

# COMPLEXITY OF ZIGZAG SAMPLING ALGORITHM FOR STRONGLY LOG-CONCAVE DISTRIBUTIONS

JIANFENG LU AND LIHAN WANG

**ABSTRACT.** We study the computational complexity of zigzag sampling algorithm for strongly log-concave distributions. The zigzag process has the advantage of not requiring time discretization for implementation, and that each proposed bouncing event requires only one evaluation of partial derivative of the potential, while its convergence rate is dimension independent. Using these properties, we prove that the zigzag sampling algorithm achieves  $\varepsilon$  error in chi-square divergence with a computational cost equivalent to  $O(\kappa^2 d^{\frac{1}{2}} (\log \frac{1}{\varepsilon})^{\frac{3}{2}})$  gradient evaluations in the regime  $\kappa \ll \frac{d}{\log d}$  under a warm start assumption, where  $\kappa$  is the condition number and  $d$  is the dimension.

## 1. INTRODUCTION AND MAIN RESULTS

Monte Carlo sampling from a high-dimensional probability distribution is a fundamental problem with applications in various areas including Bayesian statistics, machine learning, and statistical physics. Many sampling algorithms, especially those for continuous state space like  $\mathbb{R}^d$ , are based on continuous time Markov processes. Examples of these processes include the overdamped Langevin dynamics, whose invariant measure is the target measure, the underdamped Langevin dynamics and Hamiltonian Monte Carlo (HMC) [25], both augment the state space with a velocity variable  $v$ , and have the  $x$ -marginal distribution of the invariant measure as the target measure. For strongly log-concave distributions, all these processes converge to the equilibrium exponentially fast with rates independent of the dimension, making them suitable for sampling purposes. On the other hand, all of these processes require time discretizations for implementation, which not only induces further numerical errors but requires the time step to be small as well, requiring higher computational complexity if a small bias is desired. To remove such bias due to discretization, the conventional procedure is to introduce the Metropolis-Hastings acceptance-rejection step, while rejections indicate waste of computational resources.

A very different line of sampling algorithms have been recently developed in statistical physics and statistics literature [45], which are based on piecewise deterministic Markov processes (PDMPs) [20]. These processes are non-reversible, which may mix faster than reversible MCMC methods [21, 47]. Examples of such samplers include the randomized Hamiltonian Monte Carlo [11], the zigzag process [3], the bouncy particle sampler [12, 45], and some others [4, 42, 48]. The zigzag and bouncy particle samplers are appealing for big data applications, as they can be unbiased even if stochastic gradient is used [3, 12]. These algorithms, as they are still relatively new, have not yet been thoroughly analyzed. In particular, no non-asymptotic computational complexity bounds on these algorithms have been established yet, to the best of our knowledge. Our previous work [38] gives explicit exponential convergence rates for the PDMPs with log-concave potentials, which opens the possibility of deriving such complexity bounds for PDMPs, and provides the foundation of this work.

---

*Date:* December 20, 2020.

**1.1. Algorithm and Assumptions.** Let  $x$  denote the state variable in  $\mathbb{R}^d$  where  $d$  is the dimension. The target distribution we want to sample from is denoted by

$$d\mu(x) = Z^{-1} \exp(-U(x)) dx,$$

where  $U(x)$  is the potential and  $Z = \int_{\mathbb{R}^d} \exp(-U(x)) dx$  is the normalizing constant. Although the zigzag process can also be applied to sample non log-concave distributions, we will restrict our analysis to strongly log-concave distributions, namely, we make the following assumption throughout:

**Assumption 1.** *The potential function  $U(x)$  satisfies*

$$(1) \quad mI_d \leq \nabla^2 U(x) \leq LI_d,$$

where  $I_d$  denotes  $d \times d$  identity matrix. Moreover,  $U(x)$  has a unique minimizer at  $x = 0$ , and  $U(0) = 0$ .

For any random variable  $X$ , we use  $\rho(X)$  to denote its law. In this paper, we use chi-square divergence to measure the difference between two probability measures: for probability measures  $\rho_1, \rho_2$  that  $\rho_1 \ll \rho_2$ , it is defined as

$$\chi^2(\rho_1 \parallel \rho_2) := \int_{\mathbb{R}^d} \left( \frac{d\rho_1}{d\rho_2} - 1 \right)^2 d\rho_2.$$

The zigzag sampling algorithm is based on a piecewise deterministic Markov process, called zigzag process. Besides the variable  $x$ , we augment the state space by an auxiliary velocity variable  $v \in \mathbb{R}^d$ . A trajectory of the zigzag process, denoted by  $(X_t, V_t)$ , can be described as follows. Given some initial  $(X_0, V_0)$ ,  $V_t$  stays unchanged unless bouncing events or refreshing events occur at some random time following Poisson clocks. At any bouncing event on the  $j$ -th direction with rate  $(V_t^{(j)} \partial_{x_j} U(X_t))_+$ , the velocity  $V_t$  changes by flipping its  $j$ -th component. At the refreshing event with rate  $\lambda$  for some fixed  $\lambda > 0$ , the velocity  $V_t$  is completely redrawn from standard normal  $\mathcal{N}(0, \text{Id})$ . The position  $X_t$  always evolves according to the velocity as  $\frac{d}{dt} X_t = V_t$ .

It has been established [1, 8, 38] that under Assumption 1,  $\rho(X_t, V_t)$  converges to the invariant measure of the zigzag process, which is a product measure of the target measure in  $x$  and the standard Gaussian in  $v$ :

$$d\bar{\mu}(x, v) = d\mu(x) d\kappa(v) \quad \text{where} \quad d\kappa(v) = (2\pi)^{-\frac{d}{2}} \exp\left(-\frac{|v|^2}{2}\right) dv.$$

Our analysis relies on the following more quantitative convergence result for zigzag process proved in [38], which also specifies the optimal choice of refreshing rate  $\lambda$ :

**Proposition 1.1.** [38, Theorem 1] *Under Assumption 1, there exists a universal constant  $K$  independent of all parameters, such that for any initial density  $\bar{\mu}_0$ , the zigzag process with friction parameter  $\lambda = \sqrt{L}$  satisfies*

$$(2) \quad \chi^2(\rho(X_T, V_T) \parallel \bar{\mu}) \leq K \exp\left(-\frac{m}{K\sqrt{L}}T\right) \chi^2(\bar{\mu}_0 \parallel \bar{\mu}).$$

The left-hand side of (2) controls desired divergence of  $\rho(X)$  with respect to the target measure  $\mu$ , as we have

$$\begin{aligned} \chi^2(\rho(X_T, V_T) \parallel \bar{\mu}) &= \int_{\mathbb{R}^d \times \mathbb{R}^d} \left( \frac{d\rho(X_T, V_T)}{d\bar{\mu}} \right)^2 d\bar{\mu}(x, v) - 1 \\ &= \int_{\mathbb{R}^d} \left( \frac{d\rho(X_T)}{d\mu} \right)^2 \left( \int_{\mathbb{R}^d} \left( \frac{d\rho(V_T | X_T)}{d\kappa(v)} \right)^2 d\kappa(v) \right) d\mu(x) - 1 \\ &= \int_{\mathbb{R}^d} \left( \frac{d\rho(X_T)}{d\mu} \right)^2 \left( 1 + \chi^2(\rho(V_T | X_T) \parallel \kappa) \right) d\mu(x) - 1 \end{aligned}$$

$$\geq \int_{\mathbb{R}^d} \left( \frac{d\rho(X_T)}{d\mu} \right)^2 d\mu(x) - 1 = \chi^2(\rho(X_T) \parallel \mu).$$

Moreover, we would take initial condition in the form of

$$(3) \quad (X_0, V_0) \sim \bar{\mu}_0(x, v) = \mu_0(x)\kappa(v),$$

which implies that  $\chi^2(\bar{\mu}_0 \parallel \bar{\mu}) = \chi^2(\mu_0 \parallel \mu)$ . Therefore, we get

$$(4) \quad \chi^2(\rho(X_T) \parallel \mu) \leq K \exp\left(-\frac{m}{K\sqrt{L}}T\right)\chi^2(\mu_0 \parallel \mu),$$

which suggests the total time  $T$  needed to achieve control of chi-square divergence.

