

000 001 002 003 004 005 FASTER SAMPLING FROM GIBBS DISTRIBUTIONS 006 WITH QUANTUM VARIANCE REDUCTION 007 008 009

010 **Anonymous authors**
011 Paper under double-blind review
012
013
014
015
016
017

018 ABSTRACT 019 020 021 022 023 024

025 We present quantum algorithms that provide provable speedups for approximate
026 sampling from probability distributions of the form $\pi \propto e^{-f}$, where f is a potential
027 function that can be written as a finite sum, i.e., $f = \frac{1}{n} \sum_{i=1}^n f_i$. Our approach
028 focuses on stochastic gradient-based methods with only oracle access to individual
029 gradients $\{\nabla f_i\}_{i \in [n]}$. The techniques of our quantum algorithm are based
030 on a non-trivial integration of quantum mean estimation techniques and existing
031 variance reduction techniques such as SVRG and CV.
032

033 As these techniques often require occasional full-gradient calculations, the key
034 challenge is that an unbalanced weighting between variance reduction and quan-
035 tum mean estimation results in a regime where the quantum advantage is lost
036 due to frequent full-gradient computation. We overcome this difficulty by care-
037 fully optimizing the target variance level. Our algorithms improve the number of
038 gradient queries of classical samplers, such as Hamiltonian Monte Carlo (HMC)
039 and Langevin Monte Carlo (LMC), in terms of dimension, precision, and other
040 problem-dependent parameters.
041

042 1 INTRODUCTION 043 044

045 Efficient sampling from complex distributions is a fundamental problem in many scientific and
046 engineering disciplines, becoming increasingly important as modern applications deal with high-
047 dimensional data and complex probabilistic models. For example, in statistical mechanics, sampling
048 is used to analyze the thermodynamic properties of materials by exploring configurations of particle
049 systems Chandler (1987); Frenkel & Smit (2002). In convex geometry, it helps in approximating
050 volumes and studying high-dimensional structures Lovász & Vempala (2006); Cousins & Vempala
051 (2018). In probabilistic machine learning, sampling plays an important role in Bayesian inference,
052 as it facilitates posterior estimation and quantifies uncertainty in model predictions Welling & Teh
053 (2011); Wang et al. (2015); Durmus & Moulines (2018); Roy et al. (2021). Similarly, in non-convex
054 optimization, sampling allows for the exploration of complex energy landscapes and helps avoid
055 local minima, facilitating progress in tasks such as resource allocation, scheduling, and hyperpa-
056 rameter tuning in machine learning Zhang et al. (2017); Chen et al. (2020).
057

058 Given a potential function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we consider the problem of sampling from a probability
059 distribution π of the form

$$060 \pi(\mathbf{x}) = \frac{e^{-f(\mathbf{x})}}{\int e^{-f(\mathbf{x})} d\mathbf{x}}. \quad (1)$$

061 This distribution is called the Boltzmann-Gibbs distribution, and our goal is to efficiently sample
062 approximately from π while minimizing the number of gradient queries in the finite-sum setting,
063 i.e., $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$.
064

065 One widely-used method for sampling from the Gibbs distribution is to use Langevin Monte Carlo
066 (LMC) algorithm:
067

$$068 \mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) + \sqrt{2\eta_t} \epsilon_t, \quad (2)$$

069 where η_t is the step size and ϵ_t is isotropic Gaussian noise. Another method that is commonly
070 used in sampling is the Hamiltonian Monte Carlo (HMC) algorithm, which uses the principles of
071 Hamiltonian dynamics to propose new states in a Markov Chain. It introduces the Hamiltonian
072

054 $H(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \frac{1}{2}\|\mathbf{p}\|^2$ with auxiliary momentum variables and updates the position (\mathbf{x}) and
 055 momentum (\mathbf{p}) by simulating Hamiltonian dynamics, which follows the equations:
 056

$$\frac{d\mathbf{x}}{dt} = \frac{\partial H}{\partial \mathbf{p}}, \quad \frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \mathbf{x}}. \quad (3)$$

059 Similar to LMC, HMC is simulated in practice by discretizing Eq. (3). The algorithm also refreshes
 060 the momentum periodically from a random distribution, making the algorithm non-deterministic.
 061 Although effective, the computational cost of each iteration in these algorithms becomes prohibitive
 062 when the computation of the gradient is costly, such as in the finite-sum setting. To alleviate the
 063 computational burden, stochastic gradient-based samplers such as Stochastic Gradient Langevin
 064 Dynamics (SGLD) Welling & Teh (2011) and Stochastic Hamiltonian Monte Carlo (SG-HMC)
 065 Chen et al. (2014) have been proposed. Instead of computing the full gradient, these algorithms
 066 use stochastic approximation to the gradient. For example, the stochastic update for LMC becomes
 067

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{g}_t + \sqrt{2\eta_t} \epsilon_t. \quad (4)$$

069 In the finite-sum form, \mathbf{g}_t can be obtained by randomly sampling a component $i \in [n]$ and computing
 070 $\nabla f_i(\mathbf{x}_t)$. While stochastic gradient methods reduce computation at each iteration, they introduce
 071 variance into the gradient estimates, which can degrade the quality of the samples and slow down
 072 convergence.

073 Quantum computing offers a new way to address this bottleneck. By leveraging quantum primitives
 074 such as *quantum mean estimation*, it is possible to estimate averages of stochastic gradients with
 075 *quadratically fewer* oracle calls, thereby reducing the variance cost inherent in stochastic-gradient
 076 methods. Importantly, this quantum advantage can be realized within the same oracle framework
 077 used classically. Specifically, we consider quantum oracle access of the form

$$O_{\nabla f} |\mathbf{x}\rangle |i\rangle |0\rangle \mapsto |\mathbf{x}\rangle |i\rangle |\nabla f_i(\mathbf{x})\rangle. \quad (5)$$

079 This oracle can be implemented by making the classical gradient oracle reversible with additional
 080 registers, so its cost is comparable to the classical setting. While computing an exact gradient still
 081 requires $O(n)$ oracle calls, our algorithms demonstrate how quantum primitives can asymptotically
 082 reduce the number of gradient queries compared to the best classical methods.

083 With this framework, we design quantum sampling algorithms that parallel the structure of classical
 084 LMC and HMC. It is worth noting that a straightforward replacement of the stochastic gradient step
 085 in state-of-the-art classical algorithms with a quantum variance-reduction primitive does not yield
 086 asymptotic speedups, for the following reasons:

- 088 1. The state-of-the-art classical algorithms already implement variance reduction techniques
 089 such as SVRG, SARAH, SAGA, which require full gradient computations once in a while
 090 so that the variance of the stochastic gradients can be controlled more efficiently without
 091 computing the full gradient at every iteration. Therefore, directly replacing the classical
 092 stochastic gradient estimation with quantum algorithmic primitives is not meaningful un-
 093 less the full gradient computations can be done less frequently as well. In fact, using the
 094 existing classical algorithms' choice of parameters would in fact imply no speedup.
- 095 2. Quantum mean estimation algorithm requires the target variance as an input parameter,
 096 which depends on the variance of the random variable. In our setting, this corresponds to
 097 the variance of stochastic gradients at the current iteration. However, this variance depends
 098 on the trajectory of the iterates and is not bounded by a fixed constant. On the contrary, this
 099 restriction does not affect classical algorithms, as one can use a fixed batch size b and then
 100 analyze the convergence.
- 101 3. Even if the current variance of the stochastic gradients is given, it is not clear how to
 102 set the target variance level. Smaller target levels increase the cost of stochastic gradient
 103 estimation, whereas larger target levels do not give any quantum speedups. This is because
 104 in this regime, one does not effectively exploit the subroutine: the algorithm essentially
 105 becomes regular LMC or HMC with no variance reduction.

106 To address these challenges, we develop a variance upper bound that quantifies stochastic gradient
 107 variance along the LMC trajectory. This enables us to optimize the trade-off between full-gradient
 108 computation and quantum mean estimation. Then, we optimize the target variance level in such a

108

109
110
111
Table 1: Summary of the results (some of the previous results use a different scaling of f and we
convert the results to the same scaling as ours in the table). Here, we mainly focus on n and ϵ
dependency. See Theorems 4.3, 4.4 and 4.9 for explicit dependencies on L, μ, α, d .

112

Algorithm	Assumptions	Metric	Gradient Complexity
SG-HMC Zou & Gu (2021)	Strongly Convex	W_2	$\tilde{\mathcal{O}}(n\epsilon^{-2})$
SVRG-HMC Zou & Gu (2021)	Strongly Convex	W_2	$\tilde{\mathcal{O}}(n^{2/3}\epsilon^{-2/3} + \epsilon^{-1})$
SAGA-HMC Zou & Gu (2021)	Strongly Convex	W_2	$\tilde{\mathcal{O}}(n^{2/3}\epsilon^{-2/3} + \epsilon^{-1})$
CV-HMC Zou & Gu (2021)	Strongly Convex	W_2	$\tilde{\mathcal{O}}(\epsilon^{-2})$
SRVR-HMC Zou et al. (2019)	Dissipative Gradients	W_2	$\tilde{\mathcal{O}}(n + n^{1/2}\epsilon^{-2} + \epsilon^{-4})$
SVRG-LMC Kinoshita & Suzuki (2022)	LSI	KL	$\tilde{\mathcal{O}}(n + n^{1/2}\epsilon^{-1})$
SARAH-LMC Kinoshita & Suzuki (2022)	LSI	KL	$\tilde{\mathcal{O}}(n + n^{1/2}\epsilon^{-1})$
QSVRG-HMC [Theorem 4.3]	Strongly Convex	W_2	$\tilde{\mathcal{O}}(n^{1/2}\epsilon^{-3/4} + \epsilon^{-1})$
QCV-HMC [Theorem 4.4]	Strongly Convex	W_2	$\tilde{\mathcal{O}}(\epsilon^{-3/2})$
QSVRG-LMC [Theorem 4.9]	LSI	KL ¹	$\tilde{\mathcal{O}}(n + n^{1/3}\epsilon^{-1})$

123

124

125
way that both the cost of mean estimation and required full gradient computations decrease at the
126 same time. By setting these optimum parameters, we prove that our algorithms achieve asymptotically
127 fewer oracle calls than the best-known classical methods for both strongly convex and
128 nonconvex potentials (see Table 1). We focus on SVRG- and CV-based samplers because they seem
129 to be the most efficient samplers for finite-sum functions in the classical setting; although extensions
130 to other gradient based samplers are possible.

131

132
We note that our algorithms rely on fault-tolerant quantum computers to implement quantum mean
133 estimation; since such devices are not yet available, we cannot provide empirical validation at this
134 stage. Nevertheless, establishing the theoretical foundations is essential for clarifying the scope of
possible quantum speedups in sampling.

135

136

2 RELATED WORK

137

138
Non-asymptotic convergence rates for SGLD and SG-HMC have been analyzed extensively by Ra-
139 ginsky et al. (2017); Xu et al. (2018); Zou et al. (2021); Das et al. (2023) and Chen et al. (2014);
140 Zou & Gu (2021) respectively. In the finite sum setting, more sophisticated variance reduction tech-
141 niques such as SVRG Johnson & Zhang (2013), SAGA Defazio et al. (2014), SARAH Nguyen et al.
142 (2017), and Control Variates (CV) Baker et al. (2019) have been used to reduce the variance of
143 stochastic gradients by leveraging the gradient information from previous iterations. Although these
144 methods were originally introduced in the context of optimization, successive works have applied
145 these methods to improve sampling efficiency via LMC Dubey et al. (2016); Chatterji et al. (2018);
146 Baker et al. (2019); Kinoshita & Suzuki (2022) and HMC Zou et al. (2019); Zou & Gu (2021). In
147 particular, Zou & Gu (2021) has incorporated various variance reduction techniques to SG-HMC
148 and analyzed convergence in Wasserstein distance for smooth and strongly convex potentials. In the
149 non-log-concave setting, Kinoshita & Suzuki (2022) has analyzed the convergence of SVRG-LMC
150 and SARAH-LMC for target distributions that satisfy the Log-Sobolev inequality and applied their
results to optimize structured non-convex objectives.

