

1 We thank all the reviewers for the valuable comments.

2 **To Reviewer 1:**

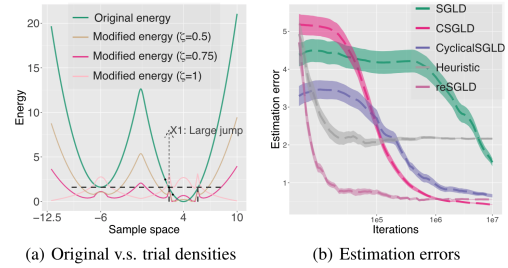
3 Q1. *Details on the vanishing gradient problem in flat-histogram SGLD:* The original step function in formula (4)
 4 leads to $\frac{\partial \log \Psi_\theta(u)}{\partial u} = \frac{1}{\Psi_\theta(u)} \frac{\partial \Psi_\theta(u)}{\partial u} = 0$ almost everywhere. Combining it with (6) leads to $\nabla_x \log \varpi_{\Psi_\theta}(x) =$
 5 $-\left[1 + \zeta \tau \frac{\partial \log \Psi_\theta(u)}{\partial u}\right] \frac{\nabla_x U(x)}{\tau} = -\frac{\nabla_x U(x)}{\tau}$. Therefore, the naive flat-histogram SGLD will behave like SGLD and fail
 6 to converge to the flattened density (2).

7 Q2. *Advantages of CSGLD over M-SGD:* (i) CSGLD belongs to the class of adaptive biasing force algorithms and
 8 has the potential to exponentially speed-up the computation [1], i.e. the fine-tuned CSGLD should outperform the
 9 fine-tuned M-SGD. (ii) CSGLD also belongs to the class of dynamic importance sampling algorithms and is able to
 10 quantify uncertainty for its estimation and prediction. While M-SGD is an optimization algorithm and produces a point
 11 estimate only. [1] Long-time convergence of an adaptive biasing force (ABF) method. Nonlinearity. 2008.

12 **To Reviewer 2:**

13 Q1. *The working distribution is different from the original one:* CSGLD belongs to the class of dynamic importance
 14 sampling algorithms, for which each sample is generated with an importance weight and they can be used for inference
 15 of the target distribution via a weighted averaging estimator. We have provided theoretical guarantee for the convergence
 16 of the weighted averaging estimator. The samples from the target distribution can also be obtained via an importance
 17 resampling step from the pool of importance samples generated by CSGLD.

18 Q2. *Asymmetric modes with a heuristic baseline:* Following section
 19 4.1, we set $\pi(x) = 1/6N(-6, 1) + 5/6N(4, 1)$ and included a
 20 baseline called *Heuristic*, which injects a Gaussian noise $N(0, 2^2)$
 21 whenever $|\nabla \tilde{U}(x)| < 0.1$. As shown by Figure (b) of this rebuttal,
 22 *Heuristic* performs quite well in the early period due to the heuristic
 23 random walk helping to escape local traps. However, it leads to a very
 24 large prediction error in the long run as the sampling equilibrium was
 25 broken by *Heuristic*, but the samples were not properly weighted.



26 Q3. *Discussions on the number of partitions:* A fine partition will
 27 lead to a smaller discretization error, but it increases the risk in stability. In particular, it will lead to large bouncy
 28 jumps around optima (a large negative learning rate, i.e., $\frac{\log \theta(2) - \log \theta(1)}{\Delta u} \ll 0$ in formula (8) will be caused there).
 29 Empirically, we suggest to partition the sample space into a moderate number of subregions, e.g. 10-1000, to balance
 30 between stability and discretization error.

31 **To Reviewer 3:**

32 Q1. *Drawbacks of simulated annealing (SA) and replica exchange SGLD (reSGLD)/parallel tempering:* SA can only be
 33 used in optimization, and it might get stuck in a poor local minimum if the temperature decreases too fast. The reSGLD
 34 requires a large correction to reduce the bias, which may lead to a large bias and insignificant accelerations.

35 Q2. *Missing baselines:* We further compared CSGLD with CyclicalSGLD and reSGLD on an asymmetric mixture
 36 distribution. All algorithms were run 10^7 iterations. For CSGLD, we set $\tilde{U}(x) \sim \mathcal{N}(U(x), 4^2)$, the default stepsize
 37 $\alpha = 0.1$ and temperature $T = 1$. For CyclicalSGLD, we set $\alpha_0 = 1$ and 100 cycles and the threshold $\beta = 0.9$; for
 38 reSGLD, we additionally include a high-temperature process with $T = 5$ and a correction of 16. Figure (b) of this
 39 rebuttal shows that CSGLD is inferior to the baselines at the beginning, but it eventually outperforms the baselines as
 40 the learning of θ_k is mature. We will include the baselines and references in the next version.

41 Q3. The gradient-vanishing problem in SGLD is not clear: Please refer to our reply to Q1 of Reviewer 1.

42 **To Reviewer 4:**

43 Q1. *Comments on bizarre peaks:* A bizarre peak always indicates that there is a local minimum of the same energy in
 44 somewhere of the sample space, as the sample space is partitioned according to the energy function in CSGLD. For
 45 example, Figure (a) of this rebuttal shows a bizarre peak at x_1 . Although x_1 is not a local minimum, it has the same
 46 energy as “-6” which is a local minimum. Note that in CSGLD, x_1 and “-6” belongs to the same subregion.

47 Q2. *In very low density level with no driving forces.* Within the same subregion, CSGLD is reduced to SGLD. Therefore,
 48 in the very low density region, it will not move farther away, but still move toward the high density region.

49 Q3. *Choices of the partition and ζ :* Regarding the partition, please refer to our response to Q3 of reviewer 2. With
 50 respect to ζ , we suggest to tune ζ to obtain the largest barrier reductions such as $\zeta = 0.75$ in Figure 1(a) of the paper.