu + u\gamma + 22\gamma^2)}, \sqrt{\frac{4u^2}{4Mu + 3\gamma^2}}, \frac{6\gamma bu}{(4Mu + 3\gamma^2)d}, \frac{1}{8\gamma}, \frac{\gamma m_2}{12(21u + \gamma)M^2}, \frac{8(\gamma^2 + 2u)}{(20u + \gamma)\gamma} \right\}.$$

354 We further assume that $\mathbb{E}\|Q_G(\nabla\tilde{U}(x))\|_2^2 \leq G^2$, let p^* be the target distribution of \mathbf{x} then after
355 K steps starting at the initial point $\mathbf{x}_0 = \mathbf{v}_0 = 0$ the output (\mathbf{x}_K) of the VC SGHMCLP-L in
356 algorithm 2 satisfies

$$\mathcal{W}_2(p(\mathbf{x}_K), p^*) = \tilde{\mathcal{O}}\left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\} \log\left(\frac{1}{\epsilon}\right) + \frac{\log\left(\frac{1}{\epsilon}\right)}{\epsilon} \sqrt{\Delta}}\right), \quad (15)$$

357 for some K satisfying

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

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 $U(\mathbf{x})$ 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):

$$d\mathbf{x}_t = -\nabla U(\mathbf{x}_t) + \sqrt{2d}\mathbf{B}_t, \quad (16)$$

363 where \mathbf{B}_t represents the standard Brownian motion in \mathbb{R}^d . Under some mild conditions, it can
364 be proved that the invariant distribution of (16) converges the target distribution $\exp(-U(\mathbf{x}))$. To

Table 3: Theoretical results of the achievable 2-Wasserstein distance and the required gradient complexity for both log-concave (*italic*) non-log-concave (**bold**) target distributions, where ϵ is any sufficiently small constant, Δ is the quantization error, and μ^* and λ^* denote the concentration rate of underdamped and overdamped Langevin dynamics respectively.

	Gradient Complexity	Achieved 2-Wasserstein
Full-precision gradient accumulators		
<i>SGLD/SGHMC</i> (Theorem 4)	$\tilde{\mathcal{O}}(\log(\epsilon^{-1})\epsilon^{-2})$	$\tilde{\mathcal{O}}(\epsilon + \Delta)$
SGLD (Theorem 10)	$\tilde{\mathcal{O}}(\epsilon^{-4}\lambda^{*-1}\log^5(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \log(\epsilon^{-1})\sqrt{\Delta}\right)$
SGHMC (Theorem 5)	$\tilde{\mathcal{O}}(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \sqrt{\log(\epsilon^{-1})\Delta}\right)$
Low-precision gradient accumulators		
<i>SGLD/SGHMC</i> (Theorem 6)	$\tilde{\mathcal{O}}(\log(\epsilon^{-1})\epsilon^{-2})$	$\tilde{\mathcal{O}}(\epsilon + \epsilon^{-1}\Delta)$
<i>VC SGLD/VC SGHMC</i> (Theorem 8)	$\tilde{\mathcal{O}}(\log(\epsilon^{-1})\epsilon^{-2})$	$\tilde{\mathcal{O}}(\epsilon + \sqrt{\Delta})$
SGLD (Theorem 11)	$\tilde{\mathcal{O}}(\epsilon^{-4}\lambda^{*-1}\log^5(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \log^5(\epsilon^{-1})\epsilon^{-4}\sqrt{\Delta}\right)$
VC SGLD (Theorem 12)	$\tilde{\mathcal{O}}(\epsilon^{-4}\lambda^{*-1}\log^3(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \log^3(\epsilon^{-1})\epsilon^{-2}\sqrt{\Delta}\right)$
SGHMC (Theorem 7)	$\tilde{\mathcal{O}}(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \log^{3/2}(\epsilon^{-1})\epsilon^{-2}\sqrt{\Delta}\right)$
VC SGHMC (Theorem 9)	$\tilde{\mathcal{O}}(\epsilon^{-2}\mu^{*-2}\log^2(\epsilon^{-1}))$	$\tilde{\mathcal{O}}\left(\epsilon + \log(\epsilon^{-1})\epsilon^{-1}\sqrt{\Delta}\right)$

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla \tilde{U}(\mathbf{x}_k) + \sqrt{2\eta} \xi_{k+1}, \quad (17)$$

367 where η is the stepsize, the ξ_{k+1} is a standard Gaussian noise, and $\nabla \tilde{U}(\mathbf{x}_k)$ is an unbiased estimation
 368 of $\nabla U(\mathbf{x}_k)$. 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:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta Q_G \left(\nabla \tilde{U}(Q_W(\mathbf{x}_k)) \right) + \sqrt{2\eta} \xi_{k+1}. \quad (18)$$

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:

$$\mathbf{x}_{k+1} = Q_W \left(\mathbf{x}_k - \eta Q_G \left(\nabla \tilde{U}(\mathbf{x}_k) \right) + \sqrt{2\eta} \xi_{k+1} \right). \quad (19)$$

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 \tilde{A} have the same definition in Theorem 5,*
 381 *and λ^* be the concentration number of (16). After K steps starting with initial point $\mathbf{x}_0 = 0$, if we*
 382 *set the stepsize to be $\eta = \tilde{\mathcal{O}}\left(\left(\frac{\epsilon}{\log(1/\epsilon)}\right)^4\right)$. The output \mathbf{x}_K of SGLDLP-F in (18) satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K), p^*) \leq \tilde{\mathcal{O}}\left(\epsilon + \tilde{A} \log\left(\frac{1}{\epsilon}\right)\right), \quad (20)$$

383 *provided*

$$K = \tilde{\mathcal{O}}\left(\frac{1}{\epsilon^4 \lambda^*} \log^5\left(\frac{1}{\epsilon}\right)\right).$$

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

given: Step size η , friction γ , inverse mass u , number of training iterations K , gradient quantizer Q_G , quantization gap Δ and upper bound of low-precision representation U . Let $\text{Var}_{\mathbf{v}}^{hmc} = u(1 - e^{-2\gamma\eta})$ and $\text{Var}_{\mathbf{x}}^{hmc} = u\gamma^{-2}(2\gamma\eta + 4e^{-\gamma\eta} - e^{-2\gamma\eta} - 3)$ and $S_{\mathbf{v}} = 1$ **{Initialize the scaling parameter}**.

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

rescale $\mathbf{v}_k = \mathbf{v}_k * S_{\mathbf{v}}$ **{Restore the velocity before update}**

update $\mu(\mathbf{v}_{k+1}) = \mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta})Q_G(\nabla\tilde{U}(\mathbf{x}_k))$

update $\mu(\mathbf{x}_{k+1}) = \mathbf{x}_k + \gamma^{-1}(1 - e^{-\gamma\eta})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-\gamma\eta} - 1)Q_G(\nabla\tilde{U}(\mathbf{x}_k))$

update $S_{\mathbf{v}} = \frac{\|\mu(\mathbf{v}_{k+1})\|_{\infty}}{U}$ **{Update the Scaling}**

update $\mathbf{v}_{k+1} \leftarrow Q^{vc}(\mu(\mathbf{v}_{k+1})/S_{\mathbf{v}}, \text{Var}_{\mathbf{v}}^{hmc}/S_{\mathbf{v}}^2, \Delta)$

update $\mathbf{x}_{k+1} \leftarrow Q^{vc}(\mu(\mathbf{x}_{k+1}), \text{Var}_{\mathbf{x}}^{hmc}, \Delta)$

end for

output: samples $\{x_k\}$

384 Theorem 10 shows that the low-precision SGLD with full-precision gradient accumulators can con-
 385 verge to the non-log-concave target distribution provided a small gradient variance and quantization
 386 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 $\eta = \tilde{O}\left(\left(\frac{\epsilon}{\log(1/\epsilon)}\right)^4\right)$,*
 388 *after K steps starting at the initial point $\mathbf{x}_0 = 0$ the output \mathbf{x}_K of the SGLDLP-L in (19) satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K), p^*) = \tilde{O}\left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\}} \log\left(\frac{1}{\epsilon}\right) + \frac{\log^5\left(\frac{1}{\epsilon}\right)}{\epsilon^4} \sqrt{\Delta}\right), \quad (21)$$

389 *provided*

$$K = \tilde{O}\left(\frac{1}{\epsilon^4 \lambda^*} \log^5\left(\frac{1}{\epsilon}\right)\right).$$

390 The VC SGLDLP-L can be done as:

$$\mathbf{x}_{k+1} = Q^{vc}\left(\mathbf{x}_k - \eta Q_G(\nabla\tilde{U}(\mathbf{x}_k)), 2\eta, \Delta\right) \quad (22)$$

391 **Theorem 12.** *Let Assumption 1, 3 and 4 hold. If we set the stepsize to be $\eta = \tilde{O}\left(\frac{\epsilon^4}{\log^4\left(\frac{1}{\epsilon}\right)}\right)$, after*
 392 *K steps from the initial point $\mathbf{x}_0 = 0$ the output \mathbf{x}_K of VC SGLDLP-L in (22) satisfies*

$$\mathcal{W}_2(p(\mathbf{x}_K), p^*) = \tilde{O}\left(\epsilon + \sqrt{\max\{\sigma^2, \sigma\}} \log\left(\frac{1}{\epsilon}\right) + \frac{\log^3\left(\frac{1}{\epsilon}\right)}{\epsilon^2} \sqrt{\Delta}\right), \quad (23)$$

393 *provided*

$$K = \tilde{O}\left(\frac{1}{\epsilon^4 \lambda^*} \log^5\left(\frac{1}{\epsilon}\right)\right).$$

394 C Technical Detail

395 In this section, we disclose more details of empirical experiments. When implementing low-
 396 precision SGHMC on classification task in the CIFAR-10 and CIFAR-100 dataset, we observed that
 397 the momentum term \mathbf{v} tend to gather in a small range around zero in which case the low-precision
 398 representations of \mathbf{v} end up in gathering only few points, thus the momentum information is seri-
 399 ously lost and cause in performance degradation. In order to tackle this problem and fully utilize all
 400 the low-precision representations, we borrow the idea of rescaling from the bit centering trick and
 401 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
 403 corrected quantization function. Bit centering trick is a technique to increase the accuracy low-
 404 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 v 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:

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

415 In this section, we prove the convergence of SGHMCLP-F in terms of 2-Wasserstein distance for
 416 strongly-log-concave target distribution via coupling argument. To simplify the notation we define
 417 the quantized stochastic gradients at \mathbf{x} as:

$$\begin{aligned} \tilde{g}(\mathbf{x}) &:= Q_G(\nabla\tilde{U}(Q_W(\mathbf{x}))) & (24) \\ &=: \nabla U(\mathbf{x}) + \xi. & (25) \end{aligned}$$

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

$$\begin{aligned} \|\mathbb{E}\xi\|^2 &\leq M^2 \frac{\Delta^2 d}{4} \\ \mathbb{E}[\|\xi\|^2] &\leq (M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2. \end{aligned}$$

420

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

427 In order to calculate the Wasserstein distance from the proposed sample $(\mathbf{x}_K, \mathbf{v}_K)$ and the target
 428 distribution sample $(\mathbf{x}^*, \mathbf{v}^*)$, we define sample $q_k = (\mathbf{x}_k, \mathbf{x}_k + \mathbf{v}_k)$ and the target distribution
 429 sample $q^* = (\mathbf{x}^*, \mathbf{x}^* + \mathbf{v}^*)$. Let $p_k = (\mathbf{x}_k, \mathbf{v}_k)$ and $\widehat{\Phi}_\eta$ be the operator that maps from p_k to p_{k+1}
 430 i.e.

$$p_{k+1} = \widehat{\Phi}_\eta p_k.$$

431 The solution $(\mathbf{x}_t, \mathbf{v}_t)$ of the continuous underdamped Langevin dynamics with exact gradient satis-
 432 fies the following equations:

$$\begin{aligned} \mathbf{v}_t &= \mathbf{v}_0 e^{-\gamma t} - u \left(\int_0^t e^{-\gamma(t-s)} \nabla U(\mathbf{x}_s) ds \right) + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s, & (26) \\ \mathbf{x}_t &= \mathbf{x}_0 + \int_0^t \tilde{\mathbf{v}}_s ds. \end{aligned}$$

433 Let Φ_η denote the operator that maps p_0 to the solution of continuous underdamped Langevin dy-
 434 namics in (26) after time step η . Notice the solution $(\tilde{\mathbf{v}}_t, \tilde{\mathbf{x}}_t)$ of the discrete underdamped Langevin
 435 dynamics with an exact gradient can be written as

$$\begin{aligned} \tilde{\mathbf{v}}_t &= \tilde{\mathbf{v}}_0 e^{-\gamma t} - u \left(\int_0^t e^{-\gamma(t-s)} \nabla U(\tilde{\mathbf{x}}_0) ds \right) + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s, & (27) \\ \tilde{\mathbf{x}}_t &= \tilde{\mathbf{x}}_0 + \int_0^t \tilde{\mathbf{v}}_s ds. \end{aligned}$$

436 We can also define a similar operator for the discrete underdamped Langevin dynamics solution
 437 $\tilde{p}_t = (\tilde{\mathbf{x}}_t, \tilde{\mathbf{v}}_t)$, let $\tilde{\Phi}_t$ be the operator that maps \tilde{p}_0 to \tilde{p}_t . Furthermore the SGHMCLP-F can be
 438 written as:

$$\begin{aligned} \mathbf{v}_t &= \mathbf{v}_0 e^{-\gamma t} - u \left(\int_0^t e^{-\gamma(t-s)} \tilde{g}(\mathbf{x}_0) ds \right) + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s, \\ \mathbf{x}_t &= \tilde{\mathbf{x}}_0 + \int_0^t \mathbf{v}_s ds. \end{aligned} \quad (28)$$

439 Given $\tilde{g}(\mathbf{x}_0) = \nabla U(\mathbf{x}_0) + \xi_0$ and $\mathbf{x}_0 = \tilde{\mathbf{x}}_0$, we know:

$$\begin{aligned} \mathbf{v}_t &= \tilde{\mathbf{v}}_t - u \left(\int_0^t e^{-\gamma(t-s)} ds \right) \xi \\ \mathbf{x}_t &= \tilde{\mathbf{x}}_t - u \left(\int_0^t \left(\int_0^r e^{-\gamma(t-s)} ds \right) dr \right) \xi. \end{aligned} \quad (29)$$

440 **Lemma 14.** Let q_0 be some initial distribution and $\tilde{\Phi}_\eta$ and Φ_η 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 stepsize $1 > \eta > 0$, then the Wasserstein distance satisfies

$$\mathcal{W}_2^2(\Phi_\eta q_0, q^*) \leq \left(\mathcal{W}_2(\tilde{\Phi}_\eta q_0, q^*) + \sqrt{5}/2u\eta\sqrt{d}M\Delta \right)^2 + 5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right).$$

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 $\mathcal{O}(\eta^2 \Delta^2)$. Furthermore from the corollary 7 in Cheng et al. [2018] we know
 446 that for any $i \in \{1, \dots, K\}$:

$$\mathcal{W}_2^2(\Phi_\eta q_i, q^*) \leq e^{-\eta/2\kappa_1} \mathcal{W}_2^2(q_i, q^*), \quad (30)$$

447 where $\kappa_1 = M/m_1$ is the condition number. Let \mathcal{E}_K denote the $26(d/m_1 + \mathcal{D}^2)$, and from the
 448 discretization error bound from Theorem 9 and Lemma 8 (sandwich inequality) in Cheng et al.
 449 [2018], we get

$$\mathcal{W}_2(\Phi_\eta q_i, \tilde{\Phi}_\eta q_i) \leq 2\mathcal{W}_2(\Phi_\eta p_i, \tilde{\Phi}_\eta p_i) \leq \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}.$$

450 By triangle inequality:

$$\begin{aligned} \mathcal{W}_2(\tilde{\Phi}_\eta q_i, q^*) &\leq \mathcal{W}_2(\Phi_\eta q_i, \tilde{\Phi}_\eta q_i) + \mathcal{W}_2(\Phi_\eta q_i, q^*) \\ &\leq \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + e^{-\eta/2\kappa_1} \mathcal{W}_2(q_i, q^*). \end{aligned}$$

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

$$\mathcal{W}_2^2(\tilde{\Phi}_\eta q_i, q^*) \leq \left(e^{-\eta/2\kappa_1} \mathcal{W}_2(q_i, q^*) + \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{5}/2u\eta\sqrt{d}M\Delta \right)^2 + 5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right).$$

452 By invoking the Lemma 7 in Dalalyan and Karagulyan [2019] we can bound the 2-Wasserstein
 453 distance by:

$$\begin{aligned} \mathcal{W}_2(q_K, q^*) &\leq e^{-K\eta/2\kappa_1} \mathcal{W}_2(q_0, q^*) + \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M\Delta\sqrt{5d}}{2}}{1 - e^{-\eta/2\kappa_1}} \\ &\quad + \frac{5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M\Delta\sqrt{5d}}{2} + \sqrt{1 - e^{-\eta/\kappa_1}} \sqrt{5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}}. \end{aligned}$$

454 Finally by sandwich inequality we have:

$$\begin{aligned} \mathcal{W}_2(p_K, p^*) &\leq 4e^{-K\eta/2\kappa} \mathcal{W}_2(p_0, p^*) + 4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M\Delta\sqrt{5d}}{2}}{1 - e^{-\eta/2\kappa}} \\ &\quad + \frac{20u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M\Delta\sqrt{5d}}{2} + \sqrt{1 - e^{-\eta/\kappa}} \sqrt{5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}}. \end{aligned}$$

455 Now we let the first term less than $\epsilon/3$, from the lemma 13 in [Cheng et al., 2018] we know that
 456 $\mathcal{W}_2(p_K, p^*) \leq 3 \left(\frac{d}{m_1} + \mathcal{D}^2 \right)$. So we can choose K as the following,

$$K \leq \frac{2\kappa_1}{\eta} \log \left(36 \left(\frac{d}{m_1} + \mathcal{D}^2 \right) \right).$$

457 Next, we choose a stepsize $\eta \leq \frac{\epsilon\kappa_1^{-1}}{\sqrt{479232/5(d/m_1 + \mathcal{D}^2)}}$ to ensure the second term is controlled below
 458 $\epsilon/3 + \frac{16\kappa_1 u M \Delta \sqrt{5d}}{2}$. Since $1 - e^{-\eta/2\kappa_1} \geq \eta/4\kappa_1$ and definition of \mathcal{E}_K ,

$$\begin{aligned} 4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M \Delta \sqrt{5d}}{2}}{1 - e^{-\eta/2\kappa}} &\leq 4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M \Delta \sqrt{5d}}{2}}{\eta/4\kappa_1} \leq 16\kappa_1 \left(\eta \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{uM \Delta \sqrt{5d}}{2} \right) \\ &\leq \epsilon/3 + \frac{16\kappa_1 u M \Delta \sqrt{5d}}{2}. \end{aligned}$$

459 Finally by choosing the stepsize satisfied that,

$$\eta \leq \frac{\epsilon M \Delta \sqrt{5d}}{120u \left[(M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right]},$$

460 the third term can be bounded as:

$$\begin{aligned} &\frac{20u^2 \eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \frac{u\eta M \Delta \sqrt{5d}}{2} + \sqrt{1 - e^{-\eta/\kappa}} \sqrt{5u^2 \eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}} \\ &\leq \frac{20u^2 \eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}{\frac{u\eta M \Delta \sqrt{5d}}{2}} = 40u\eta \frac{\left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}{M \Delta \sqrt{5d}} \leq \epsilon/3. \end{aligned}$$

461 This complete the proof.

462 D.2 Proof of Theorem 5

463 In this section we analyze the Wasserstein distance between the sample (\mathbf{x}_k, v_k) 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
 467 ideal continuous underdamped Langevin dynamics, also the distance between solutions of the ideal
 468 continuous underdamped Langevin dynamics and the target distribution.