Of course, in reality, we cannot simulate the zigzag process precisely in a direct way. To turn the zigzag process into an actual sampling algorithm, we need an upper bound estimate of the exact bouncing rate, and simulate the Poisson process for the bouncing events using Poisson thinning (see e.g., discussions in [3, Section 3]). Under Assumption 1, we will use the following upper bound estimate for the rate:

$$(5) \quad (v_i \partial_{x_i} U(x + vt))_+ \leq |v_i \partial_{x_i} U(x + vt)| \leq |v_i| |\partial_{x_i} U(x + vt)| \leq L|v_i|(|x| + t|v|).$$

This upper bound has the advantage of not involving evaluations of  $U$  and its partial derivatives, which greatly reduces the computational cost. The price to pay is the increased frequency of potential bouncing events, which scales like  $O(\sqrt{d})$  since the pessimistic bound for the partial derivative  $|\partial_{x_i} U(x)| \leq |\nabla U(x)| \leq L|x|$  typically sacrifices a factor of  $O(\sqrt{d})$  in the first inequality.

Following the above discussions, the zigzag sampling algorithm is described in Algorithm 1, where Step 12 uses the upper bound estimate in (5), while Steps 19–23 correspond to the Poisson thinning step. Note that for each potential bouncing event, the algorithm requires one evaluation of  $\partial_{x_i} U$  in Step 19. In practice, typically accessing the partial derivatives of  $U$  is the most time consuming step, therefore, in our complexity analysis, we focus on the number of access to partial derivatives.

For the zigzag sampling algorithm to be efficient, the total runtime  $T$  in the algorithm has to be chosen appropriately. It needs to be long enough so that we can reach  $\varepsilon$  accuracy according to the (4). However,  $T$  cannot not be too long for several reasons: 1) If  $T$  is too large, we will have too many velocity refreshing events, among one of those events, a very large  $V$  might be drawn from the Gaussian distribution, which will lead to a high number of bouncing events and hence computational cost; and 2) since the zigzag trajectory is ergodic with respect to  $\mu$ , eventually it will visit regions with large gradient, which also leads to extra bouncing events. In principle these might be dealt with by revising the algorithm, for example restricting  $X_t$  to a prescribed region of the state space. However, this would introduce some bias. Therefore, we focus on in this work analyzing the original zigzag sampler, and as a result, to limit the runtime  $T$ , we need to impose the following Assumption 2 on the initial distribution  $\mu_0$ . We remark that due to the same reason, we would also put a restriction on the accuracy level  $\varepsilon$  that  $\frac{1}{\varepsilon} = o(\exp(d))$  in our complexity bound.

**Assumption 2.** *The initial distribution  $\mu_0(x)$  satisfies a warm-start condition:*

$$(6) \quad \chi^2(\mu_0 \parallel \mu) \leq \exp\left(\frac{d}{4K\kappa \log d}\right),$$

where  $K$  is the same universal constant as in (2). Furthermore, the initial distribution is concentrated in the sense of

$$(7) \quad \eta := \mathbb{P}_{\mu_0}\left(|x| > \sqrt{\frac{2d}{m}}\right) < \frac{1}{4}.$$

**Remark 1.2.** *The concentration condition (7) can be easily satisfied. By Gaussian Annulus Theorem, if we pick  $\mu_0 = \mathcal{N}(0, \frac{1}{m}\text{Id})$ , then  $\mathbb{P}_{\mu_0}\left(|x| > \sqrt{\frac{2d}{m}}\right) \leq 3e^{-cd}$  for some universal constant*

**Algorithm 1** The zigzag sampling algorithm

---

**Input:** Terminal runtime  $T$ , initial distribution  $\mu_0$ .

- 1: Draw  $x \sim \mu_0$ .
- 2: Set  $t \leftarrow 0$ .
- 3: Set  $\text{refr} \leftarrow \text{true}$ .
- 4: **while**  $t < T$  **do**
- 5:   **if**  $\text{refr}$  **then**
- 6:     Draw  $v \sim \mathcal{N}(0, I_d)$ .
- 7:     Draw  $t_{\text{refr}} \sim \text{Exp}(\sqrt{L})$ .
- 8:      $t_{\text{refr}} \leftarrow \min\{t_{\text{refr}}, T - t\}$ .
- 9:      $\text{refr} \leftarrow \text{false}$ .
- 10:   **end if**
- 11:   **for**  $i = 1, \dots, d$  **do**
- 12:     Draw  $\tau_i$  such that  $\mathbb{P}(\tau_i \geq s) = \exp(-sL|v_i||x| - \frac{s^2}{2}|v_i||v|)$ .
- 13:   **end for**
- 14:   Pick  $j = \arg \min_{i=1, \dots, d} \tau_i$ .
- 15:    $\Lambda_j \leftarrow L|v_j|(|x| + \tau_j|v|)$ .
- 16:    $t \leftarrow t + \min\{\tau_j, t_{\text{refr}}\}$ .
- 17:    $x \leftarrow x + v \min\{\tau_j, t_{\text{refr}}\}$ .
- 18:   **if**  $\tau_j < t_{\text{refr}}$  **then**
- 19:      $\lambda_j \leftarrow (v_j \partial_{x_j} U(x))_+$ .
- 20:     Draw  $\alpha \sim \text{Unif}(0, 1)$ ;
- 21:     **if**  $\alpha < \frac{\lambda_j}{\Lambda_j}$  **then**
- 22:        $v_j \leftarrow -v_j$ .
- 23:     **end if**
- 24:      $t_{\text{refr}} \leftarrow t_{\text{refr}} - \tau_j$ .
- 25:   **else**
- 26:      $\text{refr} \leftarrow \text{true}$ .
- 27:   **end if**
- 28: **end while**
- 29: **return**  $x$ .

---

*c. The failure probability gets smaller if we take  $\mu_0 = \mathcal{N}(0, \frac{1}{L}\text{Id})$ . The warm start condition (6) is much more stringent though, which we will discuss further in Section 3.*

## 1.2. Main Results.

**Theorem 1.** *Under Assumptions 1 and 2, for any prescribed accuracy  $\varepsilon > 0$ , Algorithm 1 outputs a random variable  $X$  such that*

$$(8) \quad \chi^2(\rho(X) \parallel \mu) \leq \varepsilon,$$

for runtime  $T$  chosen as

$$(9) \quad T = K \left( \frac{\sqrt{L}}{m} \left( \log \frac{1}{\varepsilon} + \log \chi^2(\mu_0 \parallel \mu) + \log K \right) \right),$$

where  $K$  is the universal constant in (2).

Moreover, if  $\varepsilon \geq \exp(-\frac{d}{4K\kappa \log d})$ , with probability  $1 - \frac{C}{\sqrt{LT}} - C \log^{-\frac{3}{2}} d - \eta$ , Algorithm 1 returns an output with a computational cost of

$$O\left(d^{\frac{3}{2}} \kappa^2 \left( \log^{\frac{3}{2}} \frac{1}{\varepsilon} + \log^{\frac{3}{2}} \chi^2(\mu_0 \parallel \mu) \right)\right)$$

evaluations of partial derivatives of  $U$ , where  $\eta$  is defined in (7) and  $C$  is a universal constant.

**Remark 1.3.** *By repeated trials, the theorem implies that for any  $\delta \in (0, \frac{1}{4})$ , with probability  $1 - \delta$ , Algorithm 1 returns the desired output with a computational cost of*

$$O\left(d^{\frac{3}{2}}\kappa^2\left(\log^{\frac{3}{2}}\frac{1}{\varepsilon} + \log^{\frac{3}{2}}\chi^2(\mu_0 \parallel \mu)\right)\log\frac{1}{\delta}\left|\log^{-1}\left(\frac{1}{\sqrt{LT}} + \log^{-\frac{3}{2}}d + \eta\right)\right|\right),$$

that is  $\tilde{O}(d^{\frac{3}{2}}\kappa^2)$  evaluations of partial derivatives of  $U$ , where  $\tilde{O}(\cdot)$  hides logarithmic factors.

With the common computational model that  $d$  evaluations of partial derivatives of  $U$  is equivalent to one evaluation of  $\nabla U$  in complexity, the complexity of zigzag is equivalent to  $\tilde{O}(d^{\frac{1}{2}}\kappa^2)$  evaluations of  $\nabla U$ .

