

1 We thank the reviewers for their close reading, detailed comments, and overall positive assessment. We address the
 2 questions raised by each reviewer separately.

3 **Reviewer 1:** Thanks for the appreciation and suggestions for paper refinement. We will fix the typos and polish the
 4 figures for the final version.

5 **Reviewer 2:** We will improve the flow and formatting of the paper, and fix the references in the final version.

6 • **Empirical experiments on known exponential family distributions.**

7 We have conducted a new experiment comparing
 8 ADE with CD and SM on multivariate Gaussians with dif-
 9 ferent dimensions and banana datasets, where we know the
 10 potential functions, to investigate the effect of the number
 11 of dimensionality and complexity of the potential function
 12 on these algorithms. For fairness, we compare the ADE with

Table 1: Parameter recovering on synthetic datasets.

Dataset	SM	CD-5	ADE
2D-Gaussian	2.18×10^{-3}	5.67×10^{-3}	2.28×10^{-3}
5D-Gaussian	3.17×10^{-3}	4.19×10^{-1}	3.09×10^{-3}
10D-Gaussian	3.90×10^{-3}	6.36×10^{-1}	3.23×10^{-3}
Banana	2.33×10^{-2}	6.72×10^{-2}	6.00×10^{-3}

13 on these algorithms. For fairness, we compare the ADE with
 14 HMC-3 on 2-layer MLP to CD with HMC-5, which has the same number of operations. The models are trained with
 15 1000 samples. The 5 runs average results, in terms of RMSE between learned parameters and the true parameters,
 16 are reported in Table 1. As we can see, ADE consistently achieves comparable or the best performance. We will add
 17 these comparison and a detailed analysis to the final version.

18 • **ADE for exponential families on discrete variables.** For an enumerable discrete space, the partition function is
 19 tractable and the MLE can be computed exactly. For a combinatorial discrete space, our primal-dual MLE is still
 20 valid. The major difficulty lies in the dual sampler parametrization, where the HMC/Langevin embeddings are not
 21 applicable since the gradient is not well-defined. We are exploring alternative sampling algorithm embeddings, *e.g.*,
 22 the importance sampler and Gibbs sampler, for future work.

23 • **ADE limitations and how to overcome.** The cost of gradient computation in respect of the potential function will
 24 be proportional to the number of sampling layers T . We provide a gradient approximation which is independent of T ,
 25 however, using Danskin’s theorem. See Appendix C for details.

26 • **Parameter tuning in ADE.** When we use neural networks to parametrize $q^0(x, v)$ and the generalized HMC layers in
 27 ADE, then the parameter tuning requirements for ADE and GANs are comparable, *i.e.*, we tune the inner optimization
 28 stepsize and the schedule between the inner minimisation and outer maximisation. When we use a nonparametric
 29 $q^0(x, v)$ and a vanilla HMC layer, the parameter tuning requirements of ADE and CD are comparable, *i.e.*, we tune
 30 the optimization stepsize in ADE, and the leapfrog step size for CD.

31 **Reviewer 3:**

32 • **Main contribution.** Re: “[the authors] further conduct T vanilla HMC steps to approximately solve it.” We are most
 33 certainly *not* using a fixed (vanilla) HMC sampler to directly approximate the dual problem! The contribution of our
 34 paper is to develop a family of probability distribution parametrizations for the dual, obtained by first unfolding the
 35 HMC sampling steps and then learning neural network functions in place of these steps. The optimal solution of
 36 the coupled model and dual (Section 3.3) is exactly the target distribution. Thus, the neural network “dual sampler”
 37 is automatically adjusted to reduce the finite-step approximation error of HMC. This is not true of a fixed HMC
 38 sampler. Finally (Theorem 4) the learned “neural HMC” dual has an explicit closed form density estimate (eq.
 39 17), which can in turn be used to evaluate the entropy. To our knowledge, none of the above contributions have
 40 appeared elsewhere. We will clarify these points in the final version. We agree that a preconditioner K can be used
 41 for the kinetic energy in the HMC unfolding. Indeed, in the generalized HMC layer in Eq.(13), we considered even
 42 more general preconditioners, which include not just linear preconditioners but also nonlinear projections to latent
 43 low-dimension spaces, with automatically learned parameters.

44 • **Effect of SGD for inner problem.** As we showed in Theorem 1, the optimal solution to the inner minimisation is
 45 the target sampling distribution. Therefore, SGD will *not* “blur the precision of HMC”. Instead, solving the inner
 46 minimisation by SGD will lead the dual sampler, *i.e.*, the proposed dynamics embedding network, to the actual target
 47 distribution and reduce the finite-step approximation error.

48 Theorem 3 justifies the flexibility of the proposed unfolded dynamics-based sampler, saying that as $T \rightarrow \infty$, the
 49 neural network parametrization for the dual sampler can approach any exponential family density arbitrarily well.
 50 This is not related to the effect of SGD in the inner minimisation.

51 • **Number of sampling steps.** In our synthetic experiments, we have the initial distribution in ADE for v as a standard
 52 Gaussian and for x as $q^0(x)$, which is parametrized by a 10-layer normalizing flow, with extra 5 dynamics-based
 53 sampling layers. For fairness, we compared to CD with 15 fixed HMC-steps. The experiments demonstrate that given
 54 same number of steps, ADE will achieve a smaller approximation error by learning the HMC sampler, resulting in
 55 better performance. The leapfrog stepsize is learned in ADE and tuned in CD-HMC following [40]. We tested HMC
 56 varying the number of steps as $T = 1, 3, 5, 15$. HMC-15 performs best, as reported in the main text. This is also
 57 confirmed by our theory: as T becomes larger, the approximation error becomes smaller, and thus, the performance
 58 achieved is better.