151

152
In the context of quantum sampling, the most of the existing results make use of *quantum walks*,
153 which has been shown to provide speedups for certain Markov Chain Monte Carlo (MCMC) meth-
154 ods by improving the mixing time of the underlying Markov chain Szegedy (2004); Somma et al.
155 (2007; 2008); Wocjan & Abeyesinghe (2008); Chakrabarti et al. (2023). These methods have
156 been incorporated into various domains to improve the computation time of various tasks Magniez
157 et al. (2007); Apers & Sarlette (2019); Childs et al. (2022); Chakrabarti et al. (2023); Li & Zhang
158 (2024); Chakrabarti et al. (2024). For sampling from continuous distributions, Kinoshita & Suzuki
159 (2022) used quantum walk framework to improve the gradient queries to approximately sample from
160 strongly-convex distributions. However, a key limitation of quantum walks is that they require the
Markov chain to satisfy detailed balance condition. A Markov chain on Ω with transition density

161

¹Convergence in KL divergence implies convergence in squared TV and W_2 distances due to Pinsker's and Talagrand's inequalities.

matrix P and stationary density π needs to satisfy for all $\mathbf{x}, \mathbf{y} \in \Omega$, $\pi(\mathbf{x})P(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{y})P(\mathbf{y}, \mathbf{x})$. Unfortunately, many commonly used sampling algorithms, such as LMC and HMC, are not reversible due to the finite discretization steps involved in their implementation. The reversibilization techniques such as Metropolis reversibilization require evaluating the function exactly to compute ratio $\pi(\mathbf{x})/\pi(\mathbf{y})$ which is expensive in finite-sum case. Additionally, implementing rejection steps causes additional overheads: if a proposed move is rejected, one must revert to the previous state. However, due to the no-cloning theorem, it is not straightforward to restore the previous quantum state. Therefore, the quantum walk operator needs nontrivial techniques when it involves Metropolis correction. We refer the reviewer to recent work on this topic Claudon et al. (2025b). Moreover, even when the Markov chain is reversible, stochastic gradients introduce randomness that disrupts the coherent evolution of the quantum walk, which is a critical component of its speedup Ozgul et al. (2024). More recently, Claudon et al. (2025a) proposed a similar technique to obtain quantum speedups for nonreversible Markov chains, using the idea of geometric reversibilization with respect to the so-called “most reversible” distribution which requires $\mathcal{O}(n)$ gradient computations in this case. Another limitation of quantum walks is that they typically offer convergence guarantees in terms of total variation distance; however, many practical sampling tasks are more concerned with metrics like Wasserstein distance or Kullback-Leibler divergence.

However, hybrid algorithms that exploit quantum computing methods as subroutines are easier to implement and do not suffer from these issues. In the context of optimization, quantum algorithms such as multi-dimensional quantum mean estimation Cornelissen et al. (2022) and quantum gradient estimation Jordan (2005); Gilyén et al. (2019) have shown promise in reducing the computational cost associated with gradient-based methods van Apeldoorn et al. (2020); Chakrabarti et al. (2020); Sidford & Zhang (2023); Zhang et al. (2024); Liu et al. (2024). These techniques are particularly well-suited for addressing challenges in large-scale and noisy settings, as they can provide more accurate gradient estimates with fewer queries. However, these methods have not been considered for sampling tasks and this paper focuses on integrating these quantum techniques to enhance the efficiency of stochastic gradient-based samplers and alleviate the computational burden inherent in classical methods.

2.1 PRELIMINARIES

Notation: Bold symbols, such as \mathbf{x} and \mathbf{y} , are used to represent vectors, with $\|\cdot\|$ indicating the Euclidean or operator norm depending on the context. Given two scalars a and b , we use $a \wedge b$ to denote $\min\{a, b\}$ and use $a \vee b$ to denote $\max\{a, b\}$. The notation \tilde{O} is used to suppress the polylogarithmic dependencies on d, ϵ, L, μ and α that will be defined later in the text.

Quantum computation: Quantum computation is naturally expressed in the language of linear algebra. The *computational basis* of \mathbb{C}^d is given by $\mathbf{e}_0, \dots, \mathbf{e}_{d-1}$, where $\mathbf{e}_i = (0, \dots, 1, \dots, 0)^\top$ has a 1 in the $(i+1)^{\text{st}}$ position. In *Dirac notation*, we write $|i\rangle$ (a “ket”) for \mathbf{e}_i and $\langle i|$ (a “bra”) for \mathbf{e}_i^\top .

The *tensor product* of two quantum states is their Kronecker product. If $|u\rangle \in \mathbb{C}^{d_1}$ and $|v\rangle \in \mathbb{C}^{d_2}$, then

$$|u\rangle \otimes |v\rangle = (u_0 v_0, u_0 v_1, \dots, u_{d_1-1} v_{d_2-1})^\top \in \mathbb{C}^{d_1} \otimes \mathbb{C}^{d_2}. \quad (6)$$

The fundamental unit of quantum information is the *qubit*, a state in \mathbb{C}^2 of the form $a|0\rangle + b|1\rangle$ with $a, b \in \mathbb{C}$ and $|a|^2 + |b|^2 = 1$. An n -qubit product state takes the form $|v_1\rangle \otimes \dots \otimes |v_n\rangle \in \mathbb{C}^{2^n}$, where each $|v_i\rangle$ is a single-qubit state. Most vectors in \mathbb{C}^{2^n} , however, cannot be written as product states. For brevity, we often write $|u\rangle |v\rangle$ instead of $|u\rangle \otimes |v\rangle$.

Quantum states evolve under *unitary transformations*. In the circuit model, a *k -qubit gate* is a unitary operator in \mathbb{C}^{2^k} . Two-qubit gates are *universal*: any n -qubit unitary can be decomposed into gates acting trivially on $n-2$ qubits and nontrivially on two qubits. The *gate complexity* of an operation is defined as the number of two-qubit gates required in its circuit implementation.

Access to information of a function or a probability distribution in quantum algorithms is provided via a *quantum oracle*. Such oracles must be reversible and allow queries on superpositions of inputs. The following definition demonstrates an oracle we use in this paper for sampling from a probability distribution.

216 **Definition 2.1** (Quantum Sampling Oracle). Quantum sampling oracle O_X of a random variable
 217 $X \in \Omega$ is given by $O_X |0\rangle |0\rangle \mapsto \sum_{X \in \Omega} \sqrt{\Pr(X)} |X\rangle |\text{garbage}(X)\rangle$.
 218

219 Here, the second register contains $|\text{garbage}(X)\rangle$, which depends on X . The state in the (auxiliary)
 220 garbage register is usually generated in some intermediate step of computing X in the first register. It
 221 is important to note that the state in this quantum sampling oracle differs from the coherent quantum
 222 sample state, as the former is entangled and we cannot simply discard the garbage register.

223 Another oracle used in this paper is the *stochastic gradient oracle* as specified in Eq. (5).
 224

225 Beyond simulating classical random sampling, quantum oracles enable uniquely quantum effects
 226 such as interference. These underlie key techniques like amplitude amplification (central to Grover’s
 227 search Grover (1996)) and amplitude estimation, both of which rely on coherent oracle access.
 228 Similar considerations apply to the quantum gradient oracle Eq. (5). Whenever a classical oracle
 229 can be realized by a circuit, the corresponding quantum oracle can be implemented by a quantum
 230 circuit with little overhead. Thus, quantum oracles provide a natural framework for analyzing the
 231 complexity of tasks such as sampling from a probability distribution and optimization.

232 To sample from a distribution p over \mathbb{R}^d , it suffices to prepare the quantum state $\sum_{\mathbf{x}} \sqrt{p(\mathbf{x})} d|\mathbf{x}\rangle$
 233 and then measure it.

234 **Metrics:** We use several metrics to compare probability distributions over a state space \mathcal{X} . Let
 235 π and μ be two probability distributions on \mathcal{X} . The p -Wasserstein distance between π and μ
 236 is defined as $W_p(\pi, \mu) = (\inf_{\gamma \in \Gamma(\pi, \mu)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \|\mathbf{x} - \mathbf{y}\|^p)^{1/p}$ where $\Gamma(\pi, \mu)$ is the set of all
 237 joint distributions $\gamma(\mathbf{x}, \mathbf{y})$ whose marginals are π and μ . The KL divergence of π with respect
 238 to μ is defined as $\text{KL}(\pi \parallel \mu) = \int_{\mathcal{X}} d\mathbf{x} \pi(\mathbf{x}) \log \left(\frac{\pi(\mathbf{x})}{\mu(\mathbf{x})} \right)$ and the relative Fisher information is
 239 $\text{FI}(\pi \parallel \mu) = \int_{\mathcal{X}} d\mathbf{x} \pi(\mathbf{x}) \left\| \nabla \log \left(\frac{\pi(\mathbf{x})}{\mu(\mathbf{x})} \right) \right\|^2$. The total variation distance is defined as $\text{TV}(\pi, \mu) =$
 240 $\sup_{A \subseteq \mathcal{X}} |\pi(A) - \mu(A)| = \frac{1}{2} \int_{\mathcal{X}} d\mathbf{x} |\pi(\mathbf{x}) - \mu(\mathbf{x})|$.

241 In the next section, we analyze the trade-off between the error due to stochastic gradients and
 242 discretization to quantify how much quantum mean estimation techniques can provide speedups when
 243 combined with classical variance reduction methods such as SVRG and CV.

247 3 BACKGROUND

249 In this section, we give background on some classical and quantum algorithms for various tasks that
 250 are repeatedly referred in the main text.

252 3.1 OVERVIEW OF CLASSICAL SAMPLING ALGORITHMS

254 One widely-used method for sampling from the Gibbs distribution is through the Langevin diffusion
 255 equation, which follows the solution to the following stochastic differential equation (SDE):

$$256 \quad d\mathbf{x}_t = -\nabla f(\mathbf{x}_t) dt + \sqrt{2} dB_t, \quad (7)$$

257 where \mathbf{B}_t is the standard Brownian motion. The Euler-Maruyama discretization of this SDE results
 258 in the well-known Langevin Monte Carlo (LMC) algorithm:

$$260 \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) + \sqrt{2\eta_t} \epsilon_t, \quad (8)$$

261 In stochastic setting, once replace $\nabla f(\mathbf{x}_t)$ by stochastic gradients $\mathbf{g}(\mathbf{x}_k, \xi_k)$.

263 Hamiltonian Monte Carlo (HMC) is an advanced sampling technique designed to efficiently explore
 264 high-dimensional probability distributions by introducing auxiliary momentum variables. Given a
 265 target distribution $\pi(\mathbf{x}) \propto e^{-f(\mathbf{x})}$, HMC augments the state space with momentum variables \mathbf{p} and
 266 defines the Hamiltonian $H(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \frac{1}{2} \|\mathbf{p}\|^2$ where $\mathbf{p} \sim \mathcal{N}(0, I)$.

267 HMC alternates between updating the position \mathbf{x} and momentum \mathbf{p} by simulating Hamiltonian
 268 dynamics Eq. (3). In practice, Hamiltonian dynamics is simulated using the leapfrog integrator,
 269 which discretizes the continuous equations of motion. The key advantage of HMC is that it allows for large,
 270 efficient moves through the parameter space by leveraging gradient information

270 **Algorithm 1** SG-LMC
 271 **input** The stochastic gradient oracle $O_{\nabla f}$, initial point \mathbf{x}_0 , step size η , number of steps K
 272 **output** Approximate sample from $\pi \propto e^{-f(\mathbf{x})}$
 273 **for** $t = 0$ to K **do**
 274 Sample $\epsilon_t \sim \mathcal{N}(0, I)$
 275 $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_t \mathbf{g}(\mathbf{x}_k, \xi_k) + \sqrt{2\eta_k} \epsilon_k$,
 276 **end for**
 277 **Return** \mathbf{x}^K

280 and auxiliary momentum. This reduces the correlation between successive samples, particularly
 281 in high-dimensional spaces, resulting in faster convergence compared to simple random-walk meth-
 282 ods like the Metropolis-Hastings algorithm. In practice, Hamiltonian dynamics are simulated using
 283 the leapfrog integrator, which discretizes the continuous equations of motion.

284 After a series of updates, the momentum \mathbf{p}_{k+1} is refreshed by sampling from $\mathcal{N}(0, I)$. This dis-
 285 cretization ensures symplecticity, preserving volume in phase space and allowing the algorithm to
 286 make large, energy-conserving moves through the parameter space.