Theorem 1 guarantees that the zigzag sampling algorithm (Algorithm 1) outputs a sample from a distribution with  $\chi^2$ -divergence at most  $\varepsilon$  away from the target density for a computational complexity equivalent to  $\tilde{O}(d^{\frac{3}{2}}\kappa^2)$  partial derivative evaluations (i.e., amounts to  $\tilde{O}(d^{\frac{1}{2}}\kappa^2)$  gradient evaluations), in the regime  $\max\{\kappa, \log\frac{1}{\varepsilon}\} \ll \frac{d}{\log d}$  with a warm-start condition.

Our analysis is based on the quantitative convergence rate of the zigzag process established in [38], which is  $O(\frac{m}{\sqrt{L}})$  for  $m$ -convex and  $L$ -smooth potentials. The rest of our proof is based on estimating  $\sup|X_t|$  along a single trajectory of the zigzag process and subsequently turn this into an estimate on the number of potential bouncing events, and hence number of partial derivative evaluations. Our analysis utilizes the two important and desirable features of the zigzag sampling process:

- The implementation of the zigzag process does not need time discretization, as the velocity in deterministic portion of the trajectory remains constant, which makes it possible to simulate the exact trajectories of the zigzag process while eliminating an important source of error. This is the reason that the complexity of the zigzag process only has logarithmic dependence on  $\frac{1}{\varepsilon}$ , without Metropolis acceptance/rejection.
- Moreover, for each potential bouncing event of zigzag, only one evaluation of a *partial derivative* of the potential is required, which is  $O(d)$  cheaper than a full gradient evaluation in computational cost for usual model of computation.

We would also remark that we quantify the error of distribution in terms of  $\chi^2$ -divergence, which provides stronger guarantee than total variation, KL divergence or 2-Wasserstein distance. While  $\chi^2$ -divergence is relatively convenient for obtaining convergence rates of continuous processes [14, 38] based on Poincaré inequality, it does not seem easy to use for analyzing discretization error of SDEs. Previous works made assumptions of Poincaré inequality for the discrete invariant measure [49] or some opaque uniform warmness assumption [30], both of which are difficult to verify. We are fortunate to avoid such problem for zigzag sampler, thanks to the fact that zigzag does not need time discretization.

**1.3. Previous Works.** Here we focus on results on non-asymptotic analysis of sampling algorithms, which has been a focused research area in recent years. Many sampling algorithms have been analyzed including algorithms based on overdamped Langevin dynamics [18, 22, 26, 27, 37, 49], underdamped Langevin dynamics [17, 19, 23, 39, 43, 46], Hamiltonian Monte Carlo [10, 16, 35, 40, 41], or high order Langevin dynamics [44], among others. These methods involve discretization of ODEs or SDEs, which yields an error that scales polynomially with step size. Thus the complexity of these algorithms has polynomial dependence on  $\varepsilon^{-1}$ , where  $\varepsilon$  is the desired accuracy threshold.

Metropolized variants of sampling algorithms, including Metropolized HMC and Metropolis Adjusted Langevin Algorithm (MALA), have also been studied in [15, 29, 34], the complexities of which have only logarithmic dependence on  $\varepsilon^{-1}$ , similar to the zigzag sampling process analyzed here. In [29] the complexity upper bound for MALA is established as  $\tilde{O}(\kappa d + \kappa^{\frac{3}{2}}d^{\frac{1}{2}})$  under warm start condition, and  $\tilde{O}(\kappa d^2 + \kappa^{\frac{3}{2}}d^{\frac{3}{2}})$  with a feasible start. In [15] the complexity upper bound for MALA is improved to  $\tilde{O}(\kappa d + \kappa^{\frac{3}{2}}d^{\frac{1}{2}})$  with feasible start (where  $\mu_0 = \mathcal{N}(0, \frac{1}{L}\text{Id})$ ). The work

[15] also established bounds for Metropolized HMC, which is  $\tilde{O}(\kappa d^{\frac{11}{12}})$  with warm start (which is in fact more stringent than our Assumption 2) in the regime  $\kappa = O(d^{\frac{2}{3}})$ , and  $\tilde{O}(\kappa^{\frac{3}{4}}d + \kappa^{\frac{7}{4}}d^{\frac{1}{2}})$  with feasible start if the target potential function has a bounded Hessian. The complexity upper bound has been improved in [34] to  $\tilde{O}(\kappa d)$  for both Metropolized HMC and MALA with a feasible start, based on a refined analysis using concentration of gradient norm. In comparison, our result for zigzag relies on a warm start, while the complexity upper bound has better dependence in  $d$ . The issue of feasible start will be further discussed in Section 3.

Regarding asymptotic analysis for the convergence of zigzag process, the ergodicity was first established in [8]. Exponential convergence of the zigzag process is established in [7, 31] using a Lyapunov function argument. A central limit theorem of the zigzag process is established in [2], and a large deviation principle is established for the empirical measure in [6]. The spectrum of the zigzag process has been studied in [5, 32]. A dimension independent exponential convergence rate for the zigzag process is established in [1], using the hypocoercivity framework developed in [24]. Finally, a more quantitative convergence estimate was established in [38], for which our analysis of the sampling algorithm is based on.

## 2. STRATEGY OF THE PROOF

Since Algorithm 1 always simulates exact trajectories of the zigzag process, we see that (8) is guaranteed with the correct choice of  $T$ . Therefore we only need to estimate the computational complexity. The strategy of the proof is to first give an estimate on  $\sup_{t \in [0, T]} U(X_t)$  (Lemma 2.1), which directly controls  $\sup_{t \in [0, T]} |X_t|$ . The upper bound on  $|X_t|$  in turn provides us an estimate of upper bound on the number of partial derivative evaluations of  $U$ . The complexity upper bound we derive holds with high probability, while it does not always hold (for example, the number of proposed bouncing events from the Poisson clock might be atypically high), such events only occur with very small probability, which will be controlled in the proof.

Let  $N + 1$  be the total number of velocity refreshments (including the initial refreshment), therefore  $N$  is a Poisson random variable such that

$$(10) \quad \mathbb{P}(N = n) = \frac{(\sqrt{LT})^n}{n!} e^{-\sqrt{LT}}.$$

Let  $0 = T_0 < T_1 < T_2 < \dots < T_N \leq T < T_{N+1}$  be the refresh times, and  $V_{T_k}$  be the velocity variable after refreshment at time  $T_k$ . For  $k = 1, \dots, N$ , we use  $t_k = T_k - T_{k-1}$  to denote the time duration between refreshments. For convenience, we will also denote  $t_{N+1} = T - T_N$ .

The first step of the proof is the following lemma which controls  $\sup_{t \in [0, T]} U(X_t)$  condition on some high probability events. The proof will be deferred to the appendix.

**Lemma 2.1.** *Under Assumptions 1 and 2, suppose the following conditions hold:*

$$(11a) \quad \frac{1}{2}\sqrt{LT} \leq N \leq \frac{3}{2}\sqrt{LT};$$

$$(11b) \quad |V_{T_k} \cdot \nabla U(X_{T_k})| \leq \left(\frac{d}{\sqrt{LT}}\right)^{1/2} |\nabla U(X_{T_k})|, \quad \forall k = 1, \dots, N;$$

$$(11c) \quad |V_{T_k}| \leq 2\sqrt{d}, \quad \forall k = 1, \dots, N$$

$$(11d) \quad U(X_0) \leq \kappa d;$$

$$(11e) \quad \sum_{k=1}^{N+1} t_k^2 \leq \frac{4T}{\sqrt{L}}.$$

Then there exists a universal constant  $C$  such that

$$(12) \quad \sup_{t \in [0, T]} U(X_t) \leq C\sqrt{LT}d.$$

The next element in the proof is to control the failure event that (11) does not hold. The control of the first four events are relatively straightforward and will thus be directly carried out in the proof of theorem below; we state the probability for the event (11e) to hold as the following lemma, which will also be proved in the appendix.