469 Again let $p_k = (\mathbf{x}_k, v_k)$ denote the low-precision sample from (3) at k -th iteration, let $\hat{p}_t = (\hat{x}_t, \hat{v}_t)$
 470 denote the sample from the ideal continuous underdamped Langevin dynamics in 26 at time t . Then
 471 the Wasserstein distance between the p_k and the target distribution p^* can be bounded as:

$$\mathcal{W}_2(p_K, p^*) \leq \mathcal{W}_2(p_K, \hat{p}_{K\eta}) + \mathcal{W}_2(\hat{p}_{K\eta}, p^*).$$

472 We first bound $\mathcal{W}_2(p_K, \hat{p}_{K\eta})$ by invoking the weighted CKP inequality Bolley and Villani [2005],

$$\mathcal{W}_2^2(p_K, \hat{p}_{K\eta}) \leq \Lambda \left(\sqrt{D_{KL}(p_K || \hat{p}_{K\eta})} + \sqrt[4]{D_{KL}(p_K || \hat{p}_{K\eta})} \right),$$

473 where $\Lambda = 2 \inf_{\theta > 0} \sqrt{1/\theta (3/2 + \log \mathbb{E}_{\hat{p}_{K\eta}} [exp(\theta(\|\hat{x}_{K\eta}\|^2 + \|\hat{v}_{K\eta}\|^2))])}$. We define a Lyapunov
 474 function for every $(x, v) \in \mathbb{R}^d \times \mathbb{R}^d$

$$\mathcal{E}(\mathbf{x}, \mathbf{v}) = \|\mathbf{x}\|^2 + \|\mathbf{x} + 2\mathbf{v}/\gamma\|^2 + 8u(U(\mathbf{x}) - U(\mathbf{x}^*))/\gamma^2.$$

475 Note that $\|a\|^2 + \|b\|^2 \geq \|a - b\|^2 / 2$ and $U(x) \geq U(x^*)$, we can have:

$$\mathcal{E}(x, v) \geq \|x\|^2 + \|x + 2v/\gamma\|^2 \geq \max\{\|x\|^2, 2\|v/\gamma\|^2\}.$$

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

$$\begin{aligned} \Lambda &\leq 2 \inf_{0 < \theta \leq \min\{\frac{\gamma}{128u}, \frac{m_2}{32}\}} \sqrt{\frac{1}{\theta} \left(\frac{3}{2} + 2\theta\mathcal{E}(\mathbf{X}_0, \mathbf{V}_0) + \frac{32M\theta u(4d + 2b + m_2\|\mathbf{x}^*\|^2)}{\gamma^2 m_2} \right)} \\ &\leq 2 \sqrt{2\mathcal{E}(\mathbf{X}_0, \mathbf{V}_0) + \frac{32M\theta u(4d + 2b + m_2\|\mathbf{x}^*\|^2) + 16(12um_2 + 3\gamma^2)}{\gamma^2 m_2}} := \bar{\Lambda}. \end{aligned}$$

477 It remains to bound the divergence between the distribution p_K and $\hat{p}_{K\eta}$. We first define a continuous
478 interpolation of the low-precision sample $(\mathbf{x}_k, \mathbf{v}_k)$,

$$d\mathbf{v}_t = -\gamma\mathbf{v}_t dt - uG_t dt + \sqrt{2\gamma u} dB_t \quad (31)$$

$$d\mathbf{x}_t = \mathbf{v}_t dt, \quad (32)$$

479 where $G_t = \sum_{k=0}^K \tilde{g}(\mathbf{x}_k) \mathbf{1}_{t \in [k\eta, (k+1)\eta)}$. Integrating this equation from time 0 to t , we can get

$$\begin{aligned} \mathbf{v}_t &= \mathbf{v}_0 - \int_0^t \gamma\mathbf{v}_s ds - \int_0^t uG_s dt + \int_0^t \sqrt{2\gamma u} dB_s \\ \mathbf{x}_t &= \mathbf{x}_0 + \int_0^t \mathbf{v}_s ds. \end{aligned}$$

480 Notice that when $t = k\eta$, the solution of (31) has the same distribution with the low-precision
481 sample $(\mathbf{x}_k, \mathbf{v}_k)$. Now by Girsanov formula we can compute the Radon-Nikodym derivative of $\hat{p}_{K\eta}$
482 with respect to p_K as follow:

$$\frac{d\hat{p}_{K\eta}}{dp_K} = \exp \left\{ \sqrt{\frac{\gamma u}{2}} \int_0^t (\nabla U(\mathbf{x}_s) - G_s) dB_s - \frac{\gamma u}{4} \int_0^t \|\nabla U(\mathbf{x}_s) - G_s\|^2 ds \right\}.$$

483 It follows that

$$\begin{aligned} D_{KL}(p_K \|\hat{p}_{K\eta}) &= \mathbb{E}_{p_K} \left[\log \left(\frac{d\hat{p}_{K\eta}}{dp_K} \right) \right] \quad (33) \\ &= \frac{\gamma u}{4} \mathbb{E} \int_0^{K\eta} \|\nabla U(\mathbf{x}_s) - G_s\|^2 ds \\ &= \frac{\gamma u}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - G_s\|^2 \right] ds \\ &= \frac{\gamma u}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \tilde{g}(\mathbf{x}_k)\|^2 \right] ds. \end{aligned}$$

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

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \tilde{g}(\mathbf{x}_k)\|^2 \right] \leq 2\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] + 2\mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right]. \quad (34)$$

485 We now bound the first term in the RHS of the (34). By the smooth Assumption1, we have

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] \leq M^2 \mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right].$$

486 Notice that

$$\begin{aligned} \mathbf{x}_s &= \mathbf{x}_k + \int_{k\eta}^s \mathbf{v}_r dr \\ &= \mathbf{x}_k + \int_{k\eta}^s \left(\mathbf{v}_{k\eta} e^{-\gamma(r-k\eta)} - u \left(\int_{k\eta}^r e^{-\gamma(r-z)} \tilde{g}(\mathbf{x}_k) dz \right) + \sqrt{2\gamma u} \int_{k\eta}^r e^{-\gamma(r-z)} dB_z \right) dr. \end{aligned}$$

487 This further implies that:

$$\begin{aligned}
\|\mathbf{x}_s - \mathbf{x}_k\|^2 &= \left\| \int_{k\eta}^s \left(\mathbf{v}_{k\eta} e^{-\gamma(r-k\eta)} - u \left(\int_{k\eta}^r e^{-\gamma(r-z)} \tilde{g}(\mathbf{x}_k) dz \right) + \sqrt{2\gamma u} \int_{k\eta}^r e^{-\gamma(r-z)} dB_z \right) dr \right\|^2 \\
&\leq 3 \left\| \int_{k\eta}^s \mathbf{v}_{k\eta} e^{\gamma(k\eta-r)} dr \right\|^2 + 3 \left\| \int_{k\eta}^s \int_{k\eta}^r u \tilde{g}(\mathbf{x}_k) e^{\gamma(z-r)} dz dr \right\|^2 + 6ru \left\| \int_{k\eta}^s \int_0^s e^{-\gamma(r-z)} dB_z dr \right\|^2 \\
&\leq 3\eta^2 \|\mathbf{v}_k\|^2 + 3u^2\eta^4 \|\tilde{g}(\mathbf{x}_k)\|^2 + 3 \left[\frac{u}{\gamma^2} \left(2\gamma(s-k\eta) + 4e^{-\gamma(s-k\eta)} - e^{-2\gamma(s-k\eta)} - 3 \right) d \right] \\
&\leq 3\eta^2 \left(\|\mathbf{v}_k\|^2 + u^2\eta^2 \|\tilde{g}(\mathbf{x}_k)\|^2 + 2du \right), \tag{35}
\end{aligned}$$

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

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] \leq 3M^2\eta^2 \left(\mathbb{E} \|v_k\|^2 + u^2\eta^2 \mathbb{E} \|\tilde{g}(\mathbf{x}_k)\|^2 + 2du \right).$$

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

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right] \leq (M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2.$$

491 We need to introduce a lemma to bound the $\sup_k \|\mathbf{x}_k\|^2$, $\sup_k \|v_k\|^2$ and $\sup_k \|\tilde{g}(\mathbf{x}_k)\|^2$.

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

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

493 then for all $k \geq 0$ the $\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right]$, $\mathbb{E} \left[\|v_k\|^2 \right]$ and $\mathbb{E} \left[\|\tilde{g}(\mathbf{x}_k)\|^2 \right]$ can be bounded as

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] &\leq \bar{\mathcal{E}} + C_0 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
\mathbb{E} \left[\|v_k\|^2 \right] &\leq \gamma^2 \bar{\mathcal{E}} / 2 + \gamma^2 C_0 / 2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
\mathbb{E} \left[\|\tilde{g}(\mathbf{x}_k)\|^2 \right] &\leq 2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) + 4M^2 \bar{\mathcal{E}} + 4G^2
\end{aligned}$$

494 where $\bar{\mathcal{E}}$ and C_0 are defined as:

$$\begin{aligned}
\bar{\mathcal{E}} &= \mathbb{E} [\mathcal{E}(\mathbf{x}_0, \mathbf{v}_0)] + \frac{24(21u + \gamma)uM}{m_2\gamma^3} G^2 + \frac{96(d+b)uM}{m_2\gamma^2}, \quad G = \|\nabla U(0)\| \\
C_0 &= \frac{96u(\gamma^2 + 2u)}{m_2\gamma^4}.
\end{aligned}$$

495 We now ready to bound $\mathbb{E} \left[\|\nabla U(\mathbf{x}_s - \tilde{g}(\mathbf{x}_k))\|^2 \right]$ as:

$$\begin{aligned}
\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \tilde{g}(\mathbf{x}_k)\|^2 \right] &\leq 2\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] + 2\mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right] \\
&\leq 6M^2\eta^2 \left(\mathbb{E} \|v_k\|^2 + u^2\eta^2\mathbb{E} \|\tilde{g}(\mathbf{x}_k)\|^2 + 2du \right) + 2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
&\leq 6M^2\eta^2 \left((\gamma^2/2 + 4M^2u^2\eta^2)\bar{\mathcal{E}} + (\gamma^2C_0/2 + 2u^2\eta^2) \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) + 4u^2\eta^2G^2 + 2du \right) \\
&\quad + 2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\
&\leq 6M^2\eta^2 \left[(\gamma^2/2 + 4M^2u^2\eta^2)\bar{\mathcal{E}} + 4u^2\eta^2G^2 + 2du \right] \\
&\quad + (6M^2\eta^2(\gamma^2C_0/2 + 2u^2\eta^2) + 2) \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right).
\end{aligned}$$

496 Thus the divergence can be bounded as:

$$\begin{aligned}
D_{KL}(p_K || \hat{p}_{K\eta}) &\leq \frac{3\gamma u}{2} M^2 K \eta^3 \left[(\gamma^2/2 + 4M^2u^2\eta^2)\bar{\mathcal{E}} + 4u^2\eta^2G^2 + 2du \right] \\
&\quad + \frac{\gamma u}{4} K \eta \left(6M^2\eta^2(\gamma^2C_0/2 + 2u^2\eta^2) + 2 \right) \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right).
\end{aligned}$$

497 By the weighted CKP inequality and given $K\eta \geq 1$,

$$\begin{aligned}
\mathcal{W}_2(p_K, \hat{p}_{K\eta}) &\leq \bar{\Lambda} \left(\sqrt{D_{KL}(p_K || \hat{p}_{K\eta})} + \sqrt[4]{D_{KL}(p_K || \hat{p}_{K\eta})} \right) \\
&\leq \bar{\Lambda} \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} \right) \sqrt{K\eta},
\end{aligned}$$

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

$$\begin{aligned}
\tilde{C}_0 &= \sqrt{\frac{3\gamma u}{2} M^2 \left[(\gamma^2/2 + 4M^2u^2\eta^2)\bar{\mathcal{E}} + 4u^2\eta^2G^2 + 2du \right]} + \sqrt[4]{\frac{3\gamma u}{2} M^2 \left[(\gamma^2/2 + 4M^2u^2\eta^2)\bar{\mathcal{E}} + 4u^2\eta^2G^2 + 2du \right]} \\
\tilde{C}_1 &= \sqrt{\frac{\gamma u}{4} \left(6M^2\eta^2(\gamma^2C_0/2 + 2u^2\eta^2) + 2 \right)} + \sqrt[4]{\frac{\gamma u}{4} \left(6M^2\eta^2(\gamma^2C_0/2 + 2u^2\eta^2) + 2 \right)} \\
\tilde{A} &= \max \left\{ \sqrt{\left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)}, \sqrt[4]{\left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right)} \right\}.
\end{aligned}$$

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

$$\mathcal{W}_2(\hat{p}_{K\eta}, p^*) \leq \Gamma_0 e^{-\mu^* K\eta},$$

500 where $\mu^* = e^{-\tilde{C}(d)}$ denotes the concentration rate of the underdamped Langevin dynamics and Γ_0
501 is a constant of order $\mathcal{O}(1/\mu^*)$. Combining this inequality with previous analysis we can prove:

$$\mathcal{W}_2(p_K, p^*) \leq \bar{\Lambda} \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} \right) \sqrt{K\eta} + \Gamma_0 e^{-\mu^* K\eta}. \quad (36)$$

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

$$\bar{\Lambda} \tilde{C}_0 \sqrt{K\eta^2} = \frac{\epsilon}{2} \quad \text{and} \quad \Gamma_0 e^{-\mu^* K\eta} = \frac{\epsilon}{2}. \quad (37)$$

503 Solving the equation (37), we can have

$$K\eta = \frac{\log\left(\frac{2\Gamma_0}{\epsilon}\right)}{\mu^*} \quad \text{and} \quad \eta = \frac{\epsilon^2}{4\bar{\Lambda}^2 \tilde{C}_0^2 K\eta}.$$

504 Combining these two we can have

$$\eta = \frac{\epsilon^2 \mu^*}{4\bar{\Lambda}^2 \tilde{C}_0^2 \log\left(\frac{2\Gamma_0}{\epsilon}\right)} \quad \text{and} \quad K = \frac{4\bar{\Lambda}^2 \tilde{C}_0^2 \log^2\left(\frac{2\Gamma_0}{\epsilon}\right)}{\epsilon^2 (\mu^*)^2}.$$

505 Plugging in (36) completes the proof.

506 **D.3 Proof of Thoerem 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 \mathbf{x}_0 has a bounded and strictly positive*
 512 *density and satisfies the following:*

$$\kappa_0 := \log \int_{\mathbb{R}^d} e^{\|\mathbf{x}\|^2} p_0(\mathbf{x}) d\mathbf{x} < \infty.$$

513 Note that the for initial distribution $\mathbf{x}_0 = 0$, the value $\kappa_0 = 0$ is bounded and the assumption is
 514 satisfied. Recall the Overdamped Langevin dynamics is

$$d\mathbf{x}_t = -\nabla U(\mathbf{x}_t) dt + \sqrt{2} dB_t. \quad (38)$$

515 We further define the value of the energy function and the gradient at point 0 at 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 satisfies

$$\mathbf{x}_t = \mathbf{x}_0 - \int_0^t \nabla U(\mathbf{x}_s) ds + \sqrt{2} \int_0^t dB_s. \quad (39)$$