288 **Algorithm 2** SG-HMC
 289 **input** The stochastic gradient oracle $O_{\nabla f}$, initial point \mathbf{x}_0 , step size η , number of leapfrog steps S ,
 290 number of HMC proposals T
 291 **output** Approximate sample from $\pi \propto e^{-f(\mathbf{x})}$
 292 **for** $t = 0$ to T **do**
 293 Sample $\mathbf{p}_{St} \sim \mathcal{N}(0, I)$
 294 **for** $s = 0$ to $S - 1$ **do**
 295 $k = St + s$
 296 $\mathbf{x}_{k+1} = \mathbf{x}_k + \eta \mathbf{p}_k - \frac{\eta^2}{2} \mathbf{g}(\mathbf{x}_k, \xi_k)$
 297 $\mathbf{p}_{k+1} = \mathbf{p}_k - \frac{\eta}{2} \mathbf{g}(\mathbf{x}_k, \xi_k) - \frac{\eta}{2} \mathbf{g}(\mathbf{x}_{k+1}, \xi_{k+1/2})$
 298 **end for**
 299 **end for**
 300 **Return** \mathbf{x}^T

302 Similar to SGLD, one can replace the gradients with stochastic gradients resulting in SG-HMC (See
 303 Algorithm 2). The stochastic gradients $\mathbf{g}(\mathbf{x}, \xi)$ in Algorithm 2 can be obtained using different tech-
 304 niques such as mini-batch, SVRG, CV. In this case, we use quantum variance reduction techniques
 305 to compute $\mathbf{g}(\mathbf{x}, \xi)$.

3.2 QUANTUM MEAN ESTIMATION

310 Quantum mean estimation is a technique to estimate the mean of a d -dimensional random variable
 311 X up to ϵ accuracy using $\tilde{\mathcal{O}}(d^{1/2}/\epsilon)$ queries, which is a quadratic improvement in ϵ compared to
 312 classical algorithms Cornelissen et al. (2022). Although the quantum mean estimation algorithm
 313 is biased, Sidford & Zhang (2023) developed an unbiased quantum mean estimation algorithm.
 314 Specifically, for a multi-dimensional variable with mean μ and variance σ^2 , unbiased quantum mean
 315 estimation outputs an estimate $\hat{\mu}$ such that $\mathbb{E}[\hat{\mu}] = \mu$ and $\mathbb{E}[\|\hat{\mu} - \mu\|^2] \leq \hat{\sigma}^2$ using $\tilde{\mathcal{O}}(d^{1/2}\sigma/\hat{\sigma})$
 316 queries.

317 The following lemma shows that the mean $\mathbb{E}[X]$ for a random variable X can be computed quadra-
 318 tically faster than classical mean estimation with respect to oracle O_X .

319 **Lemma 3.1** (Unbiased Quantum Mean Estimation Sidford & Zhang (2023)). *For a d -dimensional
 320 random variable X with $\text{Var}[X] \leq \sigma^2$ and some $\hat{\sigma} \geq 0$, suppose we are given ac-
 321 cess to its quantum sampling oracle O_X (as in Definition 2.1). Then, there is a procedure
 322 $\text{QuantumMeanEstimation}(O_X, \hat{\sigma})$ that uses $\tilde{\mathcal{O}}\left(\frac{d^{1/2}\sigma}{\hat{\sigma}}\right)$ queries to O_X and outputs an unbi-
 323 ased estimate $\hat{\mu}$ of the expectation μ satisfying $\text{Var}[\hat{\mu}] \leq \hat{\sigma}^2$.*

324 **4 QUANTUM SPEEDUPS FOR FINITE-SUM SAMPLING VIA GRADIENT**
325 **ORACLE**
326

327 We assume access to the oracle defined in Eq. (5). The goal is to approximately sample from π by
328 using as few gradient computations as possible without deteriorating the convergence. To this end,
329 we introduce the first stochastic-gradient samplers that integrate unbiased quantum mean estimation
330 with classical variance-reduction frameworks, yielding provable improvements in gradient-query
331 complexity for both HMC and LMC. Our algorithms are actually quite simple: we replace the
332 stochastic gradient estimation with quantum mean estimation. However, obtaining the speedup is
333 not as simple because to lower the total computational cost, we must decrease the cost of full gradient
334 computations as well. Even though we do not modify the full gradient estimation part, the quantum
335 mean estimation allows less frequent full gradient computations due to improved variance reduction.

336 Our analysis develops a variance control results (starting with A.3) that quantifies stochastic
337 gradient variance along the sampling trajectory, enabling us to choose optimal balance between full-
338 gradient computations and quantum mean estimation. Building on this tool, we establish improved
339 complexity bounds under both strong convexity and LSI assumptions, demonstrating speedups over
340 state-of-the-art classical algorithms. For readability, we include some of the key results in the main
341 text and defer full technical details to the appendix due to the space limitation.

342 **4.1 SAMPLING UNDER STRONG CONVEXITY VIA HAMILTONIAN MONTE CARLO**
343

344 First, we consider quantum speedups for Hamiltonian Monte Carlo (HMC) algorithm using quantum
345 variance reduction techniques.

347 **Algorithm 3** QSVRG/QCV

348 **input** $O_{\nabla f}$, current iterate \mathbf{x}_k , smoothness constant L , variance scale factor b , epoch length m .
349 **output** Quantum variance reduced stochastic gradient \mathbf{g} .

351 1: **QSVRG:**
352 2: **if** $k \bmod m = 0$ **then**
353 3: $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$
354 4: $\tilde{\mathbf{x}} = \mathbf{x}_k$
355 5: **else**
356 6: Define oracle $O_{\text{SVRG}}^{\mathbf{x}_k}$:

357
$$|0\rangle |0\rangle \mapsto \frac{1}{\sqrt{n}} \sum_{i=1}^n |\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}})\rangle |i\rangle$$

361 7: $\hat{\sigma}^2 = L^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 / b^2$
362 8: $\mathbf{g}_k = \text{QuantumMeanEstimation}(O_{\text{SVRG}}^{\mathbf{x}_k}, \hat{\sigma}^2)$
363 9: **end if**

364 10: **QCV:**
365 11: Define oracle $O_{\text{CV}}^{\mathbf{x}_k}$:

367
$$|0\rangle |0\rangle \mapsto \frac{1}{\sqrt{n}} \sum_{i=1}^n |\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)\rangle |i\rangle$$

370 12: $\hat{\sigma}^2 = L^2 \|\mathbf{x}_k - \mathbf{x}_0\|^2 / b^2$
371 13: $\mathbf{g}_k = \text{QuantumMeanEstimation}(O_{\text{CV}}^{\mathbf{x}_k}, \hat{\sigma}^2)$

372 14: **Return** \mathbf{g}_k

374
375 We propose to replace the gradients in HMC (See Algorithm 2 in appendix) with quantum gradients
376 computed via Algorithm 3. Essentially Algorithm 3 combines the classical variance reduction tech-
377 niques with the unbiased quantum mean estimation algorithm in Lemma 3.1 to reduce the variance
further. The epoch length m for QSVRG determines the period where the full gradient needs to be

378 computed. The parameter b is the quantum analog of batch size and will be determined analytically. To establish the convergence of the new samplers, we make the following assumptions in this
 379 section.
 380

381 **Assumption 4.1** (Strong Convexity). There exists a positive constant μ such that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$
 382 it holds that

$$383 \quad f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2. \quad (9)$$

385 **Assumption 4.2** (Lipschitz Stochastic Gradients). There exists a positive constant L such that for
 386 all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and all functions $f_i, i = 1, \dots, n$, it holds that

$$387 \quad \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|. \quad (10)$$

389 We also define the condition number $\kappa = \frac{L}{\mu}$. These assumptions are standard and used in the
 390 classical analysis of HMC Zou & Gu (2021). Next, we give the main theorem for the quantum
 391 Hamiltonian Monte Carlo algorithm implemented with QSVRG technique.
 392

393 **Theorem 4.3** (Main Theorem for QSVRG–HMC). *Let μ_k be the distribution of \mathbf{x}_k in QSVRG–HMC
 394 algorithm. Suppose that f satisfies Assumptions 4.1 and 4.2. Given that the initial point \mathbf{x}_0 satisfies
 395 $\|\mathbf{x}_0 - \arg \min_{\mathbf{x}} f(\mathbf{x})\| \leq \frac{d}{\mu}$, then, for $\eta = \mathcal{O}(\frac{\epsilon}{L^{1/2} d^{1/2} \kappa^{3/2}})$, $S = \tilde{\mathcal{O}}\left(\frac{L d^{1/2} \kappa^{3/2}}{\epsilon}\right)$, $T = \tilde{\mathcal{O}}(1)$,
 396 $b = \mathcal{O}\left(\frac{L^{1/8} \epsilon^{1/4} n^{1/2}}{d^{1/8} \kappa^{3/8}} \vee 1\right)$, and $m = n/b$, we have*

$$397 \quad W_2(\mu_{ST}, \pi) \leq \epsilon.$$

400 The total query complexity to the stochastic gradient oracle is $\tilde{\mathcal{O}}\left(\frac{L d^{1/2} \kappa^{3/2}}{\epsilon} + \frac{L^{9/8} d^{7/8} \kappa^{3/4} n^{1/2}}{\epsilon^{3/4}}\right)$.
 401

402 The following theorem is for quantum Hamiltonian Monte Carlo algorithm implemented with QCV
 403 technique.

404 **Theorem 4.4** (Main Theorem for QCV–HMC). *Let μ_k be the distribution of \mathbf{x}_k in QCV–HMC algo-
 405 rithm. Suppose that f satisfies Assumptions 4.1 and 4.2. Given that the initial point \mathbf{x}_0 satisfies
 406 $\|\mathbf{x}_0 - \arg \min_{\mathbf{x}} f(\mathbf{x})\| \leq \frac{d}{\mu}$, then, for $\eta = \mathcal{O}(\frac{\epsilon}{L^{1/2} d^{1/2} \kappa^{3/2}})$, $S = \tilde{\mathcal{O}}\left(\frac{L d^{1/2} \kappa^{3/2}}{\epsilon}\right)$, $T = \tilde{\mathcal{O}}(1)$, and
 407 $b = \mathcal{O}\left(\frac{d^{1/4} \kappa^{3/4}}{L^{1/4} \epsilon^{1/2}} \vee 1\right)$, we have*

$$408 \quad W_2(\mu_{ST}, \pi) \leq \epsilon.$$

409 The total query complexity to the stochastic gradient oracle is $\tilde{\mathcal{O}}\left(\frac{L d^{5/4} \kappa^{9/4}}{\epsilon^{3/2}}\right)$.
 410

411 We postpone the proofs of Theorems 4.3 and 4.4 to Appendix A. The main idea in these proofs is
 412 to express the variance of stochastic gradients σ throughout the trajectory of the HMC in terms of b
 413 and the distance between current iterate and the last iterate the full gradient is computed. Then, we
 414 optimize b so that per cost of quantum mean estimation $\mathcal{O}(\frac{\sigma}{\delta})$ is equal to the per iteration cost of full
 415 gradient computation $\tilde{\mathcal{O}}(n/m)$ (because full gradient is only computed once in every m iteration)
 416 to exploit the quantum mean estimation without full gradient estimation dominating the cost.
 417

418 Theorems 4.3 and 4.4 imply that when $n = \mathcal{O}(\epsilon^{-1/2})$ the best classical (SVRG–HMC) and the best
 419 quantum (QSVRG–HMC) algorithms have $\tilde{\mathcal{O}}(\epsilon^{-1})$ gradient complexity. On the other hand, when
 420 $n = \omega(\epsilon^{-1})$, quantum algorithms have better complexity than the best classical algorithms, where
 421 the race between QSVRG–HMC and QCV–HMC depends on how large n is.

422 *Remark 4.5.* Both the classical algorithms in Zou & Gu (2021) and quantum algorithms in this paper
 423 assume that the starting point is (d/μ) -close to the minimizer $\mathbf{x}^* = \arg \min f(\mathbf{x})$. In case this point
 424 is not given, it can be obtained using $\mathcal{O}(n)$ iterations of SGD Baker et al. (2019).

427 4.2 SAMPLING UNDER LOG-SOBOLEV INEQUALITY VIA LANGEVIN MONTE CARLO

428 We use SVRG-LMC for the base algorithm in Kinoshita & Suzuki (2022) and replace the stochastic
 429 gradient calculation with unbiased quantum mean estimation. This section generalizes the strong
 430 convexity assumption with the following LSI assumption, which is common in non-logconcave
 431 sampling.