**Lemma 2.2.** *With probability  $1 - \frac{2}{\sqrt{LT}}$ , condition (11e) holds.*

The final component of the proof is to turn the estimate for  $\sup_{t \in [0, T]} U(X_t)$  to an upper bound for the number of proposed bouncing events. We will use the following large deviation result proved in [33] (see also [36, Theorem 1] for a generalization):

**Lemma 2.3.** *Let  $\{Y_n\}$  be a sequence of real-valued i.i.d. random variables. Suppose  $EY_1 = \mu > 0$  and  $M(\lambda_0) := \mathbb{E}(\exp(\lambda_0 Y_1)) < \infty$  for some  $\lambda_0 > 0$ . Let  $N_t = \inf_n \{\sum_{k=1}^n Y_k > t\}$  and  $P_t$  be the law of  $\frac{1}{t}N_t$ . Then  $(P_t)_{t \geq 0}$  satisfies the large deviation principle with rate function*

$$J(x) = \begin{cases} \sup_{\theta \in \Theta} (\theta - x \log M(\theta)) & x \geq 0, \\ \infty & x < 0. \end{cases}$$

where  $\Theta = \{\theta : M(\theta) < \infty\}$ .

*Proof of Theorem 1.* Let  $p_i$  be the probability that condition  $i$  in (11) of Lemma 2.1 fails. We start with condition (11a) of Lemma 2.1. Applying Lemma 2.3 to the Poisson process with  $t_i$  as the arrival times, we may estimate the first failure probability (here and for the rest of the proofs  $C$  denotes a universal constant that may change from line to line)

$$(13) \quad p_a \leq \exp\left(-\frac{1}{C}\sqrt{LT}\right) \leq \frac{C}{\sqrt{LT}}.$$

We now check the conditions (11b) and (11c) of Lemma 2.1. By Gaussian Annulus Theorem, for each refreshment, we have

$$(14) \quad \mathbb{P}(|V| \geq 2\sqrt{d}) \leq 3e^{-cd},$$

where  $c > 0$  is some universal constant. We also require  $V_{T_k}$  to satisfy  $|V_{T_k} \cdot n(X_{T_k})| \leq \left(\frac{d}{\sqrt{LT}}\right)^{1/2}$ , where  $n(X_{T_k}) = \frac{\nabla U(X_{T_k})}{|\nabla U(X_{T_k})|}$ , which has failure probability

$$\begin{aligned} \mathbb{P}(|V \cdot n(x)| \geq \left(\frac{d}{\sqrt{LT}}\right)^{1/2}) &= \frac{1}{\sqrt{2\pi}} \int_{\left(\frac{d}{\sqrt{LT}}\right)^{1/2}}^{\infty} \exp\left(-\frac{r^2}{2}\right) dr \\ &\leq \frac{1}{\sqrt{2\pi}} \int_{\left(\frac{d}{\sqrt{LT}}\right)^{1/2}}^{\infty} \exp\left(-\frac{r}{2}\left(\frac{d}{\sqrt{LT}}\right)^{1/2}\right) dr \\ &\leq \sqrt{\frac{2}{\pi}} \left(\frac{\sqrt{LT}}{d}\right)^{1/2} \exp\left(-\frac{d}{2\sqrt{LT}}\right). \end{aligned}$$

Since we have to draw  $V$  for  $N$  times, cumulatively this yields a failure probability

$$(15) \quad p_b + p_c \leq C \left( e^{-cd} + \left(\frac{\sqrt{LT}}{d}\right)^{1/2} \exp\left(-\frac{d}{2\sqrt{LT}}\right) \right) \mathbb{E}N.$$

Recall the assumption  $\varepsilon \geq \exp\left(-\frac{d}{4K\kappa \log d}\right)$  as well as (6), which implies that  $\sqrt{LT} \leq \frac{d}{2 \log d}$  for our choice of  $T$  as in (9). Together with condition (11a), we derive (neglecting the obviously smaller term  $e^{-cd}$ )

$$p_b + p_c \leq C\sqrt{LT} \left(\frac{\sqrt{LT}}{d}\right)^{1/2} \exp\left(-\frac{d}{2\sqrt{LT}}\right) \leq C \log^{-\frac{3}{2}} d.$$

The failure probability for condition (11d) is straightforward to estimate. Using Assumption 1, we have

$$U(x) \leq \frac{L}{2}|x|^2,$$



which indicates

$$p_d \leq \eta = \mathbb{P}(|X_0| \geq \sqrt{\frac{2d}{m}}).$$

Finally,  $p_e$  is already estimated in Lemma 2.2, which yields  $p_e \leq \frac{2}{\sqrt{LT}}$ .

In summary, the total failure probability of (11) can be bounded as

$$(16) \quad p_a + p_b + p_c + p_d + p_e \leq \frac{C}{\sqrt{LT}} + C \log^{-\frac{3}{2}} d + \eta.$$

We now assume that condition (11) holds. Thus, Lemma 2.1 together with Assumption 1 implies that

$$(17) \quad \sup_{t \in [0, T]} |X_t| \leq \left( \frac{2}{m} \sup_{t \in [0, T]} U(X_t) \right)^{1/2} \leq C \left( \frac{\sqrt{L}}{m} T d \right)^{1/2}.$$

After each refreshment or bouncing event, Algorithm 1 runs  $d$  independent Poisson clocks  $\{\tau_i\}_{i=1, \dots, d}$  defined in Step 12, thus, noticing  $\sum_i |v_i| \leq \sqrt{d}|v| \leq 2d$ ,

$$(18) \quad \mathbb{P}(\min \tau_i \geq t) \geq \exp\left(-tL|x| \sum_i |v_i| - \frac{t^2}{2}|v| \sum_i |v_i|\right) \geq \exp(-Cd^{\frac{3}{2}}(L^{\frac{5}{4}}m^{-\frac{1}{2}}T^{\frac{1}{2}}t + t^2)).$$

This motivates us to consider the following counting process  $\tilde{N}_t$ : suppose  $\tilde{t}_1, \dots$  are i.i.d. random variables with  $\mathbb{P}(\tilde{t}_i \geq s) = \exp(-As - Bs^2)$  where  $A = Cd^{\frac{3}{2}}L^{\frac{5}{4}}m^{-\frac{1}{2}}T^{\frac{1}{2}}$  and  $B = Cd^{\frac{3}{2}}$ , and let  $\tilde{N}_t = \inf_n \{\sum_{i=1}^n \tilde{t}_i > t\}$ . By construction, the probability of  $N > 2AT$  under condition (11) is controlled by  $\mathbb{P}(\tilde{N}_T > 2AT)$ . Therefore, it suffices to estimate  $\mathbb{P}(\tilde{N}_T > 2AT)$ , for which we use Lemma 2.3. To estimate the rate function, for  $\theta \in (-\infty, A)$ , we calculate

$$\begin{aligned} M(\theta) &= \int_0^\infty (A + 2Bs)e^{-(A-\theta)s - Bs^2} ds \\ &= A \int_0^\infty e^{-(A-\theta)s - Bs^2} ds - \int_0^\infty e^{-(A-\theta)s} de^{-Bs^2} \\ &= A \int_0^\infty e^{-(A-\theta)s - Bs^2} ds + 1 + \int_0^\infty e^{-Bs^2} de^{-(A-\theta)s} \\ &= 1 + \theta \int_0^\infty e^{-(A-\theta)s - Bs^2} ds \\ &\leq 1 + \theta \int_0^\infty e^{-(A-\theta)s} ds = \frac{A}{A - \theta}. \end{aligned}$$

This implies

$$\begin{aligned} J(x) &:= \sup_{\theta \in \mathbb{R}} (\theta - x \log M(\theta)) \geq \sup_{\theta \in (-\infty, A)} \left( \theta - x \log \frac{A}{A - \theta} \right) \\ &= x \log x - x \log A + A - x := \tilde{J}(x). \end{aligned}$$

It is easy to check that  $\tilde{J}(x)$  is increasing for any  $x \geq A$ , and

$$J(2A) \geq \tilde{J}(2A) = A(2 \log 2 - 1).$$

This means for some universal constant  $c > 0$ ,

$$\mathbb{P}(\tilde{N}_T \geq 2AT) \leq \exp(-cAT) = \exp(-cd^{\frac{3}{2}}L^{\frac{5}{4}}m^{-\frac{1}{2}}T^{\frac{3}{2}}),$$

which is much smaller than the previously estimated failure probability (16). In other words, we have established that with high probability the number of partial derivative evaluations is bounded by

$$O(AT) = O(d^{\frac{3}{2}}L^{\frac{5}{4}}m^{-\frac{1}{2}}T^{\frac{3}{2}}) = O\left(d^{\frac{3}{2}}\kappa^2\left(\log^{\frac{3}{2}}\frac{1}{\varepsilon} + \log^{\frac{3}{2}}\chi^2(\mu_0 \parallel \mu)\right)\right). \quad \square$$



## 3. DISCUSSIONS

We establish non-asymptotic complexity bounds for the zigzag sampling algorithm. While we focus on zigzag sampler in this work, we expect that similar analysis for other PDMPs [4, 12, 42, 48] can be carried out. We leave these for future research.