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

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

521 where $G_s = \sum_{k=0}^K \tilde{g}(\hat{x}_k) \mathbf{1}_{s \in [k\eta, (k+1)\eta)}$. The Wasserstein distance can bounded as

$$\mathcal{W}_2(p_K, p^*) \leq \mathcal{W}_2(p_K, \hat{p}_{K\eta}) + \mathcal{W}_2(\hat{p}_{K\eta}, p^*),$$

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

$$\mathcal{W}_2(p_K, \hat{p}_{K\eta}) \leq C_{\hat{p}_{K\eta}} \left[\sqrt{D_{KL}(p_K || \hat{p}_{K\eta})} + \left(\frac{D_{KL}(p_K || \hat{p}_{K\eta})}{2} \right)^{1/4} \right],$$

523 where the constant $C_{\hat{p}_{K\eta}} = 2 \inf_{\lambda > 0} \left(\frac{1}{\lambda} \left(\frac{3}{2} + \log \int_{\mathbb{R}^d} e^{\lambda \|\omega\|^2} \hat{P}_{K\eta}(d\omega) \right) \right)$. By Lemma 4 in Raginsky
 524 et al. [2017] and assuming $K\eta > 1$, we can wrtie:

$$\mathcal{W}_2^2(p_K, \hat{p}_{K\eta}) \leq (12 + 8(\kappa_0 + 2b + 2d)K\eta) \left(D_{KL}(p_K || \hat{p}_{K\eta}) + \sqrt{D_{KL}(p_K || \hat{p}_{K\eta})} \right).$$

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

$$\frac{d\hat{p}_{K\eta}}{dp_K} = exp \left\{ \frac{1}{2} \int_0^t (\nabla U(\mathbf{x}_s) - G_s) d\mathbf{B}_s - \frac{1}{4} \int_0^t \|\nabla U(\mathbf{x}_s) - G_s\|^2 ds \right\}.$$

527 Thus, we have:

$$\begin{aligned}
D_{KL}(p_K || \hat{p}_{K\eta}) &= \mathbb{E}_{p_K} \left[\log \left(\frac{d\hat{p}_{K\eta}}{dp_K} \right) \right] \\
&= \frac{1}{4} \int_0^{K\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - G_s\|^2 \right] ds \\
&= \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \tilde{g}(\mathbf{x}_k)\|^2 \right] ds \\
&\leq \frac{1}{2} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] \\
&\quad + \frac{1}{2} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right] \\
&\leq \frac{M^2}{2} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right] \\
&\quad + \frac{1}{2} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right]. \tag{41}
\end{aligned}$$

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

$$\begin{aligned}
\mathbf{x}_s - \mathbf{x}_k &= -(s - k\eta)\tilde{g}(\mathbf{x}_k) + \sqrt{2} (B_s - B_{k\eta}) \\
&= -(s - k\eta)\nabla U(\mathbf{x}_k) + (s - k\eta) (\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)) + \sqrt{2} (B_s - B_{k\eta}),
\end{aligned}$$

529 thus,

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right] &\leq 3\eta^2 \mathbb{E} \left[\|\nabla U(\mathbf{x}_k)\|^2 \right] + 3\eta^2 \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2 \right] + 6\eta d \\
&\leq 3\eta^2 (M\mathbb{E} \left[\|\mathbf{x}_k\| \right] + G)^2 + 3\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) + 6\eta d. \tag{42}
\end{aligned}$$

530 Similarly, we need a uniform bound of $\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right]$.

531 **Lemma 16.** Under assumptions 1, 3 and 4, if we set the stepsize $\eta \in (0, 1 \wedge \frac{m_2}{2M^2})$, then for all

532 $k \geq 0$, the $\mathbb{E} \left[\|\mathbf{v}_{\mathbf{x}_k}\|^2 \right]$ can be bounded as

$$\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] \leq \mathcal{E} + \frac{2(M^2 + 1)\Delta^2 d}{4m_2},$$

533 provided $\mathcal{E} = \mathbb{E} \left[\|\mathbf{x}_0\|^2 \right] + \frac{M}{m_2} (2b + 2\eta G^2 + 2d)$.

534 Using this bound, we can further bound $\mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_s\|^2 \right]$ as:

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_s\|^2 \right] &\leq 6\eta^2 M^2 \left(\mathcal{E} + \frac{2(M^2 + 1)\Delta^2 d}{m_2} \frac{\Delta^2 d}{4} \right) + 6\eta^2 G^2 + 3\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) + 6\eta d \\
&\leq 6\eta^2 M^2 \mathcal{E} + 6\eta^2 G^2 + 6\eta d + \left(\frac{12\eta^2 M^2 (M^2 + 1)}{m_2} + 3(M^2 + 1) \right) \eta^2 \frac{\Delta^2 d}{4} + 3\eta^2 \sigma^2, \\
&=: \bar{\mathcal{E}}\eta + C\eta^2 \frac{\Delta^2 d}{4} + 3\eta^2 \sigma^2
\end{aligned}$$

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

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

536 Thus the divergence can be bounded as:

$$\begin{aligned}
D_{KL}(p_K || \hat{p}_{K\eta}) &\leq \frac{M^2}{2} \left(\bar{\mathcal{E}} + C\eta \frac{\Delta^2 d}{4} + 3\eta\sigma^2 \right) K\eta^2 + \frac{1}{2} \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) K\eta \\
&= \frac{M^2}{2} \bar{\mathcal{E}} K\eta^2 + \left(\frac{M^2}{2} C\eta^2 + \frac{1}{2} (M^2 + 1) \right) \frac{\Delta^2 d}{4} K\eta + \frac{3M^2\eta^2 + 1}{2} \sigma^2 K\eta \\
&= \frac{M^2}{2} \bar{\mathcal{E}} K\eta^2 + \left(\frac{M^2}{2} C + \frac{1}{2} (M^2 + 1) \right) \frac{\Delta^2 d}{4} K\eta + \frac{3M^2 + 1}{2} \sigma^2 K\eta \\
&=: C_0 K\eta^2 + C_1 \frac{\Delta^2 d}{4} K\eta + C_2 \sigma^2 K\eta.
\end{aligned}$$

537 We are ready to bound the Wasserstein distance,

$$\begin{aligned}
\mathcal{W}_2^2(p_K, \hat{p}_{K\eta}) &\leq (12 + 8(\kappa_0 + 2b + 2d)) \left((C_0 + \sqrt{C_0})\sqrt{\eta} + (C_1 + \sqrt{C_1})A + (C_2 + \sqrt{C_2})B \right) (K\eta)^2 \\
&=: \left(\widetilde{C}_0^2 \sqrt{\eta} + \widetilde{C}_1^2 A + \widetilde{C}_2^2 B \right) (K\eta)^2,
\end{aligned}$$

538 where the constants are defined as:

$$\begin{aligned}
A &= \max \left\{ \frac{\Delta^2 d}{4}, \sqrt{\frac{\Delta^2 d}{4}} \right\} \\
B &= \max \left\{ \sigma^2, \sqrt{\sigma^2} \right\} \\
\widetilde{C}_0^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_0 + \sqrt{C_0}) \\
\widetilde{C}_1^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_1 + \sqrt{C_1}) \\
\widetilde{C}_2^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_2 + \sqrt{C_2}).
\end{aligned}$$

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

$$\begin{aligned}
\mathcal{W}_2(\hat{p}_{K\eta}, p^*) &\leq \sqrt{2C_{LS} \left(\log \|p_0\|_\infty + \frac{d}{2} \log \frac{3\pi}{m\beta} + \beta \left(\frac{M\kappa_0}{3} + B\sqrt{\kappa_0} + G_0 + \frac{b}{2} \log 3 \right) \right)} e^{-K\eta/\beta C_{LS}} \\
&=: \widetilde{C}_3 e^{-K\eta/\beta C_{LS}}
\end{aligned}$$

540 Finally, we can have

$$\mathcal{W}_2(p_K, p^*) \leq \left(\widetilde{C}_0 \eta^{1/4} + \widetilde{C}_1 \sqrt{A} + \widetilde{C}_2 \sqrt{B} \right) K\eta + \widetilde{C}_3 e^{-K\eta/\beta C_{LS}}. \quad (43)$$

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

$$\widetilde{C}_0 K\eta^{5/4} = \frac{\epsilon}{2} \quad \text{and} \quad \widetilde{C}_3 e^{-K\eta/\beta C_{LS}} = \frac{\epsilon}{2}. \quad (44)$$

542 Solving the (37), we can have

$$K\eta = C_{LS} \log \left(\frac{2\widetilde{C}_3}{\epsilon} \right) \quad \text{and} \quad \eta = \frac{\epsilon^4}{16\widetilde{C}_0^4 (K\eta)^4}.$$

543 Combining these two we can have

$$\eta = \frac{\epsilon^4}{16\widetilde{C}_0^4 C_{LS}^4 \log^4 \left(\frac{2\widetilde{C}_3}{\epsilon} \right)} \quad \text{and} \quad K = \frac{16\widetilde{C}_0^4 C_{LS}^5 \log^5 \left(\frac{2\widetilde{C}_3}{\epsilon} \right)}{\epsilon^4}.$$

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

545 **D.4 Proof of Theorem 6**

546 Recall the SGHMCLP-L the update rule:

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

547 If we let $\alpha_k^{\mathbf{x}}$ and $\alpha_k^{\mathbf{v}}$ denote the quantization error,

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

548 we can rewrite the update rule as:

$$\begin{aligned}\mathbf{v}_{k+1} &= \mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{\gamma\eta})Q_G(\nabla\tilde{U}(\mathbf{x}_s)) + \xi_k^{\mathbf{v}} + \alpha_k^{\mathbf{v}} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \gamma^{-1}(1 - e^{-\gamma\eta})\mathbf{v}_k + u\gamma^{-2}(\gamma\eta + e^{-\gamma\eta} - 1)Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{x}} + \alpha_k^{\mathbf{x}}.\end{aligned}\quad (45)$$

549 Similarly, we can define a continuous interpolation of (45) for $t \in (0, \eta]$.

$$\begin{aligned}\mathbf{v}_t &= \mathbf{v}_0 e^{-\gamma t} - u \left(\int_0^t e^{-\gamma(t-s)} (\nabla U(\mathbf{x}_0) + \zeta) ds \right) + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s + \int_0^t \alpha_v(s) ds \\ \mathbf{x}_t &= \mathbf{x}_0 + \int_0^t \mathbf{v}_s ds + \int_0^t \alpha_x(s) ds,\end{aligned}\quad (46)$$

550 where the $\zeta = Q_G(\nabla\tilde{U}(\hat{x}_0)) - \nabla\tilde{U}(\hat{x}_0)$ the function $\alpha_v(s)$, $\alpha_x(s)$ are defined as:

$$\begin{aligned}\alpha_v(s) &= \sum_{k=0}^{\infty} \alpha_k^{\mathbf{v}}/\eta \mathbf{1}_{s \in (k\eta, (k+1)\eta)} \\ \alpha_x(s) &= \sum_{k=0}^{\infty} \alpha_k^{\mathbf{x}}/\eta \mathbf{1}_{s \in (k\eta, (k+1)\eta)}.\end{aligned}$$

551 If we let $\hat{p}_0 = (\hat{x}_0, \hat{v}_0)$ be the initial sample and $\hat{p}_t = (\hat{x}_t, \hat{v}_t)$ be the sample that satisfies the
552 previous equations, we can define an operator $\hat{\Phi}_t$ that maps \hat{p}_0 to \hat{p}_t i.e., $\hat{p}_t = \hat{\Phi}_t \hat{p}_0$. Notice that
553 since \hat{p}_t is the continuous interpolation of (4), thus $\hat{p}_{k\eta} = p_k = (\mathbf{x}_k, v_k)$. Similarly, we define
554 $q_k = (\mathbf{x}_k, v_k + \mathbf{x}_k) =: (\mathbf{x}_k, \omega_k)$ as a tool to analyze the convergence of p_k .

555 We are now ready to compute the Wasserstein distance between $\hat{\Phi}_\eta q_0$ and q^* . Let Γ_1 be all of the
556 couplings between $\tilde{\Phi}_\eta q_0$ and q^* , and Γ_2 be all of the couplings between $\hat{\Phi}_\eta q_0$ and q^* . Let r_1 be the
557 optimal coupling between $\tilde{\Phi}_\eta q_0$ and q^* . By taking the difference between (46) and (27),

$$\begin{bmatrix} x \\ \omega \end{bmatrix} = \begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix} + u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \zeta + \int_0^\eta \alpha_x(s) ds \right. \\ \left. \left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr + \int_0^\eta e^{-\gamma(s-\eta)} ds \right) \zeta + \int_0^\eta \alpha_x(s) + \alpha_v(s) ds \right].$$

558 Let us now analyze the Wasserstein distance between $\hat{\Phi}_\eta q_0$ and q^* ,

$$\begin{aligned}
& \mathcal{W}_2^2(\hat{\Phi}_\eta q_0, q^*) \\
& \leq \mathbb{E}_{r_1} \left\| \begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix} + u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \zeta + \int_0^\eta \alpha_x(s) ds \right] - \begin{bmatrix} x^* \\ \omega^* \end{bmatrix} \right\|^2 \\
& \leq \mathbb{E}_{r_1} \left\| \begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix} - \begin{bmatrix} x^* \\ \omega^* \end{bmatrix} \right\|^2 + u^2 \mathbb{E} \left\| \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \zeta + \int_0^\eta (\alpha_x(s) + \alpha_v(s)) ds \right] \right\|^2 \\
& \leq \mathcal{W}_2^2(\tilde{\Phi}_\eta q_0, q^*) + 4u^2 \left(\left(\int_0^\delta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right)^2 + \left(\int_0^\delta e^{-\gamma(s-\delta)} ds \right)^2 \right) \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\
& + u^2 \mathbb{E} \left[\left\| \int_0^\eta (\alpha_x(s)) ds \right\|^2 \right] + u^2 \mathbb{E} \left[\left\| \int_0^\eta (\alpha_x(s) + \alpha_v(s)) ds \right\|^2 \right] \\
& \leq \mathcal{W}_2^2(\tilde{\Phi}_\eta q_0, q^*) + 4u^2 \left(\frac{\eta^4}{4} + \eta^2 \right) \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + u^2 \mathbb{E} [\|\alpha_k^x\|^2] + u^2 \mathbb{E} [\|\alpha_k^v\|^2] \\
& \leq \mathcal{W}_2^2(\tilde{\Phi}_\eta q_0, q^*) + 5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 \left(\mathbb{E} \|\alpha_k^x\|^2 + \mathbb{E} \|\alpha_k^v\|^2 \right) \\
& \leq \mathcal{W}_2^2(\tilde{\Phi}_\eta q_0, q^*) + 5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 (A + B),
\end{aligned}$$

559 where the constant A, B are the uniform bounds of $\mathbb{E} [\|\alpha_k^x\|]$ and $\mathbb{E} [\|\alpha_k^v\|]$ respectively. Furthermore
560 from the corollary 7 in Cheng et al. [2018] we know that for any $i \in \{1, \dots, K\}$:

$$\mathcal{W}_2^2(\Phi_\eta q_i, q^*) \leq e^{-\eta/2\kappa_1} \mathcal{W}_2^2(q_i, q^*), \quad (47)$$

561 where $\kappa_1 = M/m_1$ 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

$$\mathcal{W}_2(\Phi_\eta q_i, \tilde{\Phi}_\eta q_i) \leq 2\mathcal{W}_2(\Phi_\eta p_i, \tilde{\Phi}_\eta p_i) \leq \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}.$$

563 By triangle inequality:

$$\begin{aligned}
\mathcal{W}_2(\tilde{\Phi}_\eta q_i, q^*) & \leq \mathcal{W}_2(\Phi_\eta q_i, \tilde{\Phi}_\eta q_i) + \mathcal{W}_2(\Phi_\eta q_i, q^*) \\
& \leq \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + e^{-\eta/2\kappa_1} \mathcal{W}_2(q_i, q^*),
\end{aligned}$$

564 further implies the following inequality:

$$\mathcal{W}_2^2(\hat{\Phi}_\eta q_i, q^*) \leq \left(e^{-\eta/2\kappa_1} \mathcal{W}_2(q_i, q^*) + \eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} \right)^2 + 5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 (A + B).$$

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

$$\begin{aligned}
\mathcal{W}_2(q_K, q^*) & \leq e^{-K\eta/2\kappa_1} \mathcal{W}_2(q_0, q^*) + \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \\
& + \frac{5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 (A + B)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{1 - e^{-\eta/2\kappa_1}} \sqrt{5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 (A + B)}}.
\end{aligned}$$

567 Finally by sandwich inequality we have:

$$\begin{aligned}
\mathcal{W}_2(p_K, p^*) & \leq 4e^{-K\eta/2\kappa_1} \mathcal{W}_2(q_0, q^*) + \frac{4\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \\
& + \frac{20u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8u^2 (A + B)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{1 - e^{-\eta/2\kappa_1}} \sqrt{5u^2 \eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2 (A + B)}}.
\end{aligned} \quad (48)$$