432 **Assumption 4.6** (Log-Sobolev Inequality). We say that π satisfies the Log-Sobolev inequality with
 433 constant α if for all ρ , it holds that
 434

$$435 \quad \text{KL}(\rho||\pi) \leq \frac{1}{2\alpha} \text{FI}(\rho||\pi). \quad (11)$$

436

437 This is a sampling analog of the PL (Polyak-Łojasiewicz) condition commonly used in optimization
 438 Chewi & Stromme (2024) and standard in non-log-concave sampling literature Vempala &
 439 Wibisono (2019); Ma et al. (2019); Chewi et al. (2022); Kinoshita & Suzuki (2022). We note that
 440 LSI relaxes strong convexity in the sense that for any μ strongly convex function f , π satisfies the
 441 Log-Sobolev inequality with constant $\frac{\mu}{2}$. We also note that this assumption is weaker than the dis-
 442 ssipative gradient condition Raginsky et al. (2017); Zou et al. (2019) which is used commonly in
 443 non-log-concave sampling. We highlight the key steps in the proof idea here; full technical details
 444 appear in Appendix B. First we start by bounding the variance of the stochastic gradients along the
 445 trajectory of LMC in terms of KL divergence to the target Gibbs distribution and b .
 446

447 **Lemma 4.7** (QSVRG-LMC Variance Lemma). *Let $k' < k$ be the last iteration where the full gradient
 448 is computed in QSVRG-LMC and $\sigma_k^2 = \mathbb{E}\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2$. Then, for $\eta^2 \leq \frac{1}{6L^2m^2}$,*

$$449 \quad \sigma_{k'+l}^2 \leq \frac{16L^4\eta^2}{\alpha} \sum_{r=1}^l \text{KL}(\mu_{k'+r-1}||\pi) + \frac{8\eta dmL^2}{b^2}. \quad (12)$$

450

451 Then we prove the following theorem for LMC with stochastic gradients, which might be of inde-
 452 pendent interest that will be proved in appendix Appendix B. Although similar convergence results
 453 exist in literature, they only seem to apply for langevin algorithm with full gradients. The proof of
 454 this theorem uses a comparison between the true SDE and and approximate SDE with drift term set
 455 to stochastic gradients and we analyze this distance in terms of KL divergence between the distribu-
 456 tions using the variance bound above.
 457

458 **Theorem 4.8** (Convergence theorem for QSVRG-LMC). *Assume that $m \leq b^2$. Then, for $\eta \leq \frac{\alpha^2}{24L^2m}$,
 459 the iterates in QSVRG-LMC satisfy,*

$$460 \quad \text{KL}(\mu_k||\pi) \leq e^{-\alpha\eta k} \text{KL}(\mu_0||\pi) + \frac{64m\eta dL^2}{\alpha b^2} + \frac{24\eta dL^2}{\alpha}. \quad (13)$$

461

462 Note that the first term corresponds to convergence of continous SDE, the second term is due to
 463 stochastic gradients and the last term is due to discretization of SDE. Using the convergence theorem,
 464 we obtain the main result for LSI.

465 **Theorem 4.9** (Main Theorem for QSVRG-LMC). *Let μ_k be the distribution of \mathbf{x}_k in QSVRG-LMC
 466 algorithm. Suppose that f satisfies Assumptions 4.2 and 4.6. Then for $\eta = \mathcal{O}\left(\frac{\epsilon\alpha}{dL^2} \wedge \frac{\alpha}{L^2m}\right)$, $K =$
 467 $\tilde{\mathcal{O}}\left(\frac{L^2 \log(\text{KL}(\mu_0||\pi))}{\alpha^2} \left(n^{2/3} + \frac{d}{\epsilon}\right)\right)$, $b = \tilde{\mathcal{O}}(n^{1/3})$, and $m = \tilde{\mathcal{O}}(n^{2/3})$ we have*

$$468 \quad \left\{ \text{KL}(\mu_K||\pi), \text{TV}(\mu_K, \pi)^2, \frac{\alpha}{2} \mathbf{W}_2(\mu_K, \pi)^2 \right\} \leq \epsilon.$$

469 The total query complexity to the stochastic gradient oracle is
 470 $\tilde{\mathcal{O}}\left(\frac{L^2 \log(\text{KL}(\mu_0||\pi))}{\alpha^2} \left(nd^{1/2} + \frac{d^{3/2}n^{1/3}}{\epsilon}\right)\right)$.
 471

472 The proof of Theorem 4.9 is postponed to Appendix B. The term $n^{1/3}\epsilon^{-1}$ is only possible because
 473 the cost full gradient estimation is amortized in the sense that the total cost of stochastic gradient
 474 computations KbT which is equal to the cost of total full gradient estimation nTK/m . Hence, both
 475 costs go down thanks to quantum variance reduced gradients. We note that in classical SVRG-LMC
 476 the optimum parameters $m = b = n^{1/2}$ where b corresponds to inner batch size.
 477

478 Our algorithm improves the dominant term in gradient complexity from $\tilde{\mathcal{O}}(n^{1/2}\epsilon^{-1})$ to $\tilde{\mathcal{O}}(n^{1/3}\epsilon^{-1})$.
 479 It is also worth mentioning that recently Huang et al. (2024) proposed a proximal sampling algo-
 480 rithm that uses $\tilde{\mathcal{O}}(\sigma^2\epsilon^{-1})$ gradient queries in the LSI setting when the stochastic gradients have
 481 bounded variance σ^2 . However, this assumption is different from our setting since the variance in
 482 the stochastic gradients is not uniformly bounded by a constant, but it is bounded throughout the
 483 trajectory by a function of problem parameters such as d, b, m, L, α (See Lemma 4.7).
 484

486 REFERENCES
487

488 Simon Apers and Alain Sarlette. Quantum fast-forwarding: Markov chains and graph property
489 testing. *Quantum Info. Comput.*, 19(3–4):181–213, March 2019. ISSN 1533-7146.

490 Jack Baker, Paul Fearnhead, Emily B. Fox, and Christopher Nemeth. Control variates for stochastic
491 gradient mcmc. *Statistics and Computing*, 29(3):599–615, May 2019. ISSN 0960-3174. doi: 10.
492 1007/s11222-018-9826-2. URL <https://doi.org/10.1007/s11222-018-9826-2>.

493 Shouvanik Chakrabarti, Andrew M. Childs, Tongyang Li, and Xiaodi Wu. Quantum algorithms and
494 lower bounds for convex optimization. *Quantum*, 4:221, January 2020. ISSN 2521-327X. doi: 10.
495 22331/q-2020-01-13-221. URL <https://doi.org/10.22331/q-2020-01-13-221>.

496 Shouvanik Chakrabarti, Andrew M. Childs, Shih-Han Hung, Tongyang Li, Chunhao Wang, and
497 Xiaodi Wu. Quantum algorithm for estimating volumes of convex bodies. *ACM Transactions on
498 Quantum Computing*, 4(3), May 2023. doi: 10.1145/3588579. URL <https://doi.org/10.1145/3588579>.

499 Shouvanik Chakrabarti, Dylan Herman, Guneykan Ozgul, Shuchen Zhu, Brandon Augustino, Tianyi
500 Hao, Zichang He, Ruslan Shaydulin, and Marco Pistoia. Generalized short path algorithms:
501 Towards super-quadratic speedup over markov chain search for combinatorial optimization, 2024.
502 URL <https://arxiv.org/abs/2410.23270>.

503 David Chandler. Introduction to modern statistical. *Mechanics*. Oxford University Press, Oxford,
504 UK, 5(449):11, 1987.

505 Niladri Chatterji, Nicolas Flammarion, Yian Ma, Peter Bartlett, and Michael Jordan. On the theory
506 of variance reduction for stochastic gradient Monte Carlo. In Jennifer Dy and Andreas Krause
507 (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of
508 *Proceedings of Machine Learning Research*, pp. 764–773. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/chatterji18a.html>.

509 Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In
510 Eric P. Xing and Tony Jebara (eds.), *Proceedings of the 31st International Conference on Ma-
511 chine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 1683–1691, Be-
512 jing, China, 22–24 Jun 2014. PMLR. URL <https://proceedings.mlr.press/v32/cheni14.html>.

513 Xi Chen, Simon S. Du, and Xin T. Tong. On stationary-point hitting time and ergodicity of stochastic
514 gradient langevin dynamics. *Journal of Machine Learning Research*, 21(68):1–41, 2020. URL
515 <http://jmlr.org/papers/v21/19-327.html>.

516 Sinho Chewi and Austin J. Stromme. The ballistic limit of the log-sobolev constant equals the
517 polyak-łojasiewicz constant, 2024. URL <https://arxiv.org/abs/2411.11415>.

518 Sinho Chewi, Murat A Erdogdu, Mufan Li, Ruqi Shen, and Shunshi Zhang. Analysis of langevin
519 monte carlo from poincare to log-sobolev. In Po-Ling Loh and Maxim Raginsky (eds.), *Pro-
520 ceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine
521 Learning Research*, pp. 1–2. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/chewi22a.html>.

522 Andrew M. Childs, Tongyang Li, Jin-Peng Liu, Chunhao Wang, and Ruizhe Zhang. Quantum
523 algorithms for sampling log-concave distributions and estimating normalizing constants. In
524 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in
525 Neural Information Processing Systems*, volume 35, pp. 23205–23217. Curran Associates, Inc.,
526 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/933e953353c25ec70477ef28e45a2dcc-Paper-Conference.pdf.

527 Baptiste Claudon, Jean-Philip Piquemal, and Pierre Monmarché. Quantum speedup for nonre-
528 versible markov chains, 2025a. URL <https://arxiv.org/abs/2501.05868>.

529 Baptiste Claudon, Pablo Rodenas-Ruiz, Jean-Philip Piquemal, and Pierre Monmarché. Quantum cir-
530 cuits for the metropolis-hastings algorithm, 2025b. URL <https://arxiv.org/abs/2506.11576>.

540 Arjan Cornelissen, Yassine Hamoudi, and Sofiene Jerbi. Near-optimal quantum algorithms for
 541 multivariate mean estimation. In *Proceedings of the 54th Annual ACM SIGACT Symposium*
 542 on *Theory of Computing*, STOC '22. ACM, June 2022. doi: 10.1145/3519935.3520045. URL
 543 <http://dx.doi.org/10.1145/3519935.3520045>.

544 Ben Cousins and Santosh Vempala. Gaussian cooling and $o^*(n^3)$ algorithms for volume and gaussian
 545 volume. *SIAM Journal on Computing*, 47(3):1237–1273, 2018. doi: 10.1137/15M1054250.
 546 URL <https://doi.org/10.1137/15M1054250>.

547 Aniket Das, Dheeraj M. Nagaraj, and Anant Raj. Utilising the clt structure in stochastic gradient
 548 based sampling : Improved analysis and faster algorithms. In Gergely Neu and Lorenzo
 549 Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of
 550 *Proceedings of Machine Learning Research*, pp. 4072–4129. PMLR, 12–15 Jul 2023. URL
 551 <https://proceedings.mlr.press/v195/das23b.html>.

552 Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental
 553 gradient method with support for non-strongly convex composite objectives. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.,
 554 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/file/ede7e2b6d13a41ddf9f4bdef84fdc737-Paper.pdf.

555 Kumar Avinava Dubey, Sashank J. Reddi, Sinead A Williamson, Barnabas Poczos, Alexander J Smola, and Eric P Xing. Variance reduction in stochastic gradient langevin dynamics. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.,
 556 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/9b698eb3105bd82528f23d0c92dedfc0-Paper.pdf.

557 Alain Durmus and Eric Moulines. High-dimensional bayesian inference via the unadjusted langevin
 558 algorithm, 2018. URL <https://arxiv.org/abs/1605.01559>.

559 Daan Frenkel and Berend Smit. *Understanding Molecular Simulation: From Algorithms to Applications*, volume 1 of *Computational Science Series*. Academic Press, San Diego, second edition,
 560 2002.

561 András Gilyén, Srinivasan Arunachalam, and Nathan Wiebe. *Optimizing quantum optimization*
 562 *algorithms via faster quantum gradient computation*, pp. 1425–1444. Society for Industrial and
 563 Applied Mathematics, January 2019. ISBN 9781611975482. doi: 10.1137/1.9781611975482.87.
 564 URL <http://dx.doi.org/10.1137/1.9781611975482.87>.

565 Lov K. Grover. A fast quantum mechanical algorithm for database search. In *Proceedings of*
 566 *the Twenty-eighth Annual ACM Symposium on Theory of Computing*, pp. 212–219. ACM, 1996.
 567 arxivquant-ph/9605043.

568 Xunpeng Huang, Difan Zou, Hanze Dong, Yian Ma, and Tong Zhang. Faster sampling via stochastic
 569 gradient proximal sampler. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian
 570 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st Interna-*
 571 *tional Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*
 572 *Research*, pp. 20559–20596. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/huang24aj.html>.

573 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive vari-
 574 ance reduction. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger
 575 (eds.), *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.,
 576 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/ac1dd209cbcc5e5d1c6e28598e8cbbe8-Paper.pdf.

577 Stephen P. Jordan. Fast quantum algorithm for numerical gradient estimation. *Physical Review*
 578 *Letters*, 95(5), July 2005. ISSN 1079-7114. doi: 10.1103/physrevlett.95.050501. URL <http://dx.doi.org/10.1103/PhysRevLett.95.050501>.