We admit that our warm-start requirement (6) is much more stringent than existing results (for example [34, 39]). We observe that (6) implicitly requires the condition number  $\kappa$  to be much smaller than  $d$ , as otherwise, if  $\kappa \sim d$ , (6) requires  $\chi^2(\mu_0 \parallel \mu) = O(1)$  which is unrealistic. This restriction on condition number is not completely unexpected since the zigzag sampler does perform poorly for highly anisotropic densities (see for example numerical results in [42]).

A major issue of the warm-start assumption comes from our choice of  $\chi^2$  divergence, rather than total variation, 2-Wasserstein distance, or KL divergence as in previous works for non-asymptotic analysis of sampling algorithms. In particular, if we choose the initial condition

$$(19) \quad d\mu_0(x) = \left(\frac{L}{2\pi}\right)^{\frac{d}{2}} \exp\left(-\frac{L|x|^2}{2}\right) dx,$$

as in previous works, then for  $U(x) = \frac{m|x|^2}{2}$ , we have

$$\begin{aligned} \chi^2(\mu_0 \parallel \mu) &= Z \left(\frac{L}{2\pi}\right)^d \int_{\mathbb{R}^d} \exp(-L|x|^2 + U(x)) dx - 1 \\ &= \kappa^{\frac{d}{2}} \left(\frac{L}{2\pi}\right)^{\frac{d}{2}} \int_{\mathbb{R}^d} \exp\left(-\left(L - \frac{m}{2}\right)|x|^2\right) dx - 1 = \kappa^{\frac{d}{2}} \left(\frac{L}{2L - m}\right)^{\frac{d}{2}} - 1, \end{aligned}$$

which violates (6). On the other hand, for the same choice of  $\mu_0$ , as long as  $U$  satisfies Assumption 1, one can estimate

$$\text{KL}(\mu_0 \parallel \mu) = \left(\frac{L}{2\pi}\right)^{\frac{d}{2}} \int_{\mathbb{R}^d} \left(\frac{d}{2} \log \frac{L}{2\pi} + \log Z - \frac{L}{2}|x|^2 + U(x)\right) \exp\left(-\frac{L}{2}|x|^2\right) dx \leq \frac{d}{2} \log \kappa.$$

This means  $\log \text{KL}(\mu_0 \parallel \mu)$ , and consequently the logarithm of total variation or 2-Wasserstein distances are much smaller than any algebraic power of  $d$ , making it suitable for initialization. We hope the following conjecture is true:

**Conjecture 1.** *Under Assumption 1, there exists a universal constant  $K$  independent of all parameters, such that for any initial density  $\bar{\mu}_0$ , the zigzag process with friction parameter  $\lambda = \sqrt{L}$  satisfies*

$$\text{KL}(\rho(X_T, V_T) \parallel \bar{\mu}) \leq K \exp\left(-\frac{m}{K\sqrt{L}}T\right) \text{KL}(\bar{\mu}_0 \parallel \bar{\mu}).$$

If this is indeed true, we can establish the convergence in KL divergence of the zigzag sampler using a feasible start.

Another interesting open question is whether one can find a tighter upper bound than Step 12 of Algorithm 1 in order to reduce the computational complexity, since it magnifies proposed bouncing rate by  $O(\sqrt{d})$ . The following lemma, which might be of independent interest, provides a concentration bound for  $|\partial_{x_i} U|$  so that we might be able to give up a small probability to obtain a much sharper bouncing rate control.

**Lemma 3.1.** *Let  $U(x)$  satisfy Assumption 1, then for any  $c > 0$ ,*

$$(20) \quad \mathbb{P}_\mu \left( |\partial_{x_i} U| \geq 2\sqrt{L} + 2c\sqrt{L} \log d \right) \leq 3d^{-c}.$$

The proof of this lemma, deferred to the appendix, is inspired by [34], which uses the following Brascamp-Lieb inequality [13]:

**Lemma 3.2.** *Let  $U(x)$  satisfy Assumption 1, then for any  $g \in H^1(\mu)$ ,*

$$(21) \quad \text{Var}_\mu g \leq \int_{\mathbb{R}^d} \nabla g (\nabla^2 U)^{-1} \nabla g d\mu.$$

With Lemma 3.1, it might be possible to improve Algorithm 1 while surrendering a small probability by replacing Step 12 with  $\mathbb{P}(\tau_i \geq s) = \exp(-cs\sqrt{L}|v_i| \log d)$  since  $(v_i \partial_{x_i} U(x+vs))_+ \leq c\sqrt{L}|v_i| \log d$  with high probability. This motivates the following conjecture:

**Conjecture 2.** *Under the Assumption 1, for any  $\kappa$  and  $\log \frac{1}{\varepsilon}$  that are both smaller than some algebraic power of  $d$ , there exists an algorithm that gives a random variable  $X$  such that*

$$(22) \quad \chi^2(\rho(X) \parallel \mu) \leq \varepsilon.$$

*Moreover, with high probability, the algorithm requires  $O\left(d\kappa \log d(\log \frac{1}{\varepsilon} + \log \chi^2(\mu_0 \parallel \mu))\right)$  evaluations of partial derivatives of  $U$ .*

Unfortunately there are several difficulties for proving the conjecture. One is that although  $\partial_{x_i} U$  does not exceed  $O(\log d)$  with high probability, we are unable to control the partial derivatives for a trajectory of the zigzag process. Another issue is that since some trajectories of the zigzag process may go to regions with partial derivatives exceeding  $O(\log d)$ , we do not always simulate the exact trajectories, which introduces bias in the sampling.

**Acknowledgment.** This work is supported in part by National Science Foundation via grants CCF-1910571 and DMS-2012286.

#### REFERENCES

- [1] Christophe Andrieu, Alain Durmus, Nikolas Nüsken, and Julien Roussel, *Hypocoercivity of piecewise deterministic Markov process-Monte Carlo*, arXiv preprint arXiv:1808.08592 (2018).
- [2] Joris Bierkens and Andrew Duncan, *Limit theorems for the zig-zag process*, *Advances in Applied Probability* **49** (2017), no. 3, 791–825.
- [3] Joris Bierkens, Paul Fearnhead, and Gareth Roberts, *The zig-zag process and super-efficient sampling for Bayesian analysis of big data*, *The Annals of Statistics* **47** (2019), no. 3, 1288–1320.
- [4] Joris Bierkens, Sebastiano Grazi, Kengo Kamatani, and Gareth Roberts, *The boomerang sampler*, International conference on machine learning, 2020, pp. 908–918.
- [5] Joris Bierkens and Sjoerd M Verduyn Lunel, *Spectral analysis of the zigzag process*, arXiv preprint arXiv:1905.01691 (2019).
- [6] Joris Bierkens, Pierre Nyquist, and Mikola Schlottke, *Large deviations for the empirical measure of the zig-zag process*, arXiv preprint arXiv:1912.06635 (2019).
- [7] Joris Bierkens, Gareth Roberts, et al., *A piecewise deterministic scaling limit of lifted metropolis–hastings in the curie–weiss model*, *The Annals of Applied Probability* **27** (2017), no. 2, 846–882.
- [8] Joris Bierkens, Gareth O Roberts, and Pierre-André Zitt, *Ergodicity of the zigzag process*, *The Annals of Applied Probability* **29** (2019), no. 4, 2266–2301.
- [9] Sergey Bobkov and Michel Ledoux, *Poincaré’s inequalities and Talagrand’s concentration phenomenon for the exponential distribution*, *Probability Theory and Related Fields* **107** (1997), no. 3, 383–400.
- [10] Nawaf Bou-Rabee, Andreas Eberle, and Raphael Zimmer, *Coupling and convergence for Hamiltonian Monte Carlo*, *Annals of Applied Probability* **30** (2020), no. 3, 1209–1250.
- [11] Nawaf Bou-Rabee and Jesus Maria Sanz-Serna, *Randomized Hamiltonian Monte Carlo*, *The Annals of Applied Probability* **27** (2017), no. 4, 2159–2194.
- [12] Alexandre Bouchard-Côté, Sebastian J Vollmer, and Arnaud Doucet, *The bouncy particle sampler: A non-reversible rejection-free Markov chain Monte Carlo method*, *Journal of the American Statistical Association* **113** (2018), no. 522, 855–867.
- [13] Herm Jan Brascamp and Elliott H Lieb, *On extensions of the Brunn-Minkowski and Prékopa-Leindler theorems, including inequalities for log concave functions, and with an application to the diffusion equation*, *Journal of Functional Analysis* **22** (1976), no. 4, 366–389.
- [14] Yu Cao, Jianfeng Lu, and Lihan Wang, *On explicit  $L^2$ -convergence rate estimate for underdamped Langevin dynamics*, arXiv preprint arXiv:1908.04746 (2019).
- [15] Yuansi Chen, Raaz Dwivedi, Martin J Wainwright, and Bin Yu, *Fast mixing of Metropolized Hamiltonian Monte Carlo: Benefits of multi-step gradients*, *Journal of Machine Learning Research* **21** (2020), no. 92, 1–72.
- [16] Zongchen Chen and Santosh S Vempala, *Optimal convergence rate of Hamiltonian Monte Carlo for strongly logconcave distributions*, arXiv preprint arXiv:1905.02313 (2019).
- [17] Xiang Cheng, Niladri S Chatterji, Peter L Bartlett, and Michael I Jordan, *Underdamped Langevin MCMC: A non-asymptotic analysis*, Conference on learning theory, 2018, pp. 300–323.