568 And in this case, we know that $\mathbb{E}[\|\alpha_k^x\|]$ and $\mathbb{E}[\|\alpha_k^y\|]$ can be bounded by $\frac{\Delta^2 d}{4}$. Finally, we can have:

$$\begin{aligned} \mathcal{W}_2(p_K, p^*) &\leq 4e^{-K\eta/2\kappa_1} \mathcal{W}_2(q_0, q^*) + \frac{4\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \\ &\quad + \frac{20u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2\right) + 4u^2\Delta^2 d}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{1 - e^{-\eta/2\kappa_1}} \sqrt{5u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2\right) + u^2\Delta^2 d}}. \end{aligned}$$

569 Now we let the first term less than $\epsilon/3$, from the lemma 13 in [Cheng et al., 2018] we know that
570 $\mathcal{W}_2(q_0, q^*) \leq 3 \left(\frac{d}{m_1} + \mathcal{D}^2\right)$. So we can choose K as the following,

$$K \leq \frac{2\kappa_1}{\eta} \log \left(36 \left(\frac{d}{m_1} + \mathcal{D}^2 \right) \right).$$

571 Next, we choose a stepsize $\eta \leq \frac{\epsilon\kappa_1^{-1}}{\sqrt{479232/5(d/m_1 + \mathcal{D}^2)}}$ to ensure the second term is controlled below
572 $\epsilon/3$. Since $1 - e^{-\eta/2\kappa_1} \geq \eta/4\kappa_1$ and definition of \mathcal{E}_K ,

$$4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \leq 4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{\eta/4\kappa_1} \leq 16\kappa_1 \left(\eta \sqrt{\frac{8\mathcal{E}_K}{5}} \right) \leq \epsilon/3.$$

573 Finally by choosing the stepsize satisfied that,

$$\eta \leq \frac{\epsilon^2}{2880\kappa_1 u \left(\frac{\Delta^2 d}{4} + \sigma^2\right)},$$

574 the third term can be bounded as:

$$\begin{aligned} &\frac{20u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) + 4u^2\Delta^2 d}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{1 - e^{-\eta/2\kappa_1}} \sqrt{5u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)}} \\ &\leq \frac{20u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) + 4u^2\Delta^2 d}{\sqrt{1 - e^{-\eta/2\kappa_1}} \sqrt{5u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)}} \leq \frac{20u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) + 4u^2\Delta^2 d}{\sqrt{\eta/4\kappa_1} \sqrt{5u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)}} \\ &\leq 4\sqrt{20\kappa_1 u^2 \eta \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)} + \frac{8u^2\Delta^2 d \sqrt{\kappa_1}}{\eta^{3/2} \sqrt{5u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)}} \\ &\leq \epsilon/3 + \frac{8u^2\Delta^2 d \sqrt{\kappa_1}}{\eta^{3/2} \sqrt{5u^2\eta^2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right)}}. \end{aligned}$$

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,

$$\begin{aligned} \mathbf{v}_t &= \mathbf{v}_0 - \int_0^t \gamma \mathbf{v}_s ds - u \int_0^t G_s ds + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s + \int_0^t \alpha_v(s) ds \\ \mathbf{x}_t &= \mathbf{x}_0 + \int_0^t \mathbf{v}_s ds + \int_0^t \alpha_x(s) ds, \end{aligned}$$

579 where $G_s = \sum_{k=0}^{\infty} Q_G(\nabla U(x'_k)) \mathbf{1}_{s \in (k\eta, (k+1)\eta)}$. And we define an intermediate process by let $\mathbf{v}'_t =$

580 $\mathbf{v}_t + \alpha_x(t)$:

$$\begin{aligned} v'_t &= v'_0 - \int_0^t \gamma(v'_s - \alpha_x(s)) ds - u \int_0^t G_s ds + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s + \int_0^t \left(\alpha_v(s) + \frac{1}{t} \alpha_x(t) \right) ds \\ x'_t &= x'_0 + \int_0^t v'_s ds. \end{aligned} \quad (49)$$

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

$$\begin{aligned} \mathbf{v}_t &= \mathbf{v}_0 - \int_0^t \gamma(\mathbf{v}_s - \alpha_x(s)) ds - u \int_0^t \nabla U(\mathbf{x}_s) ds + \sqrt{2\gamma u} \int_0^t e^{-\gamma(t-s)} dB_s \\ \mathbf{x}_t &= \mathbf{x}_0 + \int_0^t \mathbf{v}_s ds. \end{aligned} \quad (50)$$

582 Notice that the process x'_t has the same distribution with \mathbf{x}_t , thus in the following analysis we study
583 the convergence of the intermediate process $p'_k = (x'_{k\eta}, v'_{k\eta})$. By taking the difference of equation

584 (49) with (50) and the Girsanov formula, we can derive the Radon-Nikodym derivative of $\hat{P}_{K\eta}$ w.r.t
585 p'_K :

$$\begin{aligned} \frac{d\hat{p}_{K\eta}}{dp'_K} &= \exp \left\{ \sqrt{\frac{u}{2\gamma}} \int_0^T (\gamma \alpha_x(s) + \alpha_v(s) + \frac{1}{T} \alpha_x(T) + \nabla U(\mathbf{x}_s) - G_s) d\mathbf{B}_s \right. \\ &\quad \left. - \frac{u}{4\gamma} \int_0^T \|\gamma \alpha_x(s) + \alpha_v(s) + \frac{1}{T} \alpha_x(T) + \nabla U(\mathbf{x}_s) - G_s\|^2 ds \right\}. \end{aligned}$$

586 Thus the divergence can be bounded as:

$$\begin{aligned} D_{KL}(p_K \| \hat{p}_{K\eta}) &= \mathbb{E}_{p_K} \left[\log \left(\frac{d\hat{p}_{K\eta}}{dp_K} \right) \right] \\ &= \frac{u}{4\gamma} \int_0^T \mathbb{E} \left[\|\gamma \alpha_x(s) + \alpha_v(s) + \frac{1}{T} \alpha_x(T) + \nabla U(\mathbf{x}_s) - G_s\|^2 \right] ds \\ &= \frac{u}{4\gamma T} \mathbb{E} \left[\|\alpha_x(T)\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_v(s) + \alpha_x(s) + \nabla U(\mathbf{x}_s) - G_s\|^2 \right] ds \\ &\leq \frac{u}{4\gamma T \eta^2} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}}\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_v(s)\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_x(s)\|^2 \right] ds \\ &\quad + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - G_s\|^2 \right] ds \\ &\leq \frac{u}{4\gamma T \eta^2} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}}\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_k^{\mathbf{v}} / \eta\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}} / \eta\|^2 \right] ds \\ &\quad + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - Q_G(\nabla U(\mathbf{x}_k))\|^2 \right] ds \\ &\leq \frac{u}{4\gamma T \eta^2} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}}\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_k^{\mathbf{v}} / \eta\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}} / \eta\|^2 \right] ds \\ &\quad + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - Q_G(\nabla U(\mathbf{x}_k))\|^2 \right] ds. \end{aligned} \quad (51)$$

587 By assumption 1, we know that:

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] \leq M^2 \mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right].$$

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

$$\mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] \leq 3M^2\eta^2 \left(\mathbb{E} \left[\|\mathbf{v}'_k\|^2 \right] + u^2\eta^2 \mathbb{E} \left[\|Q_G(\nabla U(\mathbf{x}_k))\|^2 \right] + 2du \right).$$

589 Now we need to derive a uniform bound of $\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right]$ and $\mathbb{E} \left[\|\mathbf{v}'_k\|^2 \right]$.

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

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

591 *then for all $k > 0$ $\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right]$ and $\mathbb{E} \left[\|v_k\|^2 \right]$ can be bounded as follow:*

$$\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] \leq \mathcal{E} + C\Delta^2 d, \quad \mathbb{E} \left[\|v'_k\|^2 \right] \leq \gamma^2 \mathcal{E} / 2 + \gamma^2 C \Delta^2 d / 2,$$

592 *where constants \mathcal{E} and C are defined as:*

$$\mathcal{E} = \mathbb{E}[\mathcal{E}(\mathbf{x}_0, \mathbf{v}_0)] + \frac{54(4u + \gamma^2)u}{m_2\gamma^4} \sigma^2 + \frac{12(22u + \gamma)uM^3}{m_2\gamma^3} G^2 + \frac{96(d + b)uM}{m_2\gamma^2}$$

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

593

594 Thus,

$$\begin{aligned} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] &\leq 3M^2\eta^2 \left(\mathbb{E} \left[\|v_k\|^2 \right] + u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 + 2M^2 \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + 2G^2 \right) + 2du \right) \\ &\leq 3M^2\eta^2 \left(\gamma^2 \mathcal{E} / 2 + \gamma^2 C \Delta^2 d / 2 + u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 + 2M^2 \mathcal{E} + 2M^2 C \Delta^2 d + 2G^2 \right) + 2du \right) \\ &\leq 3M^2\eta^2 \left((\gamma^2 + 2u^2M^2) \mathcal{E} + (\gamma^2 + 2u^2M^2) C \Delta^2 d + u^2\sigma^2 + 2u^2G^2 + 2du \right). \end{aligned}$$

595 Now we can go back to the divergence of p_K and $\hat{p}_{K\eta}$,

$$D_{KL}(p_K \|\hat{p}_{K\eta})$$

$$\begin{aligned} &\leq \frac{u}{4\gamma T \eta^2} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}}\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_k^{\mathbf{v}} / \eta\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k^{\mathbf{x}} / \eta\|^2 \right] ds \\ &+ \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2M^2) \mathcal{E} + (\gamma^2 + 2u^2M^2) C \Delta^2 d + u^2\sigma^2 + 2u^2G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\ &\leq \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2M^2) \mathcal{E} + (\gamma^2 + 2u^2M^2) C \Delta^2 d + u^2\sigma^2 + 2u^2G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\ &+ \frac{u\Delta^2 d}{16\gamma T \eta^2} + \frac{uK\Delta^2 d}{8\gamma \eta} \\ &\leq \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2M^2) \mathcal{E} + u^2\sigma^2 + 2u^2G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \sigma^2 \\ &+ \left(\frac{u}{4\gamma} 3M^2 K \eta^3 C (\gamma^2 + 2u^2M^2) + \frac{uK\eta}{16\gamma} + \frac{u}{16\gamma T \eta^2} + \frac{uK}{8\gamma \eta} \right) \Delta^2 d \\ &=: C_0 K \eta^3 + C_1 K \eta \sigma^2 + C_2 K \Delta^2, \end{aligned}$$

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

$$C_0 = \frac{u}{4\gamma} 3M^2 \left((\gamma^2 + 2u^2M^2) \mathcal{E} + u^2\sigma^2 + 2u^2G^2 + 2du \right)$$

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

$$C_2 = \left(\frac{u}{4\gamma} 3M^2 \eta^3 C (\gamma^2 + 2u^2M^2) + \frac{u}{16\gamma} + \frac{u}{16\gamma T^2 \eta} + \frac{u}{8\gamma \eta} \right) d.$$

597 By the weighted CKP inequality and given $K\eta \geq 1$,

$$\begin{aligned} \mathcal{W}_2(p_K, \hat{p}_{K\eta}) &\leq \bar{\Lambda} \left(\sqrt{D_{KL}(p_K || \hat{p}_{K\eta})} + \sqrt[4]{D_{KL}(p_K || \hat{p}_{K\eta})} \right) \\ &\leq \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} \right) \sqrt{K\eta} + \tilde{C}_2 \sqrt{K\Delta}, \end{aligned} \quad (52)$$

598 where the constants are defined as:

$$\begin{aligned} \tilde{C}_0 &= \left(\sqrt{C_0} + \sqrt[4]{C_0} \right) \\ \tilde{C}_1 &= \left(\sqrt{C_1} + \sqrt[4]{C_1} \right) \\ \tilde{C}_2 &= \left(\sqrt{C_2} + \sqrt[4]{C_2} \right) \\ \tilde{A} &= \max \{ \sigma, \sqrt{\sigma} \}. \end{aligned}$$

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

$$\mathcal{W}_2(p_K, p^*) \leq \bar{\Lambda} \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} \right) \sqrt{K\eta} + \tilde{C}_2 \sqrt{K\eta} + \Gamma_0 e^{-\mu^* K\eta}. \quad (53)$$

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

$$\bar{\Lambda} \tilde{C}_0 \sqrt{K\eta^2} = \frac{\epsilon}{2} \quad \text{and} \quad \Gamma_0 e^{-\mu^* K\eta} = \frac{\epsilon}{2}. \quad (54)$$

601 Solving the equation (54), we can have

$$K\eta = \frac{\log \left(\frac{2\Gamma_0}{\epsilon} \right)}{\mu^*} \quad \text{and} \quad \eta = \frac{\epsilon^2}{4\bar{\Lambda}^2 \tilde{C}_0^2 K\eta}.$$

602 Combining these two we can have

$$\eta = \frac{\epsilon^2 \mu^*}{4\bar{\Lambda}^2 \tilde{C}_0^2 \log \left(\frac{2\Gamma_0}{\epsilon} \right)} \quad \text{and} \quad K = \frac{4\bar{\Lambda}^2 \tilde{C}_0^2 \log^2 \left(\frac{2\Gamma_0}{\epsilon} \right)}{\epsilon^2 (\mu^*)^2}.$$

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,

$$\begin{aligned} \mathbf{x}_{k+1} &= Q_W(\mathbf{x}_k - \eta \nabla \tilde{U}(\mathbf{x}_k) + \sqrt{2\eta} \xi_{k+1}) \\ &=: \mathbf{x}_k - \eta \nabla \tilde{U}(\mathbf{x}_k) + \sqrt{2\eta} \xi_{k+1} + \alpha_k, \end{aligned}$$

608 where α_k is define as:

$$\alpha_k = Q_W(\mathbf{x}_k - \eta \nabla \tilde{U}(\mathbf{x}_k) + \sqrt{2\eta} \xi_{k+1}) - \mathbf{x}_k - \eta \nabla \tilde{U}(\mathbf{x}_k) + \sqrt{2\eta} \xi_{k+1}.$$

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

$$\mathbf{x}_t = \mathbf{x}_0 - \int_0^t G_s ds + \sqrt{2} \int_0^t dB(s) + \int_0^t \alpha(s) ds,$$

610 where $G_s = \sum_{k=0}^{\infty} Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \mathbf{1}_{s \in (k\eta, (k+1)\eta)}$ and $\alpha(s) = \sum_{k=0}^{\infty} \alpha_k / \eta \mathbf{1}_{s \in (k\eta, (k+1)\eta)}$. By taking the

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

$$\frac{d\hat{p}_{K\eta}}{dp_K} = \exp \left\{ \frac{1}{2} \int_0^t (\nabla U(\mathbf{x}_s) - G_s - \alpha(s)) d\mathbf{B}s - \frac{1}{4} \int_0^t \|\nabla U(\mathbf{x}_s) - G_s - \alpha(s)\|^2 ds \right\}.$$

613 Thus, the divergence can be computed as:

$$\begin{aligned}
D_{KL}(p_K \|\hat{p}_{K\eta}) &= \frac{1}{4} \int_0^{K\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - G_s - \alpha(s)\|^2 \right] ds \\
&= \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - Q_G(\nabla \tilde{U}(\mathbf{x}_k)) - \alpha_k/\eta\|^2 \right] ds \\
&= \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - Q_G(\nabla \tilde{U}(\mathbf{x}_k))\|^2 \right] ds + \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k/\eta\|^2 \right] ds \\
&= \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_s) - \nabla U(\mathbf{x}_k)\|^2 \right] ds + \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - Q_G(\nabla \tilde{U}(\mathbf{x}_k))\|^2 \right] ds \\
&\quad + \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k/\eta\|^2 \right] ds \\
&\leq \frac{M^2}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right] ds + \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - Q_G(\nabla \tilde{U}(\mathbf{x}_k))\|^2 \right] ds \\
&\quad + \frac{1}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k/\eta\|^2 \right] ds. \tag{55}
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}_s - \mathbf{x}_k\|^2 \right] &\leq 3\eta^2 \mathbb{E} \left[\|\nabla U(\mathbf{x}_k)\|^2 \right] + 3\eta^2 \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - Q_G(\nabla \tilde{U}(\mathbf{x}_k))\|^2 \right] + 6\eta d \\
&\leq 3\eta^2 \left(M \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + G \right)^2 + 3\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 6\eta d.
\end{aligned}$$

615 Again, we need to derive a uniform bound of $\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right]$,

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}_{k+1}\|^2 \right] &= \mathbb{E} \left[\left\| \mathbf{x}_k - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right] + 2\mathbb{E} \left[\|\xi_{k+1}\|^2 \right] + \mathbb{E} \left[\|\alpha_k\|^2 \right] \\
&= \mathbb{E} \left[\left\| \mathbf{x}_k - \eta \nabla U(\mathbf{x}_k) + \eta \nabla U(\mathbf{x}_k) - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right] + 2\eta d + \mathbb{E} \left[\|\alpha_k\|^2 \right] \\
&= \mathbb{E} \left[\left\| \mathbf{x}_k - \eta \nabla U(\mathbf{x}_k) + \eta \nabla U(\mathbf{x}_k) - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right] + \mathbb{E} \left[\|\alpha_k\|^2 \right] + 2\eta d \\
&= \mathbb{E} \left[\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k)\|^2 \right] + \eta^2 \mathbb{E} \left[\|\nabla U(\mathbf{x}_k) - Q_G(\nabla \tilde{U}(\mathbf{x}_k))\|^2 \right] + \mathbb{E} \left[\|\alpha_k\|^2 \right] + 2\eta d.
\end{aligned}$$