594 Yuri Kinoshita and Taiji Suzuki. Improved convergence rate of stochastic gradient langevin
 595 dynamics with variance reduction and its application to optimization. In S. Koyejo,
 596 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neu-
 597 ral Information Processing Systems*, volume 35, pp. 19022–19034. Curran Associates, Inc.,
 598 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/78e839f96568985d18463044a064ea0f-Paper-Conference.pdf.

600
 601 Tongyang Li and Ruizhe Zhang. Quantum speedups of optimizing approximately convex functions
 602 with applications to logarithmic regret stochastic convex bandits. In *Proceedings of the 36th
 603 International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY,
 604 USA, 2024. Curran Associates Inc. ISBN 9781713871088.

605
 606 Chengchang Liu, Chaowen Guan, Jianhao He, and John C.S. Lui. Quantum algorithms for non-
 607 smooth non-convex optimization. In *The Thirty-eighth Annual Conference on Neural Information
 608 Processing Systems*, 2024.

609 László Lovász and Santosh Vempala. Simulated annealing in convex bodies and an $o^*(n^4)$ volume
 610 algorithm. *Journal of Computer and System Sciences*, 72(2):392–417, 2006. ISSN 0022-0000.
 611 doi: <https://doi.org/10.1016/j.jcss.2005.08.004>. URL <https://www.sciencedirect.com/science/article/pii/S0022000005000966>. JCSS FOCS 2003 Special Issue.

613
 614 Yi-An Ma, Yuansi Chen, Chi Jin, Nicolas Flammarion, and Michael I. Jordan. Sampling can
 615 be faster than optimization. *Proceedings of the National Academy of Sciences*, 116(42):
 616 20881–20885, September 2019. ISSN 1091-6490. doi: 10.1073/pnas.1820003116. URL
 617 <http://dx.doi.org/10.1073/pnas.1820003116>.

618 Frederic Magniez, Ashwin Nayak, Jeremie Roland, and Miklos Santha. Search via quantum walk.
 619 In *Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing*, STOC
 620 '07, pp. 575–584, New York, NY, USA, 2007. Association for Computing Machinery. ISBN
 621 9781595936318. doi: 10.1145/1250790.1250874. URL <https://doi.org/10.1145/1250790.1250874>.

623
 624 Lam M. Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for
 625 machine learning problems using stochastic recursive gradient. In Doina Precup and Yee Whye
 626 Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70
 627 of *Proceedings of Machine Learning Research*, pp. 2613–2621. PMLR, 06–11 Aug 2017. URL
 628 <https://proceedings.mlr.press/v70/nguyen17b.html>.

629
 630 F. Otto and C. Villani. Generalization of an inequality by talagrand and links with the logarithmic
 631 sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000. ISSN 0022-1236.
 632 doi: <https://doi.org/10.1006/jfan.1999.3557>. URL <https://www.sciencedirect.com/science/article/pii/S0022123699935577>.

633
 634 Guneukan Ozgul, Xiantao Li, Mehrdad Mahdavi, and Chunhao Wang. Stochastic quantum sam-
 635 pling for non-logconcave distributions and estimating partition functions. In Ruslan Salakhut-
 636 dinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix
 637 Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, vol-
 638 ume 235 of *Proceedings of Machine Learning Research*, pp. 38953–38982. PMLR, 21–27 Jul
 639 2024. URL <https://proceedings.mlr.press/v235/ozgul24a.html>.

640
 641 Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic
 642 gradient langevin dynamics: a nonasymptotic analysis. In Satyen Kale and Ohad Shamir (eds.),
 643 *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine
 644 Learning Research*, pp. 1674–1703. PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/raginsky17a.html>.

645
 646 Abhishek Roy, Lingqing Shen, Krishnakumar Balasubramanian, and Saeed Ghadimi. Stochastic
 647 zeroth-order discretizations of langevin diffusions for bayesian inference, 2021. URL <https://arxiv.org/abs/1902.01373>.

648 Aaron Sidford and Chenyi Zhang. Quantum speedups for stochastic optimization. In A. Oh,
 649 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neu-
 650 ral Information Processing Systems*, volume 36, pp. 35300–35330. Curran Associates, Inc.,
 651 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/
 652 file/6ed9931d6e1fb6a85efab2c014a47e1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/6ed9931d6e1fb6a85efab2c014a47e1-Paper-Conference.pdf).

653 R. Somma, S. Boixo, and H. Barnum. Quantum simulated annealing, 2007. URL <https://arxiv.org/abs/0712.1008>.

656 R. D. Somma, S. Boixo, H. Barnum, and E. Knill. Quantum simulations of classical annealing
 657 processes. *Phys. Rev. Lett.*, 101:130504, Sep 2008. doi: 10.1103/PhysRevLett.101.130504. URL
 658 <https://link.aps.org/doi/10.1103/PhysRevLett.101.130504>.

659 M. Szegedy. Quantum speed-up of markov chain based algorithms. In *45th Annual IEEE Symposium
 660 on Foundations of Computer Science*, pp. 32–41, 2004. doi: 10.1109/FOCS.2004.53.

662 Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics.
 663 Springer New York, NY, 1 edition, 2009. ISBN 978-0-387-79051-0. doi: 10.1007/b13794.

664 Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Convex optimization
 665 using quantum oracles. *Quantum*, 4:220, January 2020. ISSN 2521-327X. doi: 10.22331/
 666 q-2020-01-13-220. URL <https://doi.org/10.22331/q-2020-01-13-220>.

667 Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted langevin algorithm:
 668 Isoperimetry suffices. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox,
 669 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Cur-
 670 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/
 671 paper/2019/file/65a99bb7a3115fdede20da98b08a370f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/65a99bb7a3115fdede20da98b08a370f-Paper.pdf).

672 Yu-Xiang Wang, Stephen Fienberg, and Alex Smola. Privacy for free: Posterior sampling and
 673 stochastic gradient monte carlo. In Francis Bach and David Blei (eds.), *Proceedings of the
 674 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine
 675 Learning Research*, pp. 2493–2502, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/wangg15.html>.

676 Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics.
 677 In *Proceedings of the 28th International Conference on International Conference on Machine
 678 Learning*, ICML’11, pp. 681–688, Madison, WI, USA, 2011. Omnipress. ISBN 9781450306195.

679 Paweł Wocjan and Anura Abeyesinghe. Speedup via quantum sampling. *Phys. Rev. A*, 78:042336,
 680 Oct 2008. doi: 10.1103/PhysRevA.78.042336. URL <https://link.aps.org/doi/10.1103/PhysRevA.78.042336>.

681 Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. Global convergence of langevin
 682 dynamics based algorithms for nonconvex optimization. In S. Bengio, H. Wal-
 683 lach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Ad-
 684 vances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
 685 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/
 686 file/9c19a2aa1d84e04b0bd4bc888792bd1e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/9c19a2aa1d84e04b0bd4bc888792bd1e-Paper.pdf).

687 Yixin Zhang, Chenyi Zhang, Cong Fang, Liwei Wang, and Tongyang Li. Quantum algorithms
 688 and lower bounds for finite-sum optimization. In Ruslan Salakhutdinov, Zico Kolter, Katherine
 689 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Pro-
 690 ceedings of the 41st International Conference on Machine Learning*, volume 235 of *Pro-
 691 ceedings of Machine Learning Research*, pp. 60244–60270. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/zhang24bz.html>.

692 Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient
 693 langevin dynamics. In Satyen Kale and Ohad Shamir (eds.), *Proceedings of the 2017 Conference
 694 on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 1980–2022.
 695 PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/zhang17b.html>.

702 Difan Zou and Quanquan Gu. On the convergence of hamiltonian monte carlo with stochas-
703 tic gradients. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International*
704 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
705 pp. 13012–13022. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/zou21b.html>.

707 Difan Zou, Pan Xu, and Quanquan Gu. Stochastic gradient hamiltonian monte carlo methods with
708 recursive variance reduction. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc,
709 E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32.
710 Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/c3535febaff29fc7c0d20cbe94391c7-Paper.pdf.

712 Difan Zou, Pan Xu, and Quanquan Gu. Faster convergence of stochastic gradient langevin dynam-
713 ics for non-log-concave sampling. In Cassio de Campos and Marloes H. Maathuis (eds.), *Pro-
714 ceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161
715 of *Proceedings of Machine Learning Research*, pp. 1152–1162. PMLR, 27–30 Jul 2021. URL
716 <https://proceedings.mlr.press/v161/zou21a.html>.

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756 A PROOFS FOR HAMILTONIAN MONTE CARLO IN STRONGLY CONVEX CASE
757758 We start with the following result in Zou & Gu (2021) that quantifies the convergence of the stochastic
759 gradient Hamiltonian Monte Carlo algorithm in Wasserstein distance.
760761 **Theorem A.1** (Theorem 4.4 in Zou & Gu (2021)). *Under Assumptions 4.1 and 4.2, let $D = \|\mathbf{x}^0 - \arg \min_{\mathbf{x}} (f(\mathbf{x}))\|$ and μ_T be the distribution of the iterate \mathbf{x}^T , then if the step size satisfies $\eta = O(L^{1/2}\sigma^{-2}\kappa^{-1} \wedge L^{-1/2})$ and $K = 1/(4\sqrt{L}\eta)$, the output of HMC satisfies*

764
$$W_2(\mu_T, \pi) \leq (1 - (128\kappa)^{-1})^{\frac{T}{2}} (2D + 2d/\mu)^{1/2} + \Gamma_1 \eta^{1/2} + \Gamma_2 \eta, \quad (14)$$

765

766 where $\Gamma_1^2 = O(L^{-3/2}\sigma^2\kappa^2)$ and $\Gamma_2^2 = O(\kappa^2(LD + \kappa d + L^{-1/2}\sigma^2\eta))$ where $\sigma^2 =$
767 $\max_{t \leq T} \mathbb{E}\|\mathbf{g}(\mathbf{x}_k, \xi_k) - \nabla f(\mathbf{x}_k)\|^2$ is the upper bound on the variance of the gradients in the
768 trajectory of SG-HMC algorithm.
769770 This is a generic result that applies to any HMC algorithm under Assumptions 4.1 and 4.2 that uses
771 stochastic gradients with variance upper bounded by σ^2 . Note that we do not assume a uniform
772 upper bound for σ that is independent of problem parameters. Instead, the variance upper bound
773 depends on the trajectory of the algorithm, which can be characterized using theoretical analysis.
774775 A.1 PROOF OF QHMC-SVRG
776777 **Lemma A.2.** *Under Assumption 4.2, if the initial point satisfies $\|\mathbf{x}^0 - \mathbf{x}^*\| \leq \frac{d}{\mu}$, then it holds that*

778
$$\mathbb{E}_i \|\mathbf{g}(\mathbf{x}_k, \xi) - \nabla f(\mathbf{x}_k)\|^2 \leq L^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2, \quad (15)$$

779

780 where $\tilde{\mathbf{x}} = \mathbf{x}_{k' < k}$ is the last iteration the full gradient is computed.
781782 *Proof.* The proof simply follows from the definition of variance in the SVRG algorithm and the
783 smoothness of each component.
784

785
$$\mathbb{E}_i \|\mathbf{g}_i(\mathbf{x}_k, \xi) - \nabla f(\mathbf{x}_k)\|^2 \leq \mathbb{E}_i \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}}) + f(\tilde{\mathbf{x}}) - f(\mathbf{x}_k)\|^2 \quad (16)$$

786

787
$$\leq \mathbb{E}_i \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}})\|^2 \quad (17)$$

788
$$\leq L^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2. \quad (18)$$

789

□

790 Lemma A.2 allows us to set the target variance in quantum mean estimation to be $L^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|/b^2$.
791 Hence, each mean estimation call takes $\mathcal{O}(d^{1/2}b)$ gradient evaluations by Lemma 3.1. The following
792 lemma characterizes the variance of the stochastic gradients along the trajectory of the algorithm.
793794 **Lemma A.3** (Modified Lemma C.2 in Zou & Gu (2021)). *Let $\mathbf{g}(\mathbf{x}_k, \xi_k)$ be the vector computed using
795 the unbiased quantum mean estimation algorithm in QHMC-SVRG. Then, under Assumption 4.2,*

796
$$\mathbb{E} \|\mathbf{g}(\mathbf{x}_k, \xi_k) - \nabla f(\mathbf{x}_k)\|^2 \leq \frac{768m^2L^2\eta^2\kappa d}{b^2}, \quad (19)$$