- [18] Arnak S Dalalyan, *Theoretical guarantees for approximate sampling from smooth and log-concave densities*, Journal of the Royal Statistical Society: Series B (Statistical Methodology) **3** (2017), no. 79, 651–676.
- [19] Arnak S Dalalyan and Lionel Riou-Durand, *On sampling from a log-concave density using kinetic Langevin diffusions*, Bernoulli **26** (2020), no. 3, 1956–1988.
- [20] Mark HA Davis, *Piecewise-deterministic Markov processes: a general class of non-diffusion stochastic models*, Journal of the Royal Statistical Society: Series B (Methodological) **46** (1984), no. 3, 353–376.
- [21] Persi Diaconis, Susan Holmes, and Radford M Neal, *Analysis of a nonreversible Markov Chain sampler*, Annals of Applied Probability (2000), 726–752.
- [22] Zhiyan Ding, Qin Li, Jianfeng Lu, and Stephen J Wright, *Random coordinate Langevin Monte Carlo*, arXiv preprint arXiv:2010.01405 (2020).
- [23] ———, *Random coordinate underdamped Langevin Monte Carlo*, arXiv preprint arXiv:2010.11366 (2020).
- [24] Jean Dolbeault, Clément Mouhot, and Christian Schmeiser, *Hypocoercivity for linear kinetic equations conserving mass*, Transactions of the American Mathematical Society **367** (2015), no. 6, 3807–3828.
- [25] Simon Duane, Anthony D Kennedy, Brian J Pendleton, and Duncan Roweth, *Hybrid monte carlo*, Physics letters B **195** (1987), no. 2, 216–222.
- [26] Alain Durmus, Szymon Majewski, and Blazej Miasojedow, *Analysis of Langevin Monte Carlo via convex optimization.*, J. Mach. Learn. Res. **20** (2019), 73–1.
- [27] Alain Durmus and Eric Moulines, *High-dimensional bayesian inference via the unadjusted Langevin algorithm*, Bernoulli **25** (2019), no. 4A, 2854–2882.
- [28] Richard Durrett, *Essentials of stochastic processes*, Vol. 1, Springer, 1999.
- [29] Raaz Dwivedi, Yuansi Chen, Martin J Wainwright, and Bin Yu, *Log-concave sampling: Metropolis-Hastings algorithms are fast!*, Conference on learning theory, 2018, pp. 793–797.
- [30] Murat A Erdogdu and Rasa Hosseinzadeh, *Convergence analysis of Langevin Monte Carlo in chi-square divergence*, arXiv preprint arXiv:2007.11612 (2020).
- [31] Joaquin Fontbona, Hélène Guérin, and Florent Malrieu, *Long time behavior of telegraph processes under convex potentials*, Stochastic Processes and their Applications **126** (2016), no. 10, 3077–3101.
- [32] Arnaud Guillin and Boris Nectoux, *Low lying eigenvalues and convergence to the equilibrium of some piecewise deterministic markov processes generators in the small temperature regime* (2020).
- [33] Tiefeng Jiang, *Large deviations for renewal processes*, Stochastic processes and their applications **50** (1994), no. 1, 57–71.
- [34] Yin Tat Lee, Ruoqi Shen, and Kevin Tian, *Logsmooth gradient concentration and tighter runtimes for Metropolized Hamiltonian Monte Carlo*, arXiv preprint arXiv:2002.04121 (2020).
- [35] Yin Tat Lee, Zhao Song, and Santosh S Vempala, *Algorithmic theory of ODEs and sampling from well-conditioned logconcave densities*, arXiv preprint arXiv:1812.06243 (2018).
- [36] Raphaël Lefevere, Mauro Mariani, and Lorenzo Zambotti, *Large deviations for renewal processes*, Stochastic Processes and their Applications **121** (2011), no. 10, 2243–2271.
- [37] Xuechen Li, Yi Wu, Lester Mackey, and Murat A Erdogdu, *Stochastic Runge-Kutta accelerates Langevin Monte Carlo and beyond*, Advances in neural information processing systems, 2019, pp. 7748–7760.
- [38] Jianfeng Lu and Lihan Wang, *On explicit  $L^2$ -convergence rate estimate for piecewise deterministic Markov processes*, arXiv preprint arXiv:2007.14927 (2020).
- [39] Yi-An Ma, Niladri Chatterji, Xiang Cheng, Nicolas Flammarion, Peter Bartlett, and Michael I Jordan, *Is there an analog of Nesterov acceleration for MCMC?*, arXiv preprint arXiv:1902.00996 (2019).
- [40] Oren Mangoubi and Aaron Smith, *Mixing of Hamiltonian Monte Carlo on strongly log-concave distributions 2: Numerical integrators*, The 22nd international conference on artificial intelligence and statistics, 2019, pp. 586–595.
- [41] Oren Mangoubi and Nisheeth Vishnoi, *Dimensionally tight bounds for second-order Hamiltonian Monte Carlo*, Advances in neural information processing systems **31** (2018), 6027–6037.
- [42] Manon Michel, Sebastian C Kapfer, and Werner Krauth, *Generalized event-chain Monte Carlo: Constructing rejection-free global-balance algorithms from infinitesimal steps*, The Journal of chemical physics **140** (2014), no. 5, 054116.
- [43] Pierre Monmarché, *High-dimensional MCMC with a standard splitting scheme for the underdamped Langevin*, arXiv preprint arXiv:2007.05455 (2020).
- [44] Wenlong Mou, Yi-An Ma, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan, *High-order Langevin diffusion yields an accelerated MCMC algorithm*, arXiv preprint arXiv:1908.10859 (2019).
- [45] Elias AJF Peters and G de With, *Rejection-free Monte Carlo sampling for general potentials*, Physical Review E **85** (2012), no. 2, 026703.
- [46] Ruoqi Shen and Yin Tat Lee, *The randomized midpoint method for log-concave sampling*, Advances in neural information processing systems, 2019, pp. 2100–2111.
- [47] Konstantin S Turitsyn, Michael Chertkov, and Marija Vucelja, *Irreversible Monte Carlo algorithms for efficient sampling*, Physica D: Nonlinear Phenomena **240** (2011), no. 4-5, 410–414.
- [48] Paul Vanetti, Alexandre Bouchard-Côté, George Deligiannidis, and Arnaud Doucet, *Piecewise-deterministic Markov chain Monte Carlo*, arXiv preprint arXiv:1707.05296 (2017).

[49] Santosh Vempala and Andre Wibisono, *Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices*, Advances in neural information processing systems, 2019, pp. 8094–8106.