616 By plugging in the inequality we derived before:

$$\mathbb{E} \left[\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k)\|^2 \right] \leq (1 - 2\eta m_2 + 2\eta^2 M^2) \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + 2\eta b + 2\eta^2 G^2.$$

617 we can have:

$$\mathbb{E} \left[\|\mathbf{x}_{k+1}\|^2 \right] \leq (1 - 2\eta m_2 + 2\eta^2 M^2) \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + 2\eta b + 2\eta^2 G^2 + \frac{\eta^2 \Delta^2 d}{4} + \eta^2 \sigma^2 + \mathbb{E} \left[\|\alpha_k\|^2 \right] + 2\eta d.$$

618 Thus for any $\eta \in (0, 1 \wedge \frac{m_2}{2M^2})$ and $1 - 2\eta m_2 + 2\eta^2 M^2 > 0$, we can bound $\mathbb{E} [\|\mathbf{x}_k\|^2]$ for any
 619 $k > 0$ as:

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_k\|^2] &\leq \mathbb{E} [\|\mathbf{x}_0\|^2] + \frac{1}{2(m_2 - \eta M^2)} \left(2b + 2G^2 + \frac{\Delta^2 d}{4} + \sigma^2 + 2d \right) + \frac{\mathbb{E} [\|\alpha_k\|^2]}{2\eta(m_2 - \eta M^2)} \\ &\leq \mathbb{E} [\|\mathbf{x}_0\|^2] + \frac{1}{m_2} \left(2b + 2G^2 + \frac{\Delta^2 d}{4} + \sigma^2 + 2d \right) + \frac{\mathbb{E} [\|\alpha_k\|^2]}{\eta m_2} \\ &\leq \mathcal{E} + \frac{\Delta^2 d}{4m_2} + \frac{\mathbb{E} [\|\alpha_k\|^2]}{\eta m_2}, \end{aligned}$$

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

$$\mathcal{E} = \mathbb{E} [\|\mathbf{x}_0\|^2] + \frac{1}{m_2} (2b + 2G^2 + \sigma^2 + 2d).$$

621 Thus, we can have,

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_s - \mathbf{x}_k\|^2] &\leq 6\eta^2 \left(\mathcal{E} + \frac{\Delta^2 d}{4m_2} + \frac{\mathbb{E} [\|\alpha_k\|^2]}{\eta m_2} \right) + 6\eta^2 G^2 + 3\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 6\eta d \\ &\leq \bar{\mathcal{E}}\eta + 3\eta^2 \sigma^2 + \frac{6 + 3m_2}{4m_2} \eta^2 \Delta^2 d + \frac{6\eta \mathbb{E} [\|\alpha_k\|^2]}{m_2}. \end{aligned}$$

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

$$\begin{aligned} D_{KL}(p_K || \hat{p}_{K\eta}) &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M\sigma^2 K\eta^3}{4} + \frac{(6 + 3m_2) M\Delta^2 d}{16m_2} K\eta^3 + \frac{6M\mathbb{E} [\|\alpha_k\|^2] K\eta^2}{4m_2} + \frac{1}{4} \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) K\eta + \frac{K\mathbb{E} [\|\alpha_k\|^2]}{4\eta} \\ &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M + 1}{4} \sigma^2 K\eta + \frac{((6 + 3m_2) M + m_2) d}{16m_2} \Delta^2 K\eta + \left(\frac{6M\eta}{4m_2} + \frac{1}{4\eta} \right) K\mathbb{E} [\|\alpha_k\|^2]. \end{aligned}$$

623 By the fact that $\mathbb{E} [\|\alpha_k\|^2] \leq \frac{\Delta^2 d}{4}$, we can further bound the divergence as:

$$\begin{aligned} D_{KL}(p_K || \hat{p}_{K\eta}) &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M + 1}{4} \sigma^2 K\eta + \left(\frac{((12 + 3m_2) M + m_2) d}{16m_2} + \frac{d}{16\eta} \right) \Delta^2 K \\ &=: C_0 K\eta^2 + C_1 \sigma^2 K\eta + C_2 \Delta^2 K, \end{aligned}$$

624 where the constants are defined as:

$$\begin{aligned} C_0 &= \frac{M\bar{\mathcal{E}}}{4} \\ C_1 &= \frac{3M + 1}{4} \\ C_2 &= \left(\frac{((12 + 3m_2) M + m_2) d}{16m_2} + \frac{d}{16\eta} \right). \end{aligned}$$

625 We are ready to bound the Wasserstein distance,

$$\begin{aligned} \mathcal{W}_2^2(p_K, \hat{p}_{K\eta}) &\leq (12 + 8(\kappa_0 + 2b + 2d)) \left[(C_0 + \sqrt{C_0} + (C_1 + \sqrt{C_1}) A) (K\eta)^2 + (C_2 + \sqrt{C_2}) \Delta K^2 \eta \right] \\ &=: (\tilde{C}_0^2 \sqrt{\eta} + \tilde{C}_1^2 A) (K\eta)^2 + \tilde{C}_2^2 \Delta K^2 \eta, \end{aligned}$$

626 where the constants are defined as:

$$\begin{aligned} A &= \max \{ \sigma^2, \sqrt{\sigma^2} \} \\ \tilde{C}_0^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_0 + \sqrt{C_0}) \\ \tilde{C}_1^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_1 + \sqrt{C_1}) \\ \tilde{C}_2^2 &= (12 + 8(\kappa_0 + 2b + 2d)) (C_2 + \sqrt{C_2}). \end{aligned}$$

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

$$\begin{aligned} \mathcal{W}_2(\hat{p}_{K\eta}, p^*) &\leq \sqrt{2C_{LS} \left(\log \|p_0\|_\infty + \frac{d}{2} \log \frac{3\pi}{m\beta} + \beta \left(\frac{M\kappa_0}{3} + B\sqrt{\kappa_0} + G_0 + \frac{b}{2} \log 3 \right) \right)} e^{-K\eta/\beta C_{LS}} \\ &=: \widetilde{C}_3 e^{-K\eta/\beta C_{LS}} \end{aligned}$$

628 Finally, we can have

$$\mathcal{W}_2(p_K, p^*) \leq \left(\widetilde{C}_0 \eta^{1/4} + \widetilde{C}_1 \sqrt{A} \right) K\eta + \widetilde{C}_2 \sqrt{\Delta} \sqrt{K^2 \eta} + \widetilde{C}_3 e^{-K\eta/\beta C_{LS}}. \quad (56)$$

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

$$\widetilde{C}_0 K \eta^{5/4} \leq \frac{\epsilon}{2} \quad \text{and} \quad \widetilde{C}_3 e^{-K\eta/\beta C_{LS}} = \frac{\epsilon}{2}. \quad (57)$$

630 Solving the (57), we can have

$$K\eta = C_{LS} \log \left(\frac{2\widetilde{C}_3}{\epsilon} \right) \quad \text{and} \quad \eta \leq \frac{\epsilon^4}{16\widetilde{C}_0^4 (K\eta)^4}.$$

631 Combining these two we can have

$$\eta \leq \frac{\epsilon^4}{16\widetilde{C}_0^4 C_{LS}^4 \log^4 \left(\frac{2\widetilde{C}_3}{\epsilon} \right)} \quad \text{and} \quad K \geq \frac{16\widetilde{C}_0^4 C_{LS}^5 \log^5 \left(\frac{2\widetilde{C}_3}{\epsilon} \right)}{\epsilon^4}.$$

632 Plugging K and η 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 update rule is the following,

$$\begin{aligned} \mathbf{v}_{k+1} &= Q^{vc} \left(v_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G \left(\nabla \tilde{U}(\mathbf{x}_k) \right), Var_v, \Delta \right) \\ \mathbf{x}_{k+1} &= Q^{vc} \left(\mathbf{x}_k + \gamma^{-1} (1 - e^{-\gamma\eta}) v_k + u\gamma^{-2} (\gamma\eta + e^{-\gamma\eta} - 1) Q_G(\nabla \tilde{U}(\mathbf{x}_k)), Var_x, \Delta \right). \end{aligned} \quad (58)$$

636 If we let $\alpha_k^{\mathbf{x}}$ and $\alpha_k^{\mathbf{v}}$ denote the quantization error,

$$\begin{aligned} \alpha_k^{\mathbf{v}} &= Q^{vc} \left(v_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G \left(\nabla \tilde{U}(\mathbf{x}_k) \right), Var_v, \Delta \right) - \left(v_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G(\nabla \tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{v}} \right) \\ \alpha_k^{\mathbf{x}} &= Q^{vc} \left(\mathbf{x}_k + \gamma^{-1} (1 - e^{-\gamma\eta}) v_k + u\gamma^{-2} (\gamma\eta + e^{-\gamma\eta} - 1) Q_G(\nabla \tilde{U}(\mathbf{x}_k)), Var_x, \Delta \right) \\ &\quad - \left(\mathbf{x}_k + \gamma^{-1} (1 - e^{-\gamma\eta}) v_k + u\gamma^{-2} (\gamma\eta + e^{-\gamma\eta} - 1) Q_G(\nabla \tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{x}} \right), \end{aligned}$$

637 we can rewrite the update rule as:

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

638 Next, we first derive a uniform bound of $\mathbb{E} \left[\|\alpha_k^{\mathbf{v}}\|^2 \right]$. 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 \mathcal{G} and the quantized stochastic gradients at
641 x satisfies the following

$$\mathbb{E} \left[\left\| Q_G(\nabla \tilde{U}(x)) \right\|^2 \right] \leq \mathcal{G}^2.$$

642 By the definition of the variance corrected quantization function Q^{vc} , when $Var_v > \rho_0 = \frac{\Delta^2}{4}$, if
 643 we let ψ_k denote $v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k))$,

$$\begin{aligned} & \mathbb{E} \left[\|\alpha_k^{\mathbf{y}}\|^2 \middle| \psi_k \right] \\ = & \mathbb{E} \left[\left\| \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) + \sqrt{Var_v} \xi_k \right. \right. \\ & \left. \left. - Q^d \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \sqrt{Var_v - \rho_0} \xi_k \right) - \text{sign}(r)c \right\|^2 \middle| \psi_k \right] \end{aligned}$$

644 Let

$$\begin{aligned} b = & Q^d \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \sqrt{Var_v - \rho_0} \xi_k \right) \\ & - \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \sqrt{Var_v - \rho_0} \xi_k \right), \end{aligned}$$

645 then

$$\begin{aligned} & \mathbb{E} \left[\|\alpha_k^{\mathbf{y}}\|^2 \middle| \psi_k \right] \\ = & \mathbb{E} \left[\left\| \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) + \sqrt{Var_v} \xi_k \right. \right. \\ & \left. \left. - \left(v_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \sqrt{Var_v - \rho_0} \xi_k \right) - b - \text{sign}(r)c \right\|^2 \middle| \psi_k \right] \\ = & \mathbb{E} \left[\left\| \sqrt{Var_v} \xi_k - \sqrt{Var_v - \rho_0} \xi_k - b - \text{sign}(r)c \right\|^2 \middle| \psi_k \right] \\ \leq & \mathbb{E} \left[\left\| \sqrt{Var_v} \xi_k - \sqrt{Var_v - \rho_0} \xi_k \right\|^2 \right] + \mathbb{E} \left[\|b + \text{sign}(r)c\|^2 \middle| \psi_k \right] \\ \leq & 2Var_v d - \rho_0 d + \rho_0 d \\ \leq & 4\gamma u d \eta. \end{aligned} \tag{59}$$

646 When $Var_v < \frac{\Delta^2}{4}$,

$$\begin{aligned} & \mathbb{E}[\|\alpha_k^{\mathbf{y}}\|^2] \\ = & \mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - \mathbf{v}_{k+1} + \sqrt{Var_v} \xi_k \right\|^2 \right] \\ = & \mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - \mathbf{v}_{k+1} \right\|^2 \right] + \mathbb{E} \left[\left\| \sqrt{Var_v} \xi_k \right\|^2 \right] \\ \leq & \max \left(2\mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - Q^s \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) \right\|^2 \right], 2Var_v \right) \end{aligned} \tag{60}$$

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

$$\begin{aligned} & \mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - Q^s \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) \right\|^2 \right] \\ \leq & \Delta (1 - e^{-\gamma\eta}) \mathbb{E} \left[\left\| v_k - u\gamma^{-1} Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\|_1 \right] \\ \leq & \Delta (1 - e^{-\gamma\eta}) \sqrt{d} \left(\mathbb{E}[\|v_k\|] + \mathbb{E} \left[\left\| Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\| \right] \right). \end{aligned}$$

648 Now we need to derive a uniform bound of $\mathbb{E}[\|v_k\|]$, by the update rule, we know that,

$$\begin{aligned} \mathbb{E}[\|\mathbf{v}_{k+1}\|^2] &= \mathbb{E} \left[\left\| \mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1}(1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) + \xi_k^{\mathbf{y}} + \alpha_k^{\mathbf{y}} \right\|^2 \right] \\ &\leq (1 + \gamma\eta/2) (1 - \gamma\eta/2)^2 \mathbb{E}[\|v_k\|^2] + \left(\frac{2}{\gamma\eta} + 1 \right) u^2 \eta^2 \mathbb{E} \left[\left\| Q_G(\nabla\tilde{U}) \right\|^2 \right] + 2\gamma u d \eta + \mathbb{E}[\|\alpha_k^{\mathbf{y}}\|^2] \\ &\leq (1 - \gamma\eta/2) \mathbb{E}[\|v_k\|^2] + 3u^2 \eta / \gamma \mathcal{G}^2 + 2\gamma u d \eta + \mathbb{E}[\|\alpha_k^{\mathbf{y}}\|^2]. \end{aligned}$$

649 When $\mathbb{E} \left[\|\alpha_k^y\|^2 \right] \leq 2Var_v d < 4\gamma ud\eta$, the inequality can be further written as:

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{v}_{k+1}\|^2 \right] &\leq (1 - \gamma\eta/2) \mathbb{E} \left[\|v_k\|^2 \right] + 3u^2\eta/\gamma\mathcal{G}^2 + 6\gamma ud\eta \\ &\leq \mathbb{E} \left[\|\mathbf{v}_0\|^2 \right] + \frac{6u^2\eta\mathcal{G}^2}{\gamma^2\eta} + \frac{12\gamma ud\eta}{\gamma\eta} \\ &\leq \mathbb{E} \left[\|\mathbf{v}_0\|^2 \right] + \frac{6u^2\eta\mathcal{G}^2}{\gamma^2} + 12ud. \end{aligned}$$

650 If $\mathbb{E} \left[\|\alpha_k^y\|^2 \right] \leq 2\mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - Q^s \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) \right\|^2 \right]$,
651 the inequality can be written as:

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{v}_{k+1}\|^2 \right] &\leq (1 - \gamma\eta/2) \mathbb{E} \left[\|v_k\|^2 \right] + 3u^2\eta/\gamma\mathcal{G}^2 + 2\gamma ud\eta + 2\Delta (1 - e^{-\gamma\eta}) \sqrt{d} \left(\mathbb{E} [\|v_k\|] + \mathbb{E} \left[\left\| Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\| \right] \right) \\ &\leq (1 - \gamma\eta/2) \mathbb{E} \left[\|v_k\|^2 \right] + 3u^2\eta/\gamma\mathcal{G}^2 + 2\gamma ud\eta + 2\Delta\gamma\eta\sqrt{d} \left(\sqrt{\mathbb{E} \left[\|v_k\|^2 \right]} + \mathcal{G} \right) \\ &\leq \left(\sqrt{1 - \gamma\eta/2} \sqrt{\mathbb{E} \left[\|v_k\|^2 \right]} + \frac{\Delta\gamma\eta\sqrt{d}}{\sqrt{1 - \gamma\eta/2}} \right)^2 + 3u^2\eta/\gamma\mathcal{G}^2 + 2\gamma ud\eta + 2\Delta\gamma\eta\sqrt{d}\mathcal{G}. \end{aligned}$$

652 Thus,

$$\begin{aligned} \mathbb{E} [\|v_k\|] &\leq \sqrt{\mathbb{E} \left[\|\mathbf{v}_0\|^2 \right]} + \frac{\Delta\gamma\eta\sqrt{d}}{(1 - \sqrt{1 - \gamma\eta/2}) \sqrt{1 - \gamma\eta/2}} + \frac{3u^2\eta/\gamma\mathcal{G}^2 + 2\gamma ud\eta + 2\Delta\gamma\eta\sqrt{d}\mathcal{G}}{\frac{\Delta\gamma\eta\sqrt{d}}{\sqrt{1 - \gamma\eta/2}} + \sqrt{\gamma\eta/2} (3u^2\eta/\gamma\mathcal{G}^2 + 2\gamma ud\eta + 2\Delta\gamma\eta\sqrt{d}\mathcal{G})} \\ &\leq \sqrt{\mathbb{E} \left[\|\mathbf{v}_0\|^2 \right]} + \frac{\Delta\gamma\eta\sqrt{d}}{1 - \gamma\eta/2} + \sqrt{6u^2/\gamma^2\mathcal{G}^2 + 4ud + 4\Delta\sqrt{d}\mathcal{G}} \\ &\leq \sqrt{\mathbb{E} \left[\|\mathbf{v}_0\|^2 \right]} + \Delta\sqrt{d} + \sqrt{6u^2/\gamma^2\mathcal{G}^2 + 4ud + 4\Delta\sqrt{d}\mathcal{G}}. \end{aligned}$$