797

798 where the expectation is over both the iterate \mathbf{x}_k and the noise in quantum mean estimation ξ_k .
799800 Next, we prove the main theorem for QSVRG-HMC.
801802 **Theorem 4.3** (Main Theorem for QSVRG-HMC). *Let μ_k be the distribution of \mathbf{x}_k in QSVRG-HMC
803 algorithm. Suppose that f satisfies Assumptions 4.1 and 4.2. Given that the initial point \mathbf{x}_0 satisfies
804 $\|\mathbf{x}_0 - \arg \min_{\mathbf{x}} f(\mathbf{x})\| \leq \frac{d}{\mu}$, then, for $\eta = \mathcal{O}(\frac{\epsilon}{L^{1/2}d^{1/2}\kappa^{3/2}})$, $S = \tilde{\mathcal{O}}\left(\frac{Ld^{1/2}\kappa^{3/2}}{\epsilon}\right)$, $T = \tilde{\mathcal{O}}(1)$,
805 $b = \mathcal{O}\left(\frac{L^{1/8}\epsilon^{1/4}n^{1/2}}{d^{1/8}\kappa^{3/8}} \vee 1\right)$, and $m = n/b$, we have*

806
$$W_2(\mu_{ST}, \pi) \leq \epsilon.$$

807

808 The total query complexity to the stochastic gradient oracle is $\tilde{\mathcal{O}}\left(\frac{Ld^{1/2}\kappa^{3/2}}{\epsilon} + \frac{L^{9/8}d^{7/8}\kappa^{3/4}n^{1/2}}{\epsilon^{3/4}}\right)$.
809

810 *Proof.* By the choice of η in the theorem statement and the variance upper bound in Lemma A.3,
811 $\eta = \mathcal{O}(L^{1/2}\sigma^{-2}\kappa^{-1} \wedge L^{-1/2})$. Therefore, by Theorem A.1, for $K = \frac{1}{4\sqrt{L}\eta}$, we have
812

$$813 \quad 814 \quad W_2(\mu_T, \pi) \leq (1 - (128\kappa)^{-1})^{\frac{T}{2}} (2D + 2d/\mu)^{1/2} + \Gamma_1 \eta^{1/2} + \Gamma_2 \eta \quad (20)$$

815 where,

$$816 \quad 817 \quad \Gamma_1^2 = \mathcal{O}\left(\frac{L^{1/2}m^2\kappa^3d\eta^2}{b^2}\right), \quad (21)$$

$$818 \quad 819 \quad \Gamma_2^2 = \mathcal{O}\left(\kappa^3d + \frac{L^{3/2}m^2\kappa^3d\eta^3}{b^2}\right). \quad (22)$$

820 We set $bm = \mathcal{O}(n)$. The first term in Eq. (20) is $\mathcal{O}(\epsilon)$ when $T = \tilde{\mathcal{O}}(\log(1/\epsilon))$. The
821 last two terms in Eq. (14) for QHMC-SVRG become $\mathcal{O}\left(\frac{L^{1/4}d^{1/2}\kappa^{3/2}\eta^{3/2}n}{b^2} + d^{1/2}\kappa^{3/2}\eta\right)$. For
822 $b = \mathcal{O}(d^{-1/8}\kappa^{-3/8}\epsilon^{1/4}n^{1/2}L^{1/8} \vee 1)$ and $\eta = \mathcal{O}(\epsilon\kappa^{-3/2}d^{-1/2})$, the bias term becomes
823 $\mathcal{O}(\epsilon)$. Using Lemma 3.1, the number of gradient calculations scales as $\tilde{\mathcal{O}}(Ld^{1/2}\kappa^{3/2}\epsilon^{-1} +$
824 $L^{9/8}d^{7/8}\kappa^{3/4}\epsilon^{-3/4}n^{1/2})$. \square

825 A.2 PROOF OF QCV-HMC

826 **Lemma A.4** (Modified Lemma C.4 in Zou & Gu (2021)). *Let $\mathbf{g}(\mathbf{x}_k, \xi_k)$ be the vector computed
827 using the unbiased quantum mean estimation algorithm in QHMC-CV. Then, under Assumption 4.2,*

$$828 \quad 829 \quad \mathbb{E}\|\mathbf{g}(\mathbf{x}_k, \xi_k) - \nabla f(\mathbf{x}_k)\|^2 \leq \frac{688Ld\kappa}{b^2},$$

830 where the expectation is over both the iterate \mathbf{x}_k and the noise in quantum mean estimation ξ_k .

831 Next we prove the main result.

832 **Theorem 4.4** (Main Theorem for QCV-HMC). *Let μ_k be the distribution of \mathbf{x}_k in QCV-HMC algo-
833 rithm. Suppose that f satisfies Assumptions 4.1 and 4.2. Given that the initial point \mathbf{x}_0 satisfies
834 $\|\mathbf{x}_0 - \arg \min_{\mathbf{x}} f(\mathbf{x})\| \leq \frac{d}{\mu}$, then, for $\eta = \mathcal{O}\left(\frac{\epsilon}{L^{1/2}d^{1/2}\kappa^{3/2}}\right)$, $S = \tilde{\mathcal{O}}\left(\frac{Ld^{1/2}\kappa^{3/2}}{\epsilon}\right)$, $T = \tilde{\mathcal{O}}(1)$, and
835 $b = \mathcal{O}\left(\frac{d^{1/4}\kappa^{3/4}}{L^{1/4}\epsilon^{1/2}} \vee 1\right)$, we have*

$$836 \quad W_2(\mu_{ST}, \pi) \leq \epsilon.$$

837 The total query complexity to the stochastic gradient oracle is $\tilde{\mathcal{O}}\left(\frac{Ld^{5/4}\kappa^{9/4}}{\epsilon^{3/2}}\right)$.

838 *Proof.* By the choice of η in the theorem statement and the variance upper bound in Lemma A.4,
839 $\eta = \mathcal{O}(L^{1/2}\sigma^{-2}\kappa^{-1} \wedge L^{-1/2})$. Therefore, by Theorem A.1, for $K = \frac{1}{4\sqrt{L}\eta}$, we have
840

$$841 \quad 842 \quad W_2(\mu_T, \pi) \leq (1 - (128\kappa)^{-1})^{\frac{T}{2}} (2D + 2d/\mu)^{1/2} + \Gamma_1 \eta^{1/2} + \Gamma_2 \eta, \quad (23)$$

843 where,

$$844 \quad 845 \quad \Gamma_1 = \mathcal{O}\left(\frac{L^{-1/2}\kappa^3d}{b^2}\right), \quad (24)$$

$$846 \quad \Gamma_2 = \mathcal{O}(\kappa^3d). \quad (25)$$

847 The first term in Eq. (23) is $\mathcal{O}(\epsilon)$ when $T = \tilde{\mathcal{O}}(1)$. The last two terms in Eq. (23) for
848 QHMC-CV become $\mathcal{O}\left(\frac{L^{-1/4}d^{1/2}\kappa^{3/2}\eta^{1/2}}{b^2} + d^{1/2}\kappa^{3/2}\eta\right)$. For $b = \mathcal{O}(L^{-1/4}d^{1/4}\kappa^{3/4}\epsilon^{-1/2} \vee 1)$ and
849 $\eta = \mathcal{O}(\epsilon d^{-1/2}\kappa^{-3/2})$, the bias term becomes $\mathcal{O}(\epsilon)$. Using Lemma 3.1, the number of gradient
850 calculations scales as $\tilde{\mathcal{O}}(Ld^{1/2}\kappa^{3/2}\epsilon^{-1} + L^{3/4}d^{5/4}\kappa^{9/4}\epsilon^{-3/2}) = \tilde{\mathcal{O}}(Ld^{5/4}\kappa^{9/4}\epsilon^{-3/2})$. \square

864 B PROOFS FOR LSI CASE
865866 **Lemma B.1** (Stochastic-LMC One Step Convergence). *Let μ_k be the distribution of the iterate \mathbf{x}_k ,
867 then if the step size satisfies $\eta = \frac{2}{3\alpha}$,*
868

869
$$\text{KL}(\mu_{k+1} || \pi) \leq e^{-3\alpha\eta/2} \left[\left(1 + \frac{32\eta^3 L^4}{\alpha} \right) \text{KL}(\mu_k || \pi) + 6\eta\sigma_k^2 + 16\eta^2 dL^2 \right], \quad (26)$$

870
871

872 where $\sigma_k^2 = \mathbb{E}_{\mathbf{x}_k, \xi_k} \|\mathbf{g}(\mathbf{x}_k, \xi_k) - \nabla f(\mathbf{x}_k)\|^2$.
873874 *Proof.* We compare one step of LMC starting at \mathbf{x}_k with stochastic gradients $\mathbf{g}(\mathbf{x}_k, \xi_k)$ to the output
875 of continuous Langevin SDE (Eq. (7)) starting at \mathbf{x}_k with true gradient $\nabla f(\mathbf{x}_t)$ after time η . This
876 technique has been used to establish the convergence of unadjusted Langevin algorithm with full
877 gradients under isoperimetry by Vempala & Wibisono (2019). We extend the analysis by Vempala
878 & Wibisono (2019) to the stochastic gradient LMC. Assume that the initial point \mathbf{x}_k and $\mathbf{g}(\mathbf{x}_k, \xi_k)$
879 obey the joint distribution μ_0 . The randomness on $\mathbf{g}(\mathbf{x}_k, \xi_k)$ depends both on the randomness on
880 \mathbf{x}_k and the randomness in the quantum mean estimation algorithm. Then, one step update of LMC
881 algorithm with stochastic gradient yields,
882

883
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \mathbf{g}(\mathbf{x}_k, \xi_k) + \sqrt{2\eta} \epsilon_k.$$

884

885 Alternatively, \mathbf{x}_{k+1} can be written as the solution of the following SDE at time $t = \eta$,
886

887
$$d\mathbf{x}_t = -\mathbf{g}_k dt + \sqrt{2} d\mathbf{W}_t$$

888

889 where $\mathbf{g}_k = \mathbf{g}(\mathbf{x}_k, \xi_k)$ and \mathbf{W}_t is the standard Brownian motion starting at $\mathbf{W}_0 = 0$. Let
890 $\mu_t(\mathbf{x}_k, \mathbf{g}_k, \mathbf{x}_t)$ be the joint distribution of \mathbf{x}_k , \mathbf{g}_k , and \mathbf{x}_t at time t . Each expectation in the proof is
891 over this joint distribution unless specified otherwise.
892893 Consider the following stochastic differential equation
894

895
$$d\mathbf{X} = \mathbf{v}(\mathbf{X}) dt + \sqrt{2} d\mathbf{W},$$

896

897 where \mathbf{v} is a smooth vector field and \mathbf{W} is the Brownian motion with $\mathbf{W}_0 = 0$. The Fokker-Planck
898 equation describes the evolution of probability density function μ_t as follows:
899

900
$$\frac{\partial \mu_t}{\partial t} = -\nabla \cdot (\mu_t \mathbf{v}) + \Delta \mu_t, \quad (27)$$

901

902 where $\nabla \cdot$ is the divergence operator and Δ is the Laplacian. Then, the Fokker Planck equation gives
903 the following evolution for the marginal density $\mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k) = \mu_t(\mathbf{x}_t = \mathbf{x}|\mathbf{x}_k, \mathbf{g}_k)$,
904

905
$$\frac{\partial \mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k)}{\partial t} = \nabla \cdot (\mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k) \mathbf{g}_k) + \Delta \mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k). \quad (28)$$

906

907 Taking the expectation over both sides with respect to $(\mathbf{x}_k, \mathbf{g}_k) \sim \mu_0$,
908

909
$$\frac{\partial \mu_t(\mathbf{x})}{\partial t} = \mathbb{E}_{(\mathbf{x}_k, \mathbf{g}_k) \sim \mu_0} [\nabla \cdot (\mu_t(\mathbf{x}|\mathbf{x}_k) \mathbf{g}_k)] + \mathbb{E}_{(\mathbf{x}_k, \mathbf{g}_k) \sim \mu_0} [\Delta \mu_t(\mathbf{x}|\mathbf{x}_k)] \quad (29)$$

910

911
$$= \int_{\mathbb{R}^d} \nabla \cdot (\mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k) \mathbf{g}_k) \mu_0(\mathbf{x}_k, \mathbf{g}_k) d\mathbf{x}_k d\mathbf{g}_k + \int_{\mathbb{R}^d} \Delta \mu_t(\mathbf{x}|\mathbf{x}_k, \mathbf{g}_k) \mu_0(\mathbf{x}_k, \mathbf{g}_k) d\mathbf{x}_k d\mathbf{g}_k \quad (30)$$