#### APPENDIX A. PROOF OF LEMMA 2.1

*Proof.* Let  $\lambda(t) = V_t \cdot \nabla_x U(X_t)$ . If no bouncing happens, then

$$\frac{d}{dt}\lambda(t) = V_t^\top \nabla_x^2 U(X_t) V_t \leq L|V_t|^2.$$

In addition,  $\lambda(t)$  decreases when bouncing happens, since there is some positive  $V_t^{(i)} \partial_{x_i} U(X_t)$  being changed to  $-V_t^{(i)} \partial_{x_i} U(X_t)$  while  $X_t$  and other  $V_t^{(j)}$ 's remain unchanged. Therefore, since  $|V_t|$  does not change between refreshments, we have for any  $t \in (0, T_{k+1} - T_k)$ ,<sup>1</sup>

$$(23) \quad \lambda(T_k + t) \leq \lambda(T_k) + tL|V_{T_k}|^2.$$

Notice for a convex function  $U(x)$  that satisfy Assumption 1, we have by co-coercivity

$$|\nabla U(x)|^2 \leq 2LU(x),$$

therefore for any  $t \in [0, T_{k+1} - T_k)$ , and any  $\alpha > 0$ ,

$$(24) \quad \begin{aligned} U(X_{T_k+t}) &= U(X_{T_k}) + \int_0^t \lambda(T_k + \tau) d\tau \\ &\leq U(X_{T_k}) + t\lambda(T_k) + \frac{Lt^2}{2}|V_{T_k}|^2 \\ &\stackrel{(11b),(11c)}{\leq} U(X_{T_k}) + t\left(\frac{d}{\sqrt{LT}}\right)^{1/2} |\nabla U(X_{T_k})| + 2Lt^2 d \\ &\leq U(X_{T_k}) + t\left(\frac{2d\sqrt{L}}{T}\right)^{1/2} \sqrt{U(X_{T_k})} + 2Lt^2 d \\ &\leq (1 + \alpha)U(X_{T_k}) + d\sqrt{L}t^2\left(\frac{1}{\sqrt{2T\alpha}} + 2\sqrt{L}\right). \end{aligned}$$

In particular,

$$U(X_{T_{k+1}}) \leq (1 + \alpha)U(X_{T_k}) + d\sqrt{L}t_{k+1}^2\left(\frac{1}{\sqrt{2T\alpha}} + 2\sqrt{L}\right).$$

Choosing  $\alpha = \frac{1}{\sqrt{LT}}$ , we have

$$U(X_{T_{k+1}}) \leq (1 + \alpha)U(X_{T_k}) + CdLt_{k+1}^2.$$

Now we apply the above formula iteratively and derive

$$\begin{aligned} U(X_T) &\leq (1 + \alpha)^{N+1}U(X_0) + CLd \sum_{k=1}^{N+1} (1 + \alpha)^{N-k+1} t_k^2 \\ &\leq (1 + \alpha)^{N+1} \left( U(X_0) + CLd \sum_{k=1}^{N+1} t_k^2 \right) \\ &\stackrel{(11d),(11e)}{\leq} C\sqrt{LT}d. \end{aligned}$$

Here we used  $\alpha = \frac{1}{\sqrt{LT}} = O(\frac{1}{N})$  so  $(1 + \alpha)^{N+1} = O(1)$ , which is true due to (11a).  $\square$

<sup>1</sup>We remark here that  $\lambda(t)$  is not well-defined at the bouncing times. Nevertheless, (23) still makes sense since  $\lambda(t)$  decreases at the bouncing events, and since we only use (23) in the time integral sense, this will not cause any problem.

## APPENDIX B. PROOF OF LEMMA 2.2

*Proof.* Let  $\Xi = \sum_{k=1}^{N+1} t_k^2$ . By properties of the Poisson process [28], if we condition on  $N$ , the distribution of  $T_1, T_2, \dots, T_N$  has the same joint distribution as that of  $N$  i.i.d. random variables uniformly distributed in  $(0, T)$ . This means

$$(25) \quad \mathbb{E}(\Xi | N) = \frac{N!}{T^N} \int_{t_1 + \dots + t_N < T} \left( \sum_{k=1}^N t_k^2 + (T - \sum_{k=1}^N t_k)^2 \right).$$

To calculate  $\mathbb{E}(\Xi | N)$ , let us define

$$I_1(N, T) = \int_{t_1 + \dots + t_N < T} \left( \sum_{k=1}^N t_k^2 + (T - \sum_{k=1}^N t_k)^2 \right)$$

and compute  $I_1(N, T)$  by induction in  $N$ . For  $N = 0$ , as the sum contains only one term,  $I_1(0, T) = T^2$ . An easy calculation shows that  $I_1(1, T) = \frac{2}{3}T^3$ . We will show in general

$$(26) \quad I_1(N, T) = \frac{2(N+1)}{(N+2)!} T^{N+2}.$$

Indeed, suppose (26) holds for  $N-1$ , we want to prove (26) for  $N$ , the starting point of which is the following observation:

$$I_1(N, T) = \int_{t_1 + \dots + t_N < T} t_1^2 dt_N \dots dt_1 + \int_0^T I_1(N-1, T-t_1) dt_1.$$

The first integral can be treated by integrating the variables one by one, from  $t_N$  to  $t_{N-1}$  and then  $t_{N-2}$ , etc.

$$(27) \quad \begin{aligned} \int_{t_1 + \dots + t_N < T} t_1^2 dt_N \dots dt_1 &= \int_{t_1 + \dots + t_{N-1} < T} t_1^2 (T - t_1 - \dots - t_{N-1}) dt_{N-1} \dots dt_1 \\ &= \frac{1}{2} \int_{t_1 + \dots + t_{N-2} < T} t_1^2 (T - t_1 - \dots - t_{N-2})^2 dt_{N-2} \dots dt_1 \\ &= \dots \\ &= \frac{1}{(N-1)!} \int_0^T t_1^2 (T - t_1)^{N-1} dt_1 = \frac{2}{(N+2)!} T^{N+2}. \end{aligned}$$

By the induction assumption (26) for  $N-1$  we have

$$\int_0^T I_1(N-1, T-t_1) dt_1 = \int_0^T \frac{2N}{(N+1)!} (T-t_1)^{N+1} dt_1 = \frac{2N}{(N+2)!} T^{N+2}.$$

Combining above with (27) we finish the proof for  $N$ . Therefore

$$\mathbb{E}(\Xi | N) = \frac{N!}{T^N} \frac{2(N+1)T^{N+2}}{(N+2)!} = \frac{2T^2}{N+2}.$$

The full expectation  $\mathbb{E}\Xi$  follows as  $N$  is a Poisson random variable

$$\begin{aligned} \mathbb{E}\Xi &= \sum_{n=0}^{\infty} \mathbb{E}(\Xi | N = n) \mathbb{P}(N = n) \\ &= \sum_{n=0}^{\infty} \frac{2T^2}{n+2} \frac{(\sqrt{LT})^n}{n!} e^{-\sqrt{LT}} \\ &= 2T^2 e^{-\sqrt{LT}} \sum_{n=0}^{\infty} \left( \frac{(\sqrt{LT})^n}{(n+1)!} - \frac{(\sqrt{LT})^n}{(n+2)!} \right) \\ &= \frac{2T}{\sqrt{L}} - \frac{2}{L} + \frac{2e^{-\sqrt{LT}}}{L} \leq \frac{2T}{\sqrt{L}}. \end{aligned}$$

To get the desired estimate, we apply Chebyshev's inequality using the second moment. By the same arguments leading towards (25), we have

$$\mathbb{E}(\Xi^2 | N) = \frac{N!}{T^N} \int_{t_1 + \dots + t_N < T} \left( \sum_{k=1}^N t_k^2 + (T - \sum_{k=1}^N t_k)^2 \right)^2.$$

Denote

$$I_2(N, T) = \int_{t_1 + \dots + t_N < T} \left( \sum_{k=1}^N t_k^2 + (T - \sum_{k=1}^N t_k)^2 \right)^2.$$