653 Finally, we can have:

$$\begin{aligned} \mathbb{E} [\|v_k\|] &\leq \max \left\{ \sqrt{\mathbb{E} \left[\|\mathbf{v}_0\|^2 \right]} + \Delta\sqrt{d} + \sqrt{6u^2/\gamma^2\mathcal{G}^2 + 4ud + 4\Delta\sqrt{d}\mathcal{G}}, \right. \\ &\quad \left. \sqrt{\mathbb{E} \left[\|\mathbf{v}_0\|^2 \right]} + \sqrt{\frac{6u^2\eta\mathcal{G}^2}{\gamma^2} + \sqrt{12ud}} \right\} =: A'. \end{aligned}$$

654 Thus, we can have,

$$\begin{aligned} \mathbb{E} \left[\left\| \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) - Q^s \left(\mathbf{v}_k e^{-\gamma\eta} - u\gamma^{-1} (1 - e^{-\gamma\eta}) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right) \right\|^2 \right] \\ \leq \Delta\gamma\eta\sqrt{d} (A' + \mathcal{G}), \end{aligned}$$

655 and we can bound the $\mathbb{E} \left[\|\alpha_k^y\|^2 \right]$ as,

$$\begin{aligned} \mathbb{E} \left[\|\alpha_k^y\|^2 \right] &\leq \max \left\{ \Delta\gamma\eta\sqrt{d} (A' + \mathcal{G}), 4\gamma ud\eta \right\} \\ &= \gamma\eta \max \left\{ \Delta\sqrt{d} (A' + \mathcal{G}), 4ud \right\} \\ &=: \gamma\eta A. \end{aligned} \tag{61}$$

656 Now we bound the $\mathbb{E} \left[\|\alpha_k^x\|^2 \right]$. When $Var_x \geq \rho_0$, as the same analysis in (59) we can show,

$$\mathbb{E} \left[\|\alpha_k^x\|^2 \right] \leq 2Var_x d \leq 4ud\eta^2.$$

657 If $Var_x < \rho_0$, and let $\mu_x = \mathbf{x}_k + \gamma^{-1} (1 - e^{-\gamma\eta}) v_k + u\gamma^{-2} (\gamma\eta + e^{-\gamma\eta} - 1) Q_G(\nabla\tilde{U}(\mathbf{x}_k))$, by
658 the same analysis in (60) we can have:

$$\begin{aligned} \mathbb{E} \left[\|\alpha_k^x\|^2 \right] \\ \leq \max \left\{ 2\mathbb{E} \left[\|\mu_x - Q^s(\mu_x)\|^2 \right], 2Var_x d \right\}. \end{aligned}$$

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

$$\begin{aligned} \mathbb{E} \left[\|\mu_x - Q^s(\mu_x)\|^2 \right] &\leq \Delta \mathbb{E} \left[\left\| \gamma^{-1} (1 - e^{-\gamma\eta}) v_k + u\gamma^{-2} (\gamma\eta + e^{-\gamma\eta} - 1) Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\|_1 \right] \\ &\leq \Delta\eta \mathbb{E} [\|v_k\|_1] + \frac{u\eta^2}{2} \mathbb{E} \left[\left\| Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\|_1 \right] \\ &\leq \Delta\eta\sqrt{d} \mathbb{E} [\|v_k\|] + \frac{u\eta^2}{2} \sqrt{d} \mathbb{E} \left[\left\| Q_G(\nabla\tilde{U}(\mathbf{x}_k)) \right\| \right] \\ &\leq \Delta\eta\sqrt{d}A' + \frac{u\eta^2}{2} \sqrt{d}\mathcal{G}. \end{aligned}$$

660 Thus, we can have,

$$\begin{aligned} \mathbb{E} \left[\|\alpha_k^x\|^2 \right] &\leq \max \left\{ 2\Delta\eta\sqrt{d}A' + u\eta^2\sqrt{d}\mathcal{G}, 4ud\eta^2 \right\} \\ &\leq \eta \max \left\{ 2\Delta\sqrt{d}A' + u\eta\sqrt{d}\mathcal{G}, 4ud\eta \right\} \\ &=: \eta B. \end{aligned} \tag{62}$$

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

$$\begin{aligned} \mathcal{W}_2(p_K, p^*) &\leq 4e^{-K\eta/2\kappa_1} \mathcal{W}_2(q_0, q^*) + \frac{4\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \\ &\quad + \frac{20u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8u^2\eta(\gamma A + B)}{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}} + \sqrt{1 - e^{-\eta/\kappa_1}} \sqrt{5u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2\eta(\gamma A + B)}}. \end{aligned}$$

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

663 $\mathcal{W}_2(q_0, q^*) \leq 3 \left(\frac{d}{m_1} + \mathcal{D}^2 \right)$. So we can choose K as the following,

$$K \leq \frac{2\kappa_1}{\eta} \log \left(36 \left(\frac{d}{m_1} + \mathcal{D}^2 \right) \right).$$

664 Next, we choose a stepsize $\eta \leq \frac{\epsilon\kappa_1^{-1}}{\sqrt{479232/5(d/m_1 + \mathcal{D}^2)}}$ to ensure the second term is controlled below

665 $\epsilon/3$. Since $1 - e^{-\eta/2\kappa_1} \geq \eta/4\kappa_1$ and definition of \mathcal{E}_K ,

$$4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{1 - e^{-\eta/2\kappa_1}} \leq 4 \frac{\eta^2 \sqrt{\frac{8\mathcal{E}_K}{5}}}{\eta/4\kappa_1} \leq 16\kappa_1 \left(\eta \sqrt{\frac{8\mathcal{E}_K}{5}} \right) \leq \epsilon/3.$$

666 Finally by choosing the stepsize satisfied that,

$$\eta \leq \frac{\epsilon^2}{2880\kappa_1 u \left(\frac{\Delta^2 d}{4} + \sigma^2 \right)},$$

667 the third term can be bounded as:

$$\begin{aligned}
& \frac{20u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8u^2\eta(\gamma A + B)}{\eta^2 \sqrt{\frac{8\mathcal{E}\kappa}{5}} + \sqrt{1 - e^{-\eta/\kappa_1}} \sqrt{5u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2\eta(\gamma A + B)}} \\
& \leq \frac{20u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8u^2\eta(\gamma A + B)}{\sqrt{1 - e^{-\eta/\kappa_1}} \sqrt{5u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2\eta(\gamma A + B)}} \leq \frac{20u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8u^2\eta(\gamma A + B)}{\sqrt{\eta/4\kappa_1} \sqrt{5u^2\eta^2 \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2u^2\eta(\gamma A + B)}} \\
& \leq 4\sqrt{20u^2\kappa_1\eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 8\kappa_1u^2(\gamma A + B)} \\
& \leq \epsilon/3 + 8\sqrt{2\kappa_1u^2(\gamma A + B)}.
\end{aligned}$$

668 This complete the proof.

669 **D.8 Proof of Theorem 9**

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

$$\mathbb{E} \left[\|\alpha_k^y\|^2 \right] \leq \gamma\eta\mathcal{A}, \quad (63)$$

671 where $A = \max \left\{ \Delta\sqrt{d}(A' + \mathcal{G}), 4ud \right\}$. By the analysis in (59), we know that if $\text{Var}_{\mathbf{x}}^{hmc} \geq \frac{\Delta^2}{4}$,
672 we can have

$$\mathbb{E} \left[\|\alpha_k^x\|^2 \right] \leq 4ud\eta^2 \quad (64)$$

673 by (62), if $\text{Var}_{\mathbf{x}}^{hmc} < \frac{\Delta^2}{4}$,

$$\mathbb{E} \left[\|\alpha_k^x\|^2 \right] \leq \eta B, \quad (65)$$

674 where $B = \max \left\{ 2\Delta\sqrt{d}A' + u\eta\sqrt{d}\mathcal{G}, 4ud\eta \right\}$. Thus, we can define the following:

$$\mathbb{E} \left[\|\alpha_k^x\|^2 \right] = \eta\mathcal{B}, \quad (66)$$

675 where \mathcal{B} is defined as:

$$\mathcal{B} = \begin{cases} 4ud\eta, & \text{if } \text{Var}_{\mathbf{x}}^{hmc} \geq \frac{\Delta^2}{4} \\ B, & \text{else.} \end{cases}$$

676 Combining the bound of $\mathbb{E} \left[\|\alpha_k^x\|^2 \right]$, $\mathbb{E} \left[\|\alpha_k^y\|^2 \right]$ with (51), we can show,

$$\begin{aligned}
& D_{KL}(p_K \|\hat{p}_{K\eta}) \\
& \leq \frac{u}{4\gamma T \eta^2} \mathbb{E} \left[\|\alpha_k^x\|^2 \right] + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\gamma \alpha_k^y / \eta\|^2 \right] ds + \frac{u}{4\gamma} \sum_{k=0}^K \int_{k\eta}^{(k+1)\eta} \mathbb{E} \left[\|\alpha_k^x / \eta\|^2 \right] ds \\
& + \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2 M^2) \mathcal{E} + (\gamma^2 + 2u^2 M^2) C \Delta^2 d + u^2 \sigma^2 + 2u^2 G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\
& \leq \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2 M^2) \mathcal{E} + (\gamma^2 + 2u^2 M^2) C \Delta^2 d + u^2 \sigma^2 + 2u^2 G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\
& + \frac{u\mathcal{B}}{4\gamma T} + \frac{uK\mathcal{A}}{4} + \frac{uK\mathcal{B}}{4\gamma} \\
& \leq \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2 M^2) \mathcal{E} + (\gamma^2 + 2u^2 M^2) C \Delta^2 d + u^2 \sigma^2 + 2u^2 G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) \\
& + \frac{uK\mathcal{A}}{4} + \frac{uK\mathcal{B}}{2\gamma} \\
& \leq \frac{u}{4\gamma} 3M^2 K \eta^3 \left((\gamma^2 + 2u^2 M^2) \mathcal{E} + u^2 \sigma^2 + 2u^2 G^2 + 2du \right) + \frac{u}{4\gamma} K \eta \sigma^2 + \frac{u}{16\gamma} K \eta \Delta^2 d + \frac{uK\mathcal{A}}{4} + \frac{uK\mathcal{B}}{2\gamma} \\
& =: C_0 K \eta^3 + C_1 K \eta \sigma^2 + C_2 K \eta \Delta^2 + C_3 K \mathcal{A} + C_4 K \mathcal{B},
\end{aligned}$$

677 where the constants are defined as

$$\begin{aligned}
C_0 &= \frac{u}{4\gamma} 3M^2 \left((\gamma^2 + 2u^2 M^2) \mathcal{E} + u^2 \sigma^2 + 2u^2 G^2 + 2du \right) \\
C_1 &= \frac{u}{4\gamma} \\
C_2 &= \frac{u}{16\gamma} d \\
C_3 &= \frac{u}{4} \\
C_4 &= \frac{u}{2\gamma}.
\end{aligned}$$

678 By the weighted CKP inequality and given $K\eta \geq 1$,

$$\begin{aligned}
\mathcal{W}_2(p_K, \hat{p}_{K\eta}) &\leq \bar{\Lambda} \left(\sqrt{D_{KL}(p_K \|\hat{p}_{K\eta})} + \sqrt[4]{D_{KL}(p_K \|\hat{p}_{K\eta})} \right) \\
&\leq \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} + \tilde{C}_2 \sqrt{\Delta} \right) \sqrt{K\eta} + \tilde{C}_3 \sqrt{K\mathcal{A}} + \tilde{C}_4 \sqrt{K\mathcal{B}},
\end{aligned}$$

679 where the constants are defined as:

$$\begin{aligned}
\tilde{C}_0 &= \bar{\Lambda} \left(\sqrt{C_0} + \sqrt[4]{C_0} \right) \\
\tilde{C}_1 &= \bar{\Lambda} \left(\sqrt{C_1} + \sqrt[4]{C_1} \right) \\
\tilde{C}_2 &= \bar{\Lambda} \left(\sqrt{C_2} + \sqrt[4]{C_2} \right) \\
\tilde{C}_3 &= \bar{\Lambda} \left(\sqrt{C_3} + \sqrt[4]{C_3} \right) \\
\tilde{C}_4 &= \bar{\Lambda} \left(\sqrt{C_4} + \sqrt[4]{C_4} \right) \\
\tilde{A}^2 &= \bar{\Lambda} \max \left\{ \sigma^2, \sqrt{\sigma^2} \right\}.
\end{aligned}$$

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

$$\mathcal{W}_2(p_K, p^*) \leq \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \tilde{A} \right) \sqrt{K\eta} + \tilde{C}_2 \sqrt{K\eta} \Delta + \tilde{C}_3 \sqrt{K\mathcal{A}} + \tilde{C}_4 \sqrt{K\mathcal{B}} + \Gamma_0 e^{-\mu^* K\eta}. \quad (67)$$

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

$$\bar{\Lambda} \widetilde{C}_0 \sqrt{K\eta^2} = \frac{\epsilon}{2} \quad \text{and} \quad \Gamma_0 e^{-\mu^* K\eta} = \frac{\epsilon}{2}. \quad (68)$$

682 Solving the equation (68), we can have

$$K\eta = \frac{\log\left(\frac{2\Gamma_0}{\epsilon}\right)}{\mu^*} \quad \text{and} \quad \eta = \frac{\epsilon^2}{4\bar{\Lambda}^2 \widetilde{C}_0^2 K\eta}.$$

683 Combining these two we can have

$$\eta = \frac{\epsilon^2 \mu^*}{4\bar{\Lambda}^2 \widetilde{C}_0^2 \log\left(\frac{2\Gamma_0}{\epsilon}\right)} \quad \text{and} \quad K = \frac{4\bar{\Lambda}^2 \widetilde{C}_0^2 \log^2\left(\frac{2\Gamma_0}{\epsilon}\right)}{\epsilon^2 (\mu^*)^2}.$$

684 Plugging in (67) completes the proof.

685 **D.9 Proof of Theorem 12**

686 Recall that the update of VC SGLDLP-L is

$$\begin{aligned} \mathbf{x}_{k+1} &= Q^{vc}\left(\mathbf{x}_k - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)), 2\eta, \Delta\right) \\ &= \mathbf{x}_k - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)) + \sqrt{2\eta} \xi_k + \alpha_k, \end{aligned}$$

687 where α_k is defined as

$$\alpha_k = Q^{vc}\left(\mathbf{x}_k - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)), 2\eta, \Delta\right) - \mathbf{x}_k - \eta Q_G(\nabla \tilde{U}(\mathbf{x}_k)) + \sqrt{2\eta} \xi_k.$$

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

$$\begin{aligned} \mathbb{E}\left[\|\alpha_k\|^2\right] &\leq \max(2\Delta\eta G, 5\eta d) \\ &=: \eta A. \end{aligned}$$

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

$$\begin{aligned} D_{KL}(p_K \|\hat{p}_{K\eta}) &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M+1}{4} \sigma^2 K\eta + \frac{((6+3m_2)M+m_2)d}{16m_2} \Delta^2 K\eta + \left(\frac{6M\eta}{4m_2} + \frac{1}{4\eta}\right) K\mathbb{E}\left[\|\alpha_k\|^2\right] \\ &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M+1}{4} \sigma^2 K\eta + \frac{((6+3m_2)M+m_2)d}{16m_2} \Delta^2 K\eta + \left(\frac{6M\eta}{4m_2} + \frac{1}{4\eta}\right) K\eta A \\ &\leq \frac{M\bar{\mathcal{E}}}{4} K\eta^2 + \frac{3M+1}{4} \sigma^2 K\eta + \frac{((6+3m_2)M+m_2)d}{16m_2} \Delta^2 K\eta + \frac{6M+m_2}{m_2} KA \\ &=: C_0 K\eta^2 + C_1 K\eta\sigma^2 + C_2 K\eta\Delta^2 + C_3 KA, \end{aligned}$$

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

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

691 We are ready to bound the Wasserstein distance,

$$\begin{aligned} \mathcal{W}_2^2(p_K, \hat{p}_{K\eta}) &\leq (12 + 8(\kappa_0 + 2b + 2d)) \left[\left((C_0 + \sqrt{C_0})\eta + (C_1 + \sqrt{C_1})\tilde{A} \right) (K\eta)^2 + \left(C_2 + \sqrt{C_2} \right) \Delta (K\eta)^2 \right. \\ &\quad \left. + \left(C_3 + \sqrt{C_3} \right) AK^2\eta \right] \\ &=: \left(\widetilde{C}_0^2 \eta + \widetilde{C}_1^2 \tilde{A} + \widetilde{C}_2^2 \Delta \right) (K\eta)^2 + \widetilde{C}_3^2 AK^2\eta, \end{aligned}$$

692 where the constants are defined as:

$$\begin{aligned}
\tilde{A} &= \max \left\{ \sigma^2, \sqrt{\sigma^2} \right\} \\
\mathcal{A} &= \max \left\{ A, \sqrt{A} \right\} \\
\tilde{C}_0^2 &= (12 + 8(\kappa_0 + 2b + 2d)) \left(C_0 + \sqrt{C_0} \right) \\
\tilde{C}_1^2 &= (12 + 8(\kappa_0 + 2b + 2d)) \left(C_1 + \sqrt{C_1} \right) \\
\tilde{C}_2^2 &= (12 + 8(\kappa_0 + 2b + 2d)) \left(C_2 + \sqrt{C_2} \right) \\
\tilde{C}_3^2 &= (12 + 8(\kappa_0 + 2b + 2d)) \left(C_3 + \sqrt{C_3} \right).
\end{aligned}$$