912

913
$$= \int_{\mathbb{R}^d} \nabla \cdot (\mu_t(\mathbf{x}) \mu(\mathbf{x}_k, \mathbf{g}_k | \mathbf{x}_t = \mathbf{x}) \mathbf{g}_k) d\mathbf{x}_k d\mathbf{g}_k + \Delta \mu_t(\mathbf{x}) \quad (31)$$

914

915
$$= \nabla \cdot \left(\mu_t(\mathbf{x}) \mathbb{E}[\mathbf{g}_k - \nabla f(\mathbf{x}_k) | \mathbf{x}_t = \mathbf{x}] + \mu_t(\mathbf{x}) \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right). \quad (32)$$

916
917

918 Consider the time derivative of KL divergence between μ_t and π ,
 919

$$920 \quad \frac{d}{dt} \text{KL}(\mu_t || \pi) = \frac{d}{dt} \int_{\mathbb{R}^d} \mu_t(\mathbf{x}) \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x} \quad (33)$$

$$923 \quad = \int_{\mathbb{R}^d} \frac{\partial \mu_t(\mathbf{x})}{\partial t} \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x}_t + \int_{\mathbb{R}^d} \mu_t(\mathbf{x}) \frac{\partial}{\partial t} \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x} \quad (34)$$

$$926 \quad = \int_{\mathbb{R}^d} \frac{\partial \mu_t(\mathbf{x})}{\partial t} \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x}_t + \int_{\mathbb{R}^d} \frac{\partial \mu_t(\mathbf{x})}{\partial t} d\mathbf{x} \quad (35)$$

$$929 \quad = \int_{\mathbb{R}^d} \frac{\partial \mu_t(\mathbf{x})}{\partial t} \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x}_t. \quad (36)$$

932 The last term in the third equality vanishes since the μ_t is probability distribution and its L_1 norm is
 933 always 1. Then the KL divergence evolves as

$$934 \quad \frac{d}{dt} \text{KL}(\mu_t || \pi) = \int_{\mathbb{R}^d} \nabla \cdot \left(\mu_t(\mathbf{x}) \mathbb{E}[\mathbf{g}_k - \nabla f(\mathbf{x}) | \mathbf{x}_t = \mathbf{x}] + \mu_t(\mathbf{x}) \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right) \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x} \quad (37)$$

$$938 \quad = - \int_{\mathbb{R}^d} \mu_t(\mathbf{x}) \left\langle \mathbb{E}[\mathbf{g}_k - \nabla f(\mathbf{x}) | \mathbf{x}_t = \mathbf{x}] + \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right), \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right\rangle d\mathbf{x} \quad (38)$$

$$942 \quad = - \int_{\mathbb{R}^d} \mu_t(\mathbf{x}) \left\| \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right\|^2 d\mathbf{x} + \mathbb{E} \left\langle \nabla f(\mathbf{x}_t) - \mathbf{g}_k, \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right\rangle. \quad (39)$$

945 The second term can be bounded as follows:
 946

$$947 \quad \mathbb{E} \left\langle \nabla f(\mathbf{x}_t) - \mathbf{g}_k, \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right\rangle \leq \mathbb{E} \left[\left\| \nabla f(\mathbf{x}_t) - \mathbf{g}_k \right\|^2 + \frac{1}{4} \left\| \nabla \log \left(\frac{\mu_t(\mathbf{x})}{\pi(\mathbf{x})} \right) \right\|^2 \right] \quad (40)$$

$$950 \quad = \mathbb{E} \left\| \nabla f(\mathbf{x}_t) - \mathbf{g}_k \right\|^2 + \frac{1}{4} \text{FI}(\mu_t || \pi) \quad (41)$$

$$952 \quad = \mathbb{E} \left\| \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k) - \mathbf{g}_k \right\|^2 + \frac{1}{4} \text{FI}(\mu_t || \pi) \quad (42)$$

$$955 \quad \leq 2\mathbb{E} \left\| \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_k) \right\|^2 + 2\mathbb{E}_{\mu_t(\mathbf{x}_t, \mathbf{x}_k)} \left\| \nabla f(\mathbf{x}_k) - \mathbf{g}_k \right\|^2 \quad (43)$$

$$957 \quad + \frac{1}{4} \text{FI}(\mu_t || \pi). \quad (44)$$

959 The first inequality holds since $\langle a, b \rangle \leq a^2 + \frac{b^2}{4}$. The last line follows from Young's inequality.
 960 Furthermore, using Lipschitzness of gradients of f , we have
 961

$$962 \quad \mathbb{E} \left\| \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_k) \right\|^2 \leq L^2 \mathbb{E} \left\| \mathbf{x}_t - \mathbf{x}_k \right\|^2 \quad (45)$$

$$963 \quad \leq L^2 \mathbb{E} \left\| -t\mathbf{g}_k + \sqrt{2t}\boldsymbol{\epsilon}_k \right\|^2 \quad (46)$$

$$964 \quad = t^2 L^2 \mathbb{E}_{\mu_0} \left\| \mathbf{g}_k \right\|^2 + 2tdL^2. \quad (47)$$

966 Plugging back these into the time derivative of KL divergence, we have
 967

$$968 \quad \frac{d}{dt} \text{KL}(\mu_t || \pi) \leq -\frac{3}{4} \text{FI}(\mu_t || \pi) + 2t^2 L^2 \mathbb{E}_{\mu_0} \left\| \mathbf{g}_k \right\|^2 + 2\mathbb{E}_{\mu_0} \left\| \nabla f(\mathbf{x}_k) - \mathbf{g}_k \right\|^2 + 4tdL^2 \quad (48)$$

$$970 \quad \leq -\frac{3}{4} \text{FI}(\mu_t || \pi) + (4t^2 L^2 + 2) \mathbb{E}_{\mu_0} \left\| \nabla f(\mathbf{x}_k) - \mathbf{g}_k \right\|^2 + 4t^2 L^2 \mathbb{E}_{\mu_0} \left\| \nabla f(\mathbf{x}_k) \right\|^2 + 4tdL^2. \quad (49)$$

972 The third term can be bounded as follows: We choose an optimal coupling $\mathbf{x}_k \sim \mu_0(\mathbf{x}_k)$ and $\mathbf{x}^* \sim \pi$
 973 so that $\mathbb{E}\|\mathbf{x}_k - \mathbf{x}^*\| = \mathbb{W}_2(\mu_0, \pi)^2$, then using Young's inequality and smoothness of f ,
 974

$$\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}_k)\|^2 \leq 2\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)\|^2 + 2\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}^*)\|^2 \quad (50)$$

$$\leq 2L^2\mathbb{E}_{\mu_0}\|\mathbf{x}_k - \mathbf{x}^*\|^2 + 2\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}^*)\|^2 \quad (51)$$

$$\leq 2L^2\mathbb{W}_2(\mu_0, \pi)^2 + 2dL \quad (52)$$

$$\leq \frac{4L^2}{\alpha}\text{KL}(\mu_0\|\pi) + 2dL. \quad (53)$$

981 The last inequality follows from Talgrand's inequality. Hence for $t \leq \eta$ and $\eta \leq \frac{1}{2L}$, we have
 982

$$\frac{d}{dt}\text{KL}(\mu_t\|\pi) \leq -\frac{3}{4}\text{FI}(\mu_t\|\pi) + (4t^2L^2 + 2)\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}_k) - \mathbf{g}_k\|^2 + \frac{16t^2L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 4tdL^2 + 8t^2dL^3 \quad (54)$$

$$\leq -\frac{3\alpha}{2}\text{KL}(\mu_t\|\pi) + (4t^2L^2 + 2)\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}_k) - \mathbf{g}_k\|^2 + \frac{16t^2L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 4tdL^2 + 8t^2dL^3 \quad (55)$$

$$\leq -\frac{3\alpha}{2}\text{KL}(\mu_t\|\pi) + 3\mathbb{E}_{\mu_0}\|\nabla f(\mathbf{x}_k) - \mathbf{g}_k\|^2 + \frac{16\eta^2L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 8\eta dL^2 \quad (56)$$

$$\leq -\frac{3\alpha}{2}\text{KL}(\mu_t\|\pi) + 3\sigma_k^2 + \frac{16\eta^2L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 8\eta dL^2. \quad (57)$$

994 The second inequality is due to Eq. (11). Equivalently, we can write,
 995

$$\frac{d}{dt}(e^{3\alpha t/2}\text{KL}(\mu_t\|\pi)) \leq e^{3\alpha t/2} \left(3\sigma_k^2 + \frac{16\eta^2L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 8\eta dL^2 \right). \quad (58)$$

998 Integrating from $t = 0$ to $t = \eta$ gives,
 999

$$e^{3\alpha\eta/2}\text{KL}(\mu_\eta\|\pi) - \text{KL}(\mu_0\|\pi) \leq 6\eta\sigma_k^2 + \frac{32\eta^3L^4}{\alpha}\text{KL}(\mu_0\|\pi) + 16\eta^2dL^2 \quad (59)$$

1002 for $\eta \leq \frac{2}{3\alpha}$. Rearranging the terms,
 1003

$$\text{KL}(\mu_\eta\|\pi) \leq e^{-3\alpha\eta/2} \left[\left(1 + \frac{32\eta^3L^4}{\alpha} \right) \text{KL}(\mu_0\|\pi) + 6\eta\sigma_k^2 + 16\eta^2dL^2 \right]. \quad (60)$$

1006 Renaming $\mu_0 = \mu_k$ and $\mu_\eta = \mu_{k+1}$, we obtain the result in the statement.
 1007

□

1010 The statement in Lemma B.1 is generic and can be applied to any LMC algorithm with stochastic
 1011 gradients with bounded variance on the trajectory of the algorithm. Note that this is different from
 1012 assuming that the variance is uniformly upper bounded. Instead, we set inner loop and variance
 1013 reduction parameters so that the variance does not explode along the trajectory of the algorithm.

1014 B.1 PROOF OF QSVRG-LMC

1016 We start with the following lemma that characterizes the variance of the quantum stochastic gradients
 1017 in QSVRG-LMC in terms of the distance between the current iterate and the reference point where
 1018 the full gradient is computed.
 1019

1020 **Lemma B.2.** *Let $\tilde{\mathbf{x}}$ be any iteration where QSVRG-LMC computes the full gradient. Then under
 1021 Assumption 4.2, the quantum stochastic gradient \mathbf{g}_k at \mathbf{x}_k that is computed using $\tilde{\mathbf{x}}$ as a reference
 1022 point in QSVRG-LMC satisfies*

$$\mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2] \leq \frac{L^2\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2}{b^2} \quad (61)$$

1023 using $\tilde{O}(d^{1/2}b)$ gradient computations.
 1024

1026 *Proof.* Recall that SVRG computes the stochastic gradient $\tilde{\mathbf{g}}$ at \mathbf{x}_k by the following.
1027

$$1028 \quad \tilde{\mathbf{g}}_k = \nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}), \quad (62)$$

1029 where $\tilde{\mathbf{x}}$ is the last iteration the full gradient is computed and i is a component randomly chosen
1030 from $[n]$. Let $\sigma_k^2 = \mathbb{E}\|\tilde{\mathbf{g}}_k - \nabla f(\mathbf{x}_k)\|^2$. Then, σ_k^2 can be bounded in terms of the distance between
1031 \mathbf{x}_k and $\tilde{\mathbf{x}}$.

$$1032 \quad \sigma_k^2 = \mathbb{E}[\|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) - \nabla f(\mathbf{x}_k)\|^2] \quad (63)$$

$$1034 \quad = \mathbb{E}[\|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}})\|^2] - (\mathbb{E}[\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}})])^2 \quad (64)$$

$$1035 \quad \leq \mathbb{E}[\|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}})\|^2] \quad (65)$$

$$1036 \quad \leq L^2\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2, \quad (66)$$

1038 where the equality follows from the fact that ∇f_i is an unbiased estimator for ∇f and the last line
1039 follows from Assumption 4.2. Hence, using unbiased quantum mean estimation in Lemma 3.1, we
1040 can obtain a random vector \mathbf{g}_k such that,

$$1041 \quad \mathbb{E}\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 \leq \frac{L^2\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2}{b^2} \quad (67)$$

1044 by using $\tilde{O}(d^{1/2}b)$ calls to the gradient oracle. \square
1045

1046 To be able to apply Lemma B.1, we need to characterize the expected upper bound on the variance
1047 of the stochastic gradients over the algorithm trajectory for SVRG.