Using the same induction argument as the proof of (26), we can prove

$$I_2(N, T) = \frac{4(N+1)(N+6)}{(N+4)!} T^{N+4}.$$

This can be easily verified for  $N = 0, 1$  and the induction follows from the calculation:

$$\begin{aligned} I_2(N, T) &= \int_{t_1 + \dots + t_N < T} t_1^4 + \int_0^T I_2(N-1, T-t_1) dt_1 + 2 \int_0^T t_1^2 I_1(N-1, T-t_1) dt_1 \\ &= \frac{1}{(N-1)!} \int_0^T t_1^4 (T-t_1)^{N-1} dt_1 + \frac{4N(N+5)}{(N+3)!} \int_0^T (T-t_1)^{N+3} dt_1 \\ &\quad + \frac{4N}{(N+1)!} \int_0^T t_1^2 (T-t_1)^{N+1} dt_1 \\ &= \frac{4(N+1)(N+6)}{(N+4)!} T^{N+4}. \end{aligned}$$

This shows  $\mathbb{E}(\Xi^2 | N) = \frac{N!}{T^N} I_2(N, T) = \frac{4(N+6)}{(N+2)(N+3)(N+4)} T^4$ , and therefore

$$\begin{aligned} \mathbb{E}\Xi^2 &= \sum_{n=0}^{\infty} \mathbb{E}(\Xi^2 | N=n) \mathbb{P}(N=n) \\ &= \sum_{n=0}^{\infty} \frac{4T^4(n+6)}{(n+2)(n+3)(n+4)} \frac{(\sqrt{LT})^n}{n!} e^{-\sqrt{LT}} \\ &= 4T^4 e^{-\sqrt{LT}} \sum_{n=0}^{\infty} \left( \frac{(\sqrt{LT})^n}{(n+2)!} - 6 \frac{(\sqrt{LT})^n}{(n+4)!} \right) \\ &= \frac{4T^2}{L} - \frac{24}{L^2} + 8e^{-\sqrt{LT}} \left( \frac{T^2}{L} + \frac{3T}{L^{\frac{3}{2}}} + \frac{3}{L^2} \right). \end{aligned}$$

This means

$$\mathbb{E}(\Xi - \mathbb{E}\Xi)^2 = \frac{8T}{L^{\frac{3}{2}}} - \frac{28}{L^2} + 8e^{-\sqrt{LT}} \left( \frac{T^2}{L} + \frac{2T}{L^{\frac{3}{2}}} + \frac{4}{L^2} - \frac{4e^{-\sqrt{LT}}}{L^2} \right) \leq \frac{8T}{L^{\frac{3}{2}}},$$

where the inequality above holds for  $\sqrt{LT}$  larger than some universal constant (which we would assume as it is the interesting parameter regime).

Finally, to conclude the proof, we apply Chebyshev inequality to estimate the failure probability as

$$\mathbb{P}(\Xi \geq \frac{4T}{\sqrt{L}}) \leq \mathbb{P}(\Xi - \mathbb{E}\Xi \geq \frac{2T}{\sqrt{L}}) \leq \frac{L\mathbb{E}(\Xi - \mathbb{E}\Xi)^2}{4T^2} \leq \frac{2}{\sqrt{LT}}. \quad \square$$

#### APPENDIX C. PROOF OF LEMMA 3.1

*Proof.* The first step is to show that

$$(28) \quad \mathbb{E}_\mu |\partial_{x_i} U| \leq \sqrt{L}.$$

This is straightforward, since using integration by parts,

$$(29) \quad \mathbb{E}_\mu |\partial_{x_i} U|^2 = \int_{\mathbb{R}^d} (\partial_{x_i} U)^2 d\mu = \int_{\mathbb{R}^d} \partial_{x_i x_i} U d\mu \leq L,$$

and (28) then follows from Cauchy-Schwarz inequality.

The next step is to establish a concentration bound. Let  $G(x) = \psi(\partial_{x_i} U)$ , where  $\psi(a) = \psi(|a|)$  is a smooth nonnegative increasing function satisfying

$$\psi(0) = \psi'(0) = 0, \quad \psi(a) = |a| \quad \text{for } |a| \geq 1, \quad \text{and } |\psi'(a)| \leq 2,$$

and  $g(x) = \exp(\frac{\lambda}{2} \lambda G(x))$ . By the construction of  $G$ , we have

$$(30) \quad \mathbb{E}_\mu G = \mathbb{E}_\mu \psi(\partial_{x_i} U) \leq 2\mathbb{E}_\mu |\partial_{x_i} U| \leq 2\sqrt{L}.$$

Then  $\nabla g(x) = \frac{\lambda}{2} \psi'(\partial_{x_i} U) \nabla(\partial_{x_i} U) g(x)$ . By Lemma 3.2 for  $g(x)$ , we have

$$\begin{aligned} \mathbb{E}_\mu \exp(\lambda G) - \left(\mathbb{E}_\mu \exp\left(\frac{\lambda G}{2}\right)\right)^2 &= \text{Var}_\mu g(x) \\ &\leq \frac{\lambda^2}{4} \int_{\mathbb{R}^d} (\psi'(\partial_{x_i} U))^2 \nabla(\partial_{x_i} U) (\nabla^2 U)^{-1} \nabla(\partial_{x_i} U) g^2(x) d\mu \\ &\leq \lambda^2 \int_{\mathbb{R}^d} \nabla(\partial_{x_i} U) (\nabla^2 U)^{-1} \nabla(\partial_{x_i} U) g^2(x) d\mu \\ &= \lambda^2 \int_{\mathbb{R}^d} \partial_{x_i x_i} U g^2(x) d\mu \leq \lambda^2 L \mathbb{E}_\mu \exp(\lambda G). \end{aligned}$$

Thus for  $\lambda \leq \frac{1}{2\sqrt{L}}$  we have

$$(31) \quad \mathbb{E}_\mu \exp(\lambda G) \leq \frac{1}{1 - \lambda^2 L} \left(\mathbb{E}_\mu \exp\left(\frac{\lambda G}{2}\right)\right)^2.$$

Now we use (31) recursively, and we obtain for  $H(\lambda) := \mathbb{E}_\mu \exp(\lambda G)$ ,

$$(32) \quad H(\lambda) \leq \prod_{k=0}^{\infty} \left(\frac{1}{1 - \frac{\lambda^2 L}{4^k}}\right)^{2^k} \lim_{\ell \rightarrow \infty} H\left(\frac{\lambda}{\ell}\right)^\ell.$$

Notice

$$(33) \quad \lim_{\ell \rightarrow \infty} H\left(\frac{\lambda}{\ell}\right)^\ell = \lim_{\ell \rightarrow \infty} \left(\mathbb{E}_\mu \exp\left(\frac{\lambda G}{\ell}\right)\right)^\ell = \lim_{\ell \rightarrow \infty} \left(1 + \mathbb{E}_\mu \frac{\lambda G}{\ell}\right)^\ell = \exp(\lambda \mathbb{E}_\mu G).$$

Moreover, by [9, Proposition 4.1],

$$(34) \quad \prod_{k=0}^{\infty} \left(\frac{1}{1 - \frac{\lambda^2 L}{4^k}}\right)^{2^k} \leq \frac{1 + \lambda\sqrt{L}}{1 - \lambda\sqrt{L}}.$$

Substituting (33) and (34) into (32), we obtain

$$H(\lambda) \leq \frac{1 + \lambda\sqrt{L}}{1 - \lambda\sqrt{L}} \exp(\lambda \mathbb{E}_\mu G).$$

Finally, combining the above exponential moment bound of  $G$  with Chebyshev inequality, we get

$$\mathbb{P}_\mu \left( G(x) \geq \mathbb{E}_\mu G + r \right) \leq \exp(-\lambda r) \frac{1 + \lambda\sqrt{L}}{1 - \lambda\sqrt{L}}.$$

Now take  $\lambda = 1/2\sqrt{L}$ , and  $r = 2c\sqrt{L} \log d$ , and using (30) (noticing  $G(x) = |\partial_{x_i} U|$  when  $G(x) \geq r$ ), we arrive at

$$\mathbb{P}_\mu \left( |\partial_{x_i} U| \geq 2\sqrt{L} + 2c\sqrt{L} \log d \right) \leq 3d^{-c}. \quad \square$$



DEPARTMENT OF MATHEMATICS, DEPARTMENT OF PHYSICS, AND DEPARTMENT OF CHEMISTRY, DUKE UNIVERSITY, DURHAM NC 27708

*Email address:* [jianfeng@math.duke.edu](mailto:jianfeng@math.duke.edu)

DEPARTMENT OF MATHEMATICS, DUKE UNIVERSITY, DURHAM NC 27708

*Email address:* [lihan@math.duke.edu](mailto:lihan@math.duke.edu)