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

$$\begin{aligned}
\mathcal{W}_2(\hat{p}_{K\eta}, p^*) &\leq \sqrt{2C_{LS} \left(\log \|p_0\|_\infty + \frac{d}{2} \log \frac{3\pi}{m\beta} + \beta \left(\frac{M\kappa_0}{3} + B\sqrt{\kappa_0} + G_0 + \frac{b}{2} \log 3 \right) \right)} e^{-K\eta/\beta C_{LS}} \\
&=: \tilde{C}_4 e^{-K\eta/\beta C_{LS}}
\end{aligned}$$

694 Finally, we can have

$$\mathcal{W}_2(p_K, p^*) \leq \left(\tilde{C}_0 \sqrt{\eta} + \tilde{C}_1 \sqrt{A} + \tilde{C}_2 \sqrt{\Delta} \right) K\eta + \tilde{C}_3 \sqrt{\mathcal{A}} \sqrt{K^2 \eta} + \tilde{C}_4 e^{-K\eta/\beta C_{LS}}. \quad (69)$$

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

$$\tilde{C}_0 K \eta^{5/4} = \frac{\epsilon}{2} \quad \text{and} \quad \tilde{C}_3 e^{-K\eta/\beta C_{LS}} = \frac{\epsilon}{2}. \quad (70)$$

696 Solving the (70), we can have

$$K\eta = C_{LS} \log \left(\frac{2\tilde{C}_3}{\epsilon} \right) \quad \text{and} \quad \eta = \frac{\epsilon^4}{16\tilde{C}_0^4 (K\eta)^4}.$$

697 Combining these two we can have

$$\eta = \frac{\epsilon^4}{16\tilde{C}_0^4 C_{LS}^4 \log^4 \left(\frac{2\tilde{C}_3}{\epsilon} \right)} \quad \text{and} \quad K = \frac{16\tilde{C}_0^4 C_{LS}^5 \log^5 \left(\frac{2\tilde{C}_3}{\epsilon} \right)}{\epsilon^4}.$$

698 Plugging K and η into (69) completes the proof.

699 E Technical Proofs

700 E.1 Proof of Lemma 13

701 *Proof.* By the definition of ξ in (25)

$$\begin{aligned}
\|\mathbb{E}\xi\|^2 &= \|\mathbb{E}\tilde{g}(\mathbf{x}) - \mathbb{E}\nabla U(\mathbf{x})\|^2 \\
&= \|\mathbb{E}\nabla U(Q_w(\mathbf{x})) - \mathbb{E}\nabla U(\mathbf{x})\|^2 \\
&\leq \mathbb{E} \left[\|\nabla U(Q_w(\mathbf{x})) - \nabla U(\mathbf{x})\|^2 \right] \\
&\leq M^2 \mathbb{E} \left[\|Q_w(\mathbf{x}) - \nabla U(\mathbf{x})\|^2 \right] \\
&\leq M \frac{\Delta^2 d}{4}.
\end{aligned}$$

702 We also know that from the definition that

$$\begin{aligned}
\mathbb{E} \|\xi\|^2 &= \mathbb{E} \|\tilde{g}(\mathbf{x}) - \nabla U(\mathbf{x})\|^2 \\
&= \mathbb{E} \left\| Q_G(\nabla \tilde{U}(Q_W(\mathbf{x}))) - \nabla \tilde{U}(Q_W(\mathbf{x})) + \nabla \tilde{U}(Q_W(\mathbf{x})) - \nabla U(Q_W(\mathbf{x})) + \nabla U(Q_W(\mathbf{x})) - \nabla U(\mathbf{x}) \right\|^2 \\
&= \mathbb{E} \left\| Q_G(\nabla \tilde{U}(Q_W(\mathbf{x}))) - \nabla \tilde{U}(Q_W(\mathbf{x})) \right\|^2 + \mathbb{E} \left\| \nabla \tilde{U}(Q_W(\mathbf{x})) - \nabla U(Q_W(\mathbf{x})) \right\|^2 + \mathbb{E} \|\nabla U(Q_W(\mathbf{x})) - \nabla U(\mathbf{x})\|^2 \\
&\leq \frac{\Delta^2 d}{4} + \sigma^2 + M^2 \mathbb{E} \|Q_W(\mathbf{x}) - \mathbf{x}\|^2 \\
&\leq (M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2,
\end{aligned}$$

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

704

□

705 E.2 Proof of Lemma 14

706 *Proof.* Let Γ_1 be the set of all couplings between $\tilde{\Phi}_\eta q_0$ and q^* and Γ_2 be the set of all couplings
707 between $\hat{\Phi}_\eta q_0$ and q^* . Let r_1 be the optimal coupling between $\tilde{\Phi}_\eta q_0$ and q^* , i.e.

$$\mathbb{E}_{(\theta, \phi) \sim r_1} [\|\theta - \phi\|^2] = \mathcal{W}_2^2(\tilde{\Phi}_\eta q_0, q^*).$$

708 Let $\left(\begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix}, \begin{bmatrix} x^* \\ \omega^* \end{bmatrix} \right) \sim r_1$. We define the random variable $\begin{bmatrix} x \\ \omega \end{bmatrix}$ as

$$\begin{bmatrix} x \\ \omega \end{bmatrix} = \begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix} + u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \xi \right. \\ \left. \left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr + \int_0^\delta e^{-\gamma(s-\eta)} ds \right) \xi \right].$$

709 By equation (29), $\left(\begin{bmatrix} x \\ \omega \end{bmatrix}, \begin{bmatrix} x^* \\ \omega^* \end{bmatrix} \right)$ define a valid coupling between $\Phi_\eta q_0$ and q^* . Now we can analyze
710 the Wasserstein distance between $\Phi_\eta q_0$ and q^* .

$$\begin{aligned}
\mathcal{W}_2^2(\hat{\Phi}_\eta q_0, q^*) &\leq \mathbb{E}_{r_1} \left[\left\| \begin{bmatrix} \tilde{x} \\ \tilde{\omega} \end{bmatrix} + u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \xi \right. \right. \right. \\ &\quad \left. \left. \left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr + \int_0^\delta e^{-\gamma(s-\eta)} ds \right) \xi \right] - \begin{bmatrix} x^* \\ \omega^* \end{bmatrix} \right\|^2 \right] \\ &\leq \mathbb{E}_{r_1} \left[\left\| \begin{bmatrix} \tilde{x} - x^* \\ \tilde{\omega} - \omega^* \end{bmatrix} + u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) \mathbb{E}\xi \right. \right. \right. \\ &\quad \left. \left. \left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr + \int_0^\delta e^{-\gamma(s-\eta)} ds \right) \mathbb{E}\xi \right] \right\|^2 \right] \\ &\quad + \mathbb{E}_{r_1} \left[\left\| u \left[\left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr \right) (\xi - \mathbb{E}\xi) \right. \right. \right. \\ &\quad \left. \left. \left(\int_0^\eta \left(\int_0^r e^{-\gamma(s-r)} ds \right) dr + \int_0^\delta e^{-\gamma(s-\eta)} ds \right) (\xi - \mathbb{E}\xi) \right] \right\|^2 \right] \\ &\leq \left(\mathcal{W}_2(\tilde{\Phi}_\eta q_0, q^*) + 2u\sqrt{\eta^4/4 + \eta^2} \|\mathbb{E}\xi\| \right)^2 + 4u^2(\eta^4/4 + \eta^2) \mathbb{E}_{r_1} [\|\xi - \mathbb{E}\xi\|^2] \\ &\leq \left(\mathcal{W}_2(\tilde{\Phi}_\eta q_0, q^*) + \sqrt{5}/2u\eta\sqrt{d}M\Delta \right)^2 + 5u^2\eta^2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right).
\end{aligned} \tag{71}$$

711

□

712 E.3 Proof of Lemma 15

713 *Proof.* In order to get the upper bound of $\|\mathbf{x}_k\|$ and $\|\mathbf{v}_k\|$, we bound the Lyapunov function
714 $\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)$. By the smooth Assumption 1, we know

$$U(\mathbf{x}_{k+1}) - U(x^*) \leq U(\mathbf{x}_k) + \langle \nabla U(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + M^2/2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 - U(x^*).$$

715 Recall the definition of the Lyapunov function

$$\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1}) = \|\mathbf{x}_{k+1}\|^2 + \|\mathbf{x}_{k+1} + 2\mathbf{v}_{k+1}/\gamma\|^2 + 8u(U(\mathbf{x}_{k+1}) - U(x^*))/\gamma^2.$$

716 For the first two terms we have

$$\begin{aligned}\|\mathbf{x}_{k+1}\|^2 &= \|\mathbf{x}_k\|^2 + 2\langle \mathbf{x}_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \\ \|\mathbf{x}_{k+1} + 2\mathbf{v}_{k+1}/\gamma\|^2 &= \|\mathbf{x}_k + 2\mathbf{v}_k/\gamma\|^2 + 2\langle \mathbf{x}_k + 2\mathbf{v}_k/\gamma, \mathbf{x}_{k+1} - \mathbf{x}_k + 2(\mathbf{v}_{k+1} - \mathbf{v}_k)/\gamma \rangle \\ &\quad + \|\mathbf{x}_{k+1} - \mathbf{x}_k + 2(\mathbf{v}_{k+1} - \mathbf{v}_k)/\gamma\|^2.\end{aligned}$$

717 This implies the following:

$$\begin{aligned}\mathbb{E}[\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \mathbb{E}[\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] + 4\mathbb{E}[\langle \mathbf{x}_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle] + \frac{4}{\gamma}\mathbb{E}[\langle \mathbf{x}_k, \mathbf{v}_{k+1} - \mathbf{v}_k \rangle] + \frac{4}{\gamma}\mathbb{E}[\langle \mathbf{v}_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle] \\ &\quad (72) \\ &\quad + \frac{8}{\gamma^2}\mathbb{E}[\langle \mathbf{v}_k, \mathbf{v}_{k+1} - \mathbf{v}_k \rangle] + \frac{8u}{\gamma^2}\mathbb{E}[\langle \nabla U(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + M/2\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2] \\ &\quad + \mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2] + \mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}_k + 2(\mathbf{v}_{k+1} - \mathbf{v}_k)/\gamma\|^2].\end{aligned}$$

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

$$\begin{aligned}\mathbb{E}[\langle \mathbf{x}_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle] &= \frac{1 - e^{-\gamma\eta}}{\gamma}\mathbb{E}[\langle \mathbf{x}_k, \mathbf{v}_k \rangle] + \frac{u(\gamma\eta + e^{-\gamma\eta} - 1)}{\gamma^2}\mathbb{E}[\langle \mathbf{x}_k, \tilde{g}(\mathbf{x}_k) \rangle], \\ \mathbb{E}[\langle \mathbf{x}_k, \mathbf{v}_{k+1} - \mathbf{v}_k \rangle] &= -(1 - e^{-\gamma\eta})\mathbb{E}[\langle \mathbf{x}_k, \mathbf{v}_k \rangle] - \frac{u(1 - e^{-\gamma\eta})}{\gamma}\mathbb{E}[\langle \mathbf{x}_k, \tilde{g}(\mathbf{x}_k) \rangle], \\ \mathbb{E}[\langle \mathbf{v}_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle] &= \frac{1 - e^{-\gamma\eta}}{\gamma}\mathbb{E}[\|\mathbf{v}_k\|^2] + \frac{u(\gamma\eta + e^{-\gamma\eta} - 1)}{\gamma^2}\mathbb{E}[\langle \mathbf{v}_k, \tilde{g}(\mathbf{x}_k) \rangle], \\ \mathbb{E}[\langle \mathbf{v}_k, \mathbf{v}_{k+1} - \mathbf{v}_k \rangle] &= -(1 - e^{-\gamma\eta})\mathbb{E}[\|\mathbf{v}_k\|^2] - \frac{u(1 - e^{-\gamma\eta})}{\gamma}\mathbb{E}[\langle \mathbf{v}_k, \tilde{g}(\mathbf{x}_k) \rangle].\end{aligned}$$

719 Plug into the (72) yields:

$$\begin{aligned}\mathbb{E}[\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \mathbb{E}[\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{4u(2 - \gamma\eta - 2e^{-\gamma\eta})}{\gamma^2}\mathbb{E}[\langle \mathbf{x}_k, \tilde{g}(\mathbf{x}_k) \rangle] - \frac{4(1 - e^{-\gamma\eta})}{\gamma^2}\mathbb{E}[\|\mathbf{v}_k\|^2] \\ &\quad + \frac{4u(\gamma\eta + e^{-\gamma\eta} - 1)}{\gamma^3}\mathbb{E}[\langle \mathbf{v}_k, \tilde{g}(\mathbf{x}_k) \rangle] + \frac{8u(1 - e^{-\gamma\eta})}{\gamma^3}\mathbb{E}[\langle \mathbf{v}_k, \nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k) \rangle] \\ &\quad + \frac{8u^2(\gamma\eta + e^{-\gamma\eta} - 1)}{\gamma^4}\mathbb{E}[\langle \nabla U(\mathbf{x}_k), \tilde{g}(\mathbf{x}_k) \rangle] + \left(\frac{4Mu}{\gamma^2} + 3\right)\mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2] \\ &\quad + \frac{8}{\gamma^2}\mathbb{E}[\|\mathbf{v}_{k+1} - \mathbf{v}_k\|^2].\end{aligned}\quad (73)$$

720 By Assumption 3, we know that $\langle \mathbf{x}_k, \nabla U(\mathbf{x}_k) \rangle \geq m_2\|\mathbf{x}_k\|^2 - b$. We then assume $\eta \leq 1/(8\gamma)$ and
721 use the inequality $-x \leq e^{-x} - 1 \leq x^2/2 - x$ for any $x \geq 0$, it follows that

$$\begin{aligned}& - \frac{4u(2 - \gamma\eta - 2e^{-\gamma\eta})}{\gamma^2}\mathbb{E}[\langle \mathbf{x}_k, \tilde{g}(\mathbf{x}_k) \rangle] \\ &= - \frac{4u(2 - \gamma\eta - 2e^{-\gamma\eta})}{\gamma^2}(\mathbb{E}[\langle \mathbf{x}_k, \nabla U(\mathbf{x}_k) \rangle] + \mathbb{E}[\langle \mathbf{x}_k, \tilde{g}(\mathbf{x}_k) - \nabla U(\mathbf{x}_k) \rangle]) \\ &\leq - \frac{4u(2 - \gamma\eta - 2e^{-\gamma\eta})}{\gamma^2}(m_2\mathbb{E}[\|\mathbf{x}_k\|^2] - b) + \frac{4u(2 - \gamma\eta - 2e^{-\gamma\eta})}{\gamma^2}\left(\frac{1}{8}\mathbb{E}[\|\mathbf{x}_k\|^2] + 2\mathbb{E}[\|\tilde{g}(\mathbf{x}_k) - \nabla U(\mathbf{x}_k)\|^2]\right) \\ &\leq - \frac{3m_2u\eta}{\gamma}\mathbb{E}[\|\mathbf{x}_k\|^2] + \frac{4u\eta b}{\gamma} + \frac{8u\eta}{\gamma}\mathbb{E}[\|\tilde{g}(\mathbf{x}_k) - \nabla U(\mathbf{x}_k)\|^2],\end{aligned}$$

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 $\gamma\eta - (\gamma\eta)^2 \leq 2 - \gamma\eta - 2e^{-\gamma\eta} \leq \gamma\eta$. Again by Young's
724 inequality and the update rule in (3) we have:

$$\begin{aligned}\mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2] &\leq 2\eta^2\mathbb{E}[\|\mathbf{v}_k\|^2] + u^2\eta^4/2\mathbb{E}[\|\tilde{g}(\mathbf{x}_k)\|^2] + \mathbb{E}[\|\xi_k^x\|^2] \\ \mathbb{E}[\|\mathbf{v}_{k+1} - \mathbf{v}_k\|^2] &\leq 2\gamma^2\eta^2\mathbb{E}[\|\mathbf{v}_k\|^2] + 2u^2\eta^2\mathbb{E}[\|\tilde{g}(\mathbf{x}_k)\|^2] + \mathbb{E}[\|\xi_k^v\|^2].\end{aligned}$$

725 It is easy to verify the fact that $\mathbb{E} \left[\|\xi_k^v\|^2 \right] \leq 2\gamma u d \eta$ and $\mathbb{E} \left[\|\xi_k^x\|^2 \right] \leq 2u d \eta^2$. Thus,

$$\begin{aligned} & \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] \\ & \leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um\eta^2}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{3(1 - e^{-\gamma\eta}) - \eta^2(8Mu + u\gamma + 22\gamma^2)}{\gamma^2} \mathbb{E} [\|\mathbf{v}_k\|^2] \\ & \quad + \frac{36u^2\eta^2 + 2\gamma u\eta^2 + (4Mu + 3\gamma^2)\eta^4}{2\gamma^2} \mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2] + \frac{2u^2\eta^2}{\gamma^2} \mathbb{E} [\|\nabla U(\mathbf{x}_k)\|^2] \\ & \quad + \frac{8u\eta(\gamma^2 + 2u)}{\gamma^3} \mathbb{E} [\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2] + \frac{(8Mu + 6\gamma^2)u d \eta^2 + 4(4d + b)u\gamma\eta}{\eta^2}. \end{aligned}$$

726 If we set

$$\eta \leq \min \left\{ \frac{\gamma}{4(8Mu + u\gamma + 22\gamma^2)}, \sqrt{\frac{4u^2}{4Mu + 3\gamma^2}}, \frac{6\gamma b u}{(4Mu + 3\gamma^2)d} \right\},$$