1048 **Lemma 4.7** (QSVRG–LMC Variance Lemma). *Let $k' < k$ be the last iteration where the full gradient
1049 is computed in QSVRG–LMC and $\sigma_k^2 = \mathbb{E}\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2$. Then, for $\eta^2 \leq \frac{1}{6L^2m^2}$,*

$$1051 \quad \sigma_{k'+l}^2 \leq \frac{16L^4\eta^2}{\alpha} \sum_{r=1}^l \text{KL}(\mu_{k'+r-1} || \pi) + \frac{8\eta dm L^2}{b^2}. \quad (12)$$

1054 *Proof.* Let $\tilde{\mathbf{x}} = \mathbf{x}_{k'}$. Then, by Lemma B.2, quantum stochastic gradient \mathbf{g}_k satisfies
1055

$$1056 \quad \mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2] \leq \frac{L^2\mathbb{E}\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2}{b^2}. \quad (68)$$

1058 Let $\tilde{\mathbf{x}} = \mathbf{y}_0$ and $\mathbf{x}_k = \mathbf{y}_k$, then using the update rule of Langevin Monte Carlo,
1059

$$1060 \quad \mathbb{E}[\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2] = \mathbb{E}\left[\left\|\sum_{r=1}^l (\mathbf{y}_r - \mathbf{y}_{r-1})\right\|^2\right] = \mathbb{E}\left[\left\|\sum_{r=1}^l -\eta \mathbf{g}_{r-1} + \sqrt{2\eta}\epsilon_{r-1}\right\|^2\right] \quad (69)$$

$$1063 \quad \leq \mathbb{E}\left[2\eta^2\left\|\sum_{r=1}^l \mathbf{g}_{r-1}\right\|^2 + 4\eta\left\|\sum_{r=1}^l \epsilon_{r-1}\right\|^2\right] \quad (70)$$

$$1067 \quad \leq 2\eta^2 m \sum_{r=1}^l \mathbb{E}\|\mathbf{g}_{r-1}\|^2 + 4\eta \sum_{r=1}^l \|\epsilon_{r-1}\|^2 \quad (71)$$

$$1070 \quad \leq 2\eta^2 m \sum_{r=1}^l \mathbb{E}\|\mathbf{g}_{r-1}\|^2 + 4\eta dm. \quad (72)$$

1073 The first inequality is due to Young's inequality and the second inequality follows from the fact
1074 that the Gaussian noises at different iterations are independent and the fact that $l \leq m$. Defining
1075 $\sigma_{\max}^2 = \max_k \mathbb{E}\|\sigma_k\|^2$, we can write the first term on the right-hand side in terms of σ_{\max}^2 ,

$$1076 \quad \mathbb{E}[\|\mathbf{g}_r\|^2] = \mathbb{E}\|\mathbf{g}_r - \nabla f(\mathbf{x}_r) + \nabla f(\mathbf{x}_r)\|^2 \quad (73)$$

$$1077 \quad \leq 2\mathbb{E}\|\mathbf{g}_r - \nabla f(\mathbf{x}_r)\|^2 + 2\|\nabla f(\mathbf{x}_r)\|^2 \quad (74)$$

$$1078 \quad \leq 2\sigma_{\max}^2 + \frac{8L^2}{\alpha} \text{KL}(\mu_r || \pi) + 4dL, \quad (75)$$

1080 and using Eq. (68),
 1081

1082
 1083
$$\sigma_{\max}^2 \leq \frac{4L^2m^2\eta^2\sigma_{\max}^2}{b^2} + \frac{16L^4\eta^2m}{b^2\alpha} \sum_{r=1}^l \text{KL}(\mu_{r-1} \parallel \pi) + \frac{8dL^3\eta^2m^2}{b^2} + \frac{4\eta dmL^2}{b^2}. \quad (76)$$

 1084
 1085

1086 If we set $\eta^2 \leq \frac{1}{6L^2m^2}$, we obtain
 1087

1088
 1089
$$\sigma_{k'+l}^2 \leq \frac{32L^4\eta^2m}{b^2\alpha} \sum_{r=1}^l \text{KL}(\mu_{r-1} \parallel \pi) + \frac{8\eta dmL^2}{b^2}. \quad (77)$$

 1090
 1091
 1092
 1093 \square
 1094
 1095
 1096

1097 **Theorem 4.8** (Convergence theorem for QSVRG-LMC). *Assume that $m \leq b^2$. Then, for $\eta \leq \frac{\alpha^2}{24L^2m}$,
 1098 the iterates in QSVRG-LMC satisfy,*
 1099

1100
 1101
$$\text{KL}(\mu_k \parallel \pi) \leq e^{-\alpha\eta k} \text{KL}(\mu_0 \parallel \pi) + \frac{64m\eta dL^2}{\alpha b^2} + \frac{24\eta dL^2}{\alpha}. \quad (13)$$

 1102
 1103
 1104
 1105

1106 *Proof.* Let $l < k$ be the last iteration the full gradient is computed. Then, using Lemmas 4.7 and B.1,
 1107 we can write one step bound as follows.
 1108

1109
 1110
$$\text{KL}(\mu_{k+1} \parallel \pi) \leq e^{-3\alpha\eta/2} \left[\left(1 + \frac{32\eta^3L^4}{\alpha}\right) \text{KL}(\mu_k \parallel \pi) + \frac{192m\eta^3L^4}{b^2\alpha} \sum_{r=l}^k \text{KL}(\mu_r \parallel \pi) + \frac{48m\eta^2dL^2}{b^2} + 16\eta^2dL^2 \right]. \quad (78)$$

 1111
 1112
 1113

1114 First, we claim that the following inequality is true.
 1115

1116
 1117
$$\text{KL}(\mu_{k+1} \parallel \pi) \leq e^{-\alpha\eta k} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2dL^2 + 16\eta^2dL^2b^2}{b^2(1 - e^{-\alpha\eta})}. \quad (79)$$

 1118
 1119

1120 To prove Eq. (79), we use induction. For $k = 1$, the statement holds due to Eq. (78). That is,
 1121

1122
 1123
$$\text{KL}(\mu_1 \parallel \pi) \leq e^{-3\alpha\eta/2} \left[\left(1 + \frac{224\eta^3L^4}{\alpha}\right) \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2dL^2}{b^2} + 16\eta^2dL^2 \right] \quad (80)$$

 1124
 1125

1126
$$\leq e^{-\alpha\eta} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2dL^2}{b^2} + 16\eta^2dL^2 \quad (81)$$

 1127

1128
$$\leq e^{-\alpha\eta} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2dL^2 + 16\eta^2dL^2b^2}{b^2(1 - e^{-\alpha\eta})}. \quad (82)$$

 1129
 1130

1131 The first inequality is due to the fact that $m \leq b^2$. The second inequality holds since
 1132
$$\left(1 + \frac{224\eta^3L^4}{\alpha}\right) \leq \left(1 + \frac{\eta\alpha}{2}\right) \leq e^{\alpha\eta/2}$$
 since $\eta \leq \frac{\alpha}{24L^2m}$. The third inequality follows from the
 1133 fact that $1 - e^{-\alpha\eta} \leq 1$. Next, assume that the statement holds for $k - 1$, and then we prove the k -th

1134 step of induction.

1135

$$1136 \text{KL}(\mu_k \parallel \pi) \leq e^{-3\alpha\eta/2} \left[\left(1 + \frac{32\eta^3 L^4}{\alpha}\right) \text{KL}(\mu_{k-1} \parallel \pi) + \frac{192\eta^3 L^4}{\alpha} \sum_{r=\ell}^{k-1} \text{KL}(\mu_r \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \right] \quad (83)$$

1137

1138

$$1139 \leq e^{-3\alpha\eta/2} \left(1 + \frac{32\eta^3 L^4}{\alpha}\right) \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) \quad (84)$$

1140

1141

$$1142 + e^{-3\alpha\eta/2} \frac{192\eta^3 L^4}{\alpha} \sum_{r=l}^{k-1} \left(e^{-\alpha\eta r} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \quad (85)$$

1143

1144

$$1145 \leq e^{-3\alpha\eta/2} \left(1 + \frac{32\eta^3 L^4}{\alpha}\right) \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) \quad (86)$$

1146

1147

$$1148 + e^{-3\alpha\eta/2} \frac{192m\eta^3 L^4}{\alpha} e^{m\alpha\eta} \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \quad (87)$$

1149

1150

$$1151 \leq e^{-3\alpha\eta/2} \left(1 + \frac{32\eta^3 L^4}{\alpha}\right) \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) \quad (88)$$

1152

1153

$$1154 + e^{-3\alpha\eta/2} \frac{96m\eta^3 L^4}{\alpha} \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \quad (89)$$

1155

1156

$$1157 \leq e^{-3\alpha\eta/2} \left(1 + \frac{128\eta^3 L^4}{\alpha}\right) \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \quad (90)$$

1158

1159

$$1160 \leq e^{-\alpha\eta} \left(e^{-\alpha\eta(k-1)} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})}\right) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2} \quad (91)$$

1161

1162

$$1163 \leq e^{-\alpha\eta k} \text{KL}(\mu_0 \parallel \pi) + \frac{48m\eta^2 dL^2 + 16\eta^2 dL^2 b^2}{b^2(1 - e^{-\alpha\eta})} \quad (92)$$

1164

1165

$$1166 \leq e^{-\alpha\eta k} \text{KL}(\mu_0 \parallel \pi) + \frac{64m\eta dL^2 + 24\eta dL^2 b^2}{\alpha b^2}. \quad (93)$$

1167

1168 The first two inequalities are due to Eq. (78). The third and fourth inequality follow from the fact
1169 that $k - l \leq m$ and $e^{m\alpha\eta} \leq e^{\frac{\alpha^2}{8L^2}} \leq e^{\frac{1}{8}} \leq \frac{1}{2}$ for $\eta \leq \frac{\alpha}{8mL^2}$ and the fifth inequality holds since
1170 $\left(1 + \frac{128\eta^3 L^4}{\alpha}\right) \leq \left(1 + \frac{\eta\alpha}{2}\right) \leq e^{\alpha\eta/2}$ for $\eta \leq \frac{\alpha}{24L^2 m}$. The final inequality follows from the fact that
1171 $1 - e^{-\alpha\eta} \geq \frac{3}{4}\alpha\eta$ when $\alpha\eta \leq \frac{1}{4}$. This concludes the proof. \square

1172 **Theorem 4.9** (Main Theorem for QSVRG–LMC). *Let μ_k be the distribution of \mathbf{x}_k in QSVRG–LMC
1173 algorithm. Suppose that f satisfies Assumptions 4.2 and 4.6. Then for $\eta = \mathcal{O}\left(\frac{\epsilon\alpha}{dL^2} \wedge \frac{\alpha}{L^2 m}\right)$, $K = \tilde{\mathcal{O}}\left(\frac{L^2 \log(\text{KL}(\mu_0 \parallel \pi))}{\alpha^2} (n^{2/3} + \frac{d}{\epsilon})\right)$, $b = \tilde{\mathcal{O}}(n^{1/3})$, and $m = \tilde{\mathcal{O}}(n^{2/3})$ we have*

1174

$$1175 \left\{ \text{KL}(\mu_K \parallel \pi), \text{TV}(\mu_K, \pi)^2, \frac{\alpha}{2} \text{W}_2(\mu_K, \pi)^2 \right\} \leq \epsilon.$$

1176

1177 The total query complexity to the stochastic gradient oracle is
1178 $\tilde{\mathcal{O}}\left(\frac{L^2 \log(\text{KL}(\mu_0 \parallel \pi))}{\alpha^2} \left(nd^{1/2} + \frac{d^{3/2} n^{1/3}}{\epsilon}\right)\right)$.

1179 **Proof.** Setting $b = \tilde{\mathcal{O}}(n^{1/3})$ and $m = \tilde{\mathcal{O}}(n^{2/3})$ and $\eta \leq \frac{\epsilon\alpha}{176dL^2}$ the second term on the right
1180 hand side of Theorem 4.8 becomes smaller than $\epsilon/2$. By the step size requirement of Theorem 4.8,

1188 we have $\eta \leq \frac{\epsilon\alpha}{176dL^2} \wedge \frac{\alpha}{24L^2m}$. The first term in Theorem 4.8 is smaller than $\epsilon/2$ when $K \leq$
 1189 $\frac{\alpha\eta}{\log(2\text{KL}(\mu_0\|\pi)/\epsilon)}$. Hence TV distance is smaller than ϵ . The results for W_2 distance and TV distance
 1190 hold due to Talagrand’s inequality Otto & Villani (2000) and Pinsker’s inequality Tsybakov (2009)
 1191 respectively. The total gradient complexity is $bK = \tilde{\mathcal{O}}\left(\frac{L^2\text{KL}(\mu_0\|\pi)}{\alpha^2}\left(nd^{1/2} + \frac{d^{3/2}n^{1/3}}{\epsilon}\right)\right)$. \square
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241