727 we can obtain the following,

$$\begin{aligned} \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] & \leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] + \frac{(20u + \gamma)u\eta^2}{\gamma^2} \mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2] \\ & \quad + \frac{2u^2\eta^2}{\gamma^2} \mathbb{E} [\|\nabla U(\mathbf{x}_k)\|^2] + \frac{8u\eta(\gamma^2 + 2u)}{\gamma^3} \mathbb{E} [\|\nabla U(\mathbf{x}_k) - \tilde{g}(\mathbf{x}_k)\|^2] + \frac{16(d + b)u\eta}{\gamma}. \end{aligned} \quad (74)$$

728 Furthermore we can bound $\mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2]$ by the following analysis:

$$\begin{aligned} \mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2] & \leq 2\mathbb{E} [\|\tilde{g}(\mathbf{x}_k) - \nabla U(\mathbf{x}_k)\|^2] + 2\mathbb{E} [\|\nabla U(\mathbf{x}_k)\|^2] \\ & \leq 2 \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) + 4M^2 \mathbb{E} [\|\mathbf{x}_k\|^2] + 4G^2, \end{aligned} \quad (75)$$

729 where G^2 is the bound of the gradient at 0, i.e. $\|\nabla U(0)\|^2 \leq G^2$. Thus we can have:

$$\begin{aligned} \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] & \leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] + \frac{(21u + \gamma)4M^2u\eta^2}{\gamma^2} \mathbb{E} [\|\mathbf{x}_k\|^2] \\ & \quad + \left(\frac{2(20u + \gamma)u\eta^2}{\gamma^2} + \frac{8u\eta(\gamma^2 + 2u)}{\gamma^3} \right) \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\ & \quad + \frac{(21u + \gamma)4u\eta^2}{\gamma^2} G^2 + \frac{16(d + b)u\eta}{\gamma}. \end{aligned}$$

730 If we set the stepsize

$$\eta \leq \min \left\{ \frac{\gamma m_2}{12(21u + \gamma)M^2}, \frac{8(\gamma^2 + 2u)}{(20u + \gamma)\gamma} \right\},$$

731 then we have:

$$\begin{aligned} \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] & \leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{8um_2\eta}{3\gamma} \mathbb{E} [\|\mathbf{x}_k\|^2] - \frac{2\eta}{\gamma} \mathbb{E} [\|\mathbf{v}_k\|^2] \\ & \quad + \left(\frac{16u\eta(\gamma^2 + 2u)}{\gamma^3} \right) \left((M^2 + 1) \frac{\Delta^2 d}{4} + \sigma^2 \right) \\ & \quad + \frac{(21u + \gamma)4u\eta^2}{\gamma^2} G^2 + \frac{16(d + b)u\eta}{\gamma}. \end{aligned}$$

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

$$\mathcal{E}(x, v) \leq 5/2 \|x\|^2 + \frac{12}{\gamma^2} + \frac{2uM}{\gamma^2} \left(3\|x\|^2 + 6\|x^*\|^2 \right).$$

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

$$\mathcal{E}(x, v) \leq \frac{16uM}{\gamma^2} \|x\|^2 + \frac{12}{\gamma^2} \|v\|^2 + \frac{12uM}{\gamma^2} \|x^*\|^2. \quad (76)$$

735 Thus,

$$\begin{aligned} \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \left(1 - \frac{\gamma m_2 \eta}{6M}\right) \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] + \left(\frac{16u\eta(\gamma^2 + 2u)}{\gamma^3}\right) \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) \\ &\quad + \frac{(21u + \gamma)4u\eta^2}{\gamma^2} G^2 + \frac{16(d+b)u\eta}{\gamma}. \end{aligned}$$

736 Finally we show that

$$\begin{aligned} \sup_{k \geq 0} \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] &\leq \mathbb{E} [\mathcal{E}(x_0, v_0)] + \frac{6M}{\gamma m_2 \eta} \left(\frac{16u\eta(\gamma^2 + 2u)}{\gamma^3}\right) \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) \\ &\quad + \frac{6M}{\gamma m_2 \eta} \frac{(21u + \gamma)4u\eta^2}{\gamma^2} G^2 + \frac{6M}{\gamma m_2 \eta} \frac{16(d+b)u\eta}{\gamma} \\ &\leq \mathbb{E} [\mathcal{E}(x_0, v_0)] + \frac{96u(\gamma^2 + 2u)}{m_2 \gamma^4} \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) + \frac{24(21u + \gamma)uM}{m_2 \gamma^3} G^2 + \frac{96(d+b)uM}{m_2 \gamma^2} \\ &\leq \bar{\mathcal{E}} + C_0 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right), \end{aligned} \quad (77)$$

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

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_k\|^2] &\leq \bar{\mathcal{E}} + C_0 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) \\ \mathbb{E} [\|\mathbf{v}_k\|^2] &\leq \gamma^2 \bar{\mathcal{E}}/2 + \gamma^2 C_0/2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right). \end{aligned}$$

740 Combining with equation (75) we can bound $\mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2]$ as:

$$\mathbb{E} [\|\tilde{g}(\mathbf{x}_k)\|^2] \leq 2 \left((M^2 + 1)\frac{\Delta^2 d}{4} + \sigma^2\right) + 4M^2 \bar{\mathcal{E}} + 4G^2. \quad (78)$$

741 □

742 E.4 Proof of Lemma 16

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

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_{k+1}\|^2] &= \mathbb{E} [\|\mathbf{x}_k - \eta \tilde{g}(\mathbf{x}_k)\|^2] + \sqrt{8\eta} \mathbb{E} [\langle \mathbf{x}_k - \eta \tilde{g}(\mathbf{x}_k), \xi_{k+1} \rangle] + 2\eta \mathbb{E} [\|\xi_{k+1}\|^2] \\ &= \mathbb{E} [\|\mathbf{x}_k - \eta \tilde{g}(\mathbf{x}_k)\|^2] + 2\eta d \\ &= \mathbb{E} [\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k) - \eta (\tilde{g}(\mathbf{x}_k) - \nabla U(Q_W(\mathbf{x}_k))) - \eta (\nabla U(Q_W(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k))\|^2] + 2\eta d \\ &= \mathbb{E} [\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k) - \eta (\nabla U(Q_W(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k))\|^2] + \eta^2 \mathbb{E} [\|\tilde{g}(\mathbf{x}_k) - \nabla U(Q_W(\mathbf{x}_k))\|^2] + 2\eta d \\ &= (\mathbb{E} [\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k)\|] + \eta \mathbb{E} [\|\nabla U(Q_W(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k)\|])^2 + \eta^2 \frac{\Delta^2 d}{4} + 2\eta d. \end{aligned}$$

744 We know the fact that:

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k)\|^2] &= \mathbb{E} [\|\mathbf{x}_k\|^2] - 2\eta \mathbb{E} [\langle \mathbf{x}_k, \nabla U(\mathbf{x}_k) \rangle] + \eta^2 \mathbb{E} [\|\nabla U(\mathbf{x}_k)\|^2] \\ &= \mathbb{E} [\|\mathbf{x}_k\|^2] + 2\eta (b - m_2 \mathbb{E} [\|\mathbf{x}_k\|^2]) + 2\eta^2 (M^2 \mathbb{E} [\|\mathbf{x}_k\|^2] + G^2) \\ &= (1 - 2\eta m_2 + 2\eta^2 M^2) \mathbb{E} [\|\mathbf{x}_k\|^2] + 2\eta b + 2\eta^2 G^2. \end{aligned}$$

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

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{x}_{k+1}\|^2 \right] &\leq (1+c) \mathbb{E} \left[\|\mathbf{x}_k - \eta \nabla U(\mathbf{x}_k)\|^2 \right] + \left(1 + \frac{1}{c}\right) \eta^2 \mathbb{E} \left[\|\nabla U(Q_W(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k)\|^2 \right] + \eta^2 \frac{\Delta^2 d}{4} + 2\eta d \\ &\leq (1 - \eta m_2 + \eta^2 M^2) \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + \frac{1 - \eta m_2 + \eta^2 M^2}{\eta m_2 - \eta^2 M^2} \frac{M^2 \eta^2 \Delta^2 d}{4} + \frac{1 - \eta m_2 + \eta^2 M^2}{1 - 2\eta m_2 + 2\eta^2 M^2} (2\eta b + 2\eta^2 G^2) \\ &\quad + \eta^2 \frac{\Delta^2 d}{4} + 2\eta d. \end{aligned}$$

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

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] &\leq \mathbb{E} \left[\|x_0\|^2 \right] + \frac{1 - \eta m_2 + \eta^2 M^2}{\eta^2 (m_2 - \eta M^2)^2} \frac{M^2 \eta^2 \Delta^2 d}{4} + \frac{1 - \eta m_2 + \eta^2 M^2}{\eta (1 - 2\eta m_2 + 2\eta^2 M^2) (m_2 - \eta M^2)} (2\eta b + 2\eta^2 G^2) \\ &\quad + \frac{1}{\eta (m_2 - \eta M^2)} \left(\eta^2 \frac{\Delta^2 d}{4} + 2\eta d \right) \\ &= \mathbb{E} \left[\|x_0\|^2 \right] + \frac{1 - \eta m_2 + \eta^2 M^2}{(m_2 - \eta M^2)^2} \frac{M^2 \Delta^2 d}{4} + \frac{1 - \eta m_2 + \eta^2 M^2}{(1 - 2\eta m_2 + 2\eta^2 M^2) (m_2 - \eta M^2)} (2b + 2\eta G^2) \\ &\quad + \frac{1}{m_2 - \eta M^2} \left(\eta \frac{\Delta^2 d}{4} + 2d \right) \\ &\leq \mathbb{E} \left[\|x_0\|^2 \right] + \frac{2M^2}{m_2} \frac{\Delta^2 d}{4} + \frac{2}{m_2} (2b + 2\eta G^2) + \frac{2}{m_2} \left(\eta \frac{\Delta^2 d}{4} + 2d \right). \end{aligned}$$

748 Now if we let $\mathcal{E} = \mathbb{E} \left[\|x_0\|^2 \right] + \frac{M}{m_2} (2b + 2\eta G^2 + 2d)$, then we can write:

$$\mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] \leq \mathcal{E} + \frac{2(M^2 + 1)}{m_2} \frac{\Delta^2 d}{4}.$$

749

□

750 E.5 Proof of Lemma 17

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

$$\eta \leq \min \left\{ \frac{\gamma}{4(8Mu + u\gamma + 22\gamma^2)}, \sqrt{\frac{4u^2}{4Mu + 3\gamma^2}}, \frac{6\gamma bu}{(4Mu + 3\gamma^2)d} \right\},$$

752 we can obtain the following,

$$\begin{aligned} \mathbb{E} [\mathcal{E}(\mathbf{x}_{k+1}, \mathbf{v}_{k+1})] &\leq \mathbb{E} [\mathcal{E}(\mathbf{x}_k, \mathbf{v}_k)] - \frac{3um_2\eta}{\gamma} \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] - \frac{2\eta}{\gamma} \mathbb{E} \left[\|\mathbf{v}_k\|^2 \right] + \frac{(20u + \gamma)u\eta^2}{\gamma^2} \mathbb{E} \left[\left\| Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right] \\ &\quad + \frac{2u^2\eta^2}{\gamma^2} \mathbb{E} \left[\|\nabla U(\mathbf{x}_k)\|^2 \right] + \frac{8u\eta(\gamma^2 + 2u)}{\gamma^3} \mathbb{E} \left[\left\| \nabla U(\mathbf{x}_k) - Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right] + \frac{16(d+b)u\eta}{\gamma}. \end{aligned} \tag{79}$$

753 By assumption 1, we can bound $\mathbb{E} \left[\left\| Q_G(\nabla \tilde{U}(\mathbf{x}_k)) \right\|^2 \right]$ by the following,

$$\begin{aligned} \mathbb{E} \left[\left\| Q_G(\nabla U(\mathbf{x}_k)) \right\|^2 \right] &= \mathbb{E} \left[\left\| Q_G(\nabla \tilde{U}(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k) + \nabla U(\mathbf{x}_k) - \nabla U(0) + \nabla U(0) \right\|^2 \right] \\ &\leq \mathbb{E} \left[\left\| Q_G(\nabla \tilde{U}(\mathbf{x}_k)) - \nabla U(\mathbf{x}_k) \right\|^2 \right] + 2\mathbb{E} \left[\left\| \nabla U(\mathbf{x}_k) - \nabla U(0) \right\|^2 \right] + 2\mathbb{E} \left[\left\| \nabla U(0) \right\|^2 \right] \\ &\leq \left(\frac{\Delta^2 d}{4} + \sigma^2 \right) + 2M^2 \mathbb{E} \left[\|\mathbf{x}_k\|^2 \right] + 2G^2. \end{aligned}$$

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 $\eta \leq \frac{\gamma m_2}{6(22u + \gamma)M^2}$, 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.

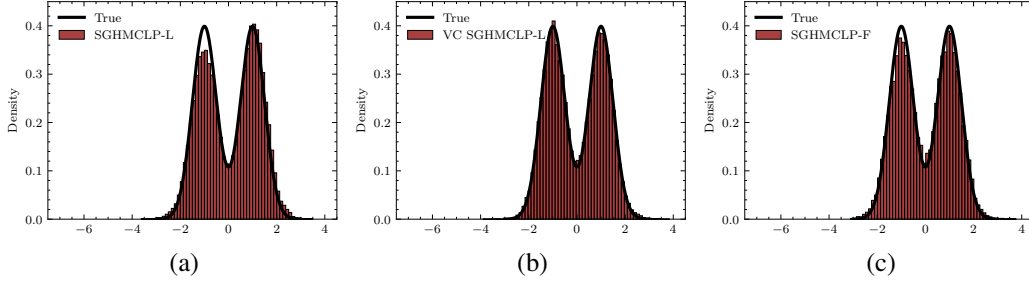


Figure 5: Low-precision SGHMC with stepsize equal to 0.01 on a Gaussian mixture distribution. (a): SGHMCLP-L. (b): VC SGHMCLP-L. (c): SGHMCLP-F.

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 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
 779 function. 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 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 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.

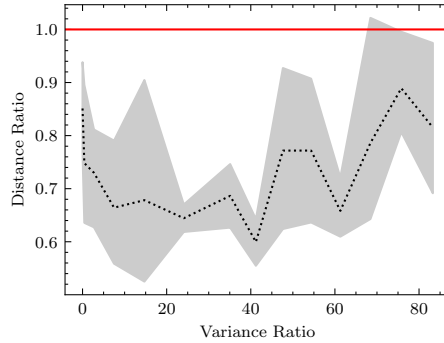


Figure 4: Wasserstein Distance Ratio of VC SGHMCLP-L & SGHMCLP-L (Smaller is better). The dashed line is the 2-Wasserstein distance to the target distribution ratio between the sample obtained by VC SGHMCLP-L and SGHMCLP-L.

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.

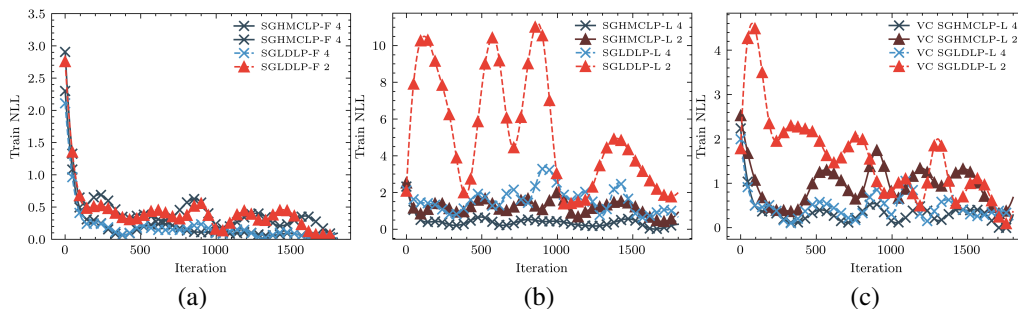


Figure 6: Training NLL of low-precision SGHMC and SGLD on MLP with MNIST in terms of different numbers of fractional bits. (a): Methods with Full-Precision Gradients Accumulators. (b): Methods with Low-Precision Gradients Accumulators. (c): Variance corrected quantization. The low-precision SGHMC adopted with full-precision gradient accumulators achieves comparable results with SGLD. However, when adopted with low-precision gradient accumulators and variance-corrected quantization SGHMC shows more robustness to quantization error especially when the number of representable bits is low.

Table 4: Test errors (%) of Low-precision gradient accumulators on CIFAR with ResNet-18.

	CIFAR-10	CIFAR-100
32-bit Float		
SGD	4.73 ± 0.10	22.34 ± 0.22
SGLD	4.52 ± 0.07	22.40 ± 0.04
SGHMC	4.78 ± 0.08	22.37 ± 0.04
8-bit Fixed Point		
SGD	8.50 ± 0.22	28.42 ± 0.35
SGLD	7.81 ± 0.07	27.15 ± 0.35
VC SGLD	7.03 ± 0.23	26.73 ± 0.12
SGHMC	6.63 ± 0.10	26.57 ± 0.10
VC SGHMC	6.60 ± 0.06	26.43 ± 0.19
8-bit Block Float Point		
SGD	5.86 ± 0.18	26.75 ± 0.11
SGLD	5.75 ± 0.05	26.11 ± 0.38
VC SGLD	5.51 ± 0.01	25.14 ± 0.11
SGHMC	5.38 ± 0.06	25.29 ± 0.03
VC SGHMC	5.15 ± 0.08	24.45 ± 0.16