

LeapfrogLayers: A Trainable Framework for Effective Topological Sampling

Sam Foreman,^{a,*} Xiao-Yong Jin^{a,b} and James C. Osborn^{a,b}

^a*Leadership Computing Facility, Argonne National Laboratory,
Lemont, IL, USA*

^b*Computational Science Division, Argonne National Laboratory,
Lemont, IL, USA*

E-mail: foremans@anl.gov, xjin@anl.gov, osborn@alcf.anl.gov

We introduce LeapfrogLayers, an invertible neural network architecture that can be trained to efficiently sample the topology of a 2D $U(1)$ lattice gauge theory. We show an improvement in the integrated autocorrelation time of the topological charge when compared with traditional HMC, and propose methods for scaling our model to larger lattice volumes. Our implementation is open source, and is publicly available on github at github.com/saforem2/l2hmc-qcd.

*The 38th International Symposium on Lattice Field Theory
26-30 July 2021
Zoom / Gather @ MIT, Cambridge MA, USA*

*Speaker

1. Introduction

The main task in lattice field theory calculations is to evaluate integrals of the form

$$\langle O \rangle \propto \int [D\mathbf{x}] O(\mathbf{x}) p(\mathbf{x}), \quad (1)$$

for some multivariate target distribution $p(\mathbf{x}) \propto e^{-S(\mathbf{x})}/Z$. We can approximate the integral using Markov Chain Monte Carlo (MCMC) sampling techniques. This is done by sequentially generating a chain of configurations $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, with $\mathbf{x}_i \sim p(\mathbf{x})$ and averaging the value of the function $O(\mathbf{x})$ over the chain

$$\bar{O} = \frac{1}{N} \sum_{n=1}^N O(\mathbf{x}_n). \quad (2)$$

Accounting for correlations between states in the chain, the sampling variance of this estimator is given by

$$\sigma^2 = \frac{\tau_O^{\text{int}}}{N} \sum_{n=1}^N \left[O(\mathbf{x}_n) - \bar{O} \right]^2 \quad (3)$$

where τ_O^{int} is the integrated autocorrelation time. This quantity can be interpreted as the additional time required for these induced correlations to become negligible.

1.1 Charge Freezing

The ability to efficiently generate independent configurations is currently a major bottleneck for lattice simulations. In this work, we consider a $U(1)$ gauge model on a 2D periodic lattice. The theory is defined in terms of the link variables $U_\mu(x) = e^{ix} \in U(1)$ with $x \in [-\pi, \pi)$. Our target distribution is given by $p(\mathbf{x}) \propto e^{-S_\beta(\mathbf{x})}$, where $S_\beta(\mathbf{x})$ is the Wilson action

$$S_\beta(\mathbf{x}) = \beta \sum_P 1 - \cos x_P, \quad (4)$$

defined in terms of the sum of the gauge variables around the elementary plaquette, $x_P \equiv x_\mu(n) + x_\nu(n + \hat{\mu}) - x_\mu(n + \hat{\nu}) - x_\nu(n)$, starting at site n of the lattice.¹ For a given lattice configuration, we can calculate the topological charge $Q \in \mathbb{Z}$ using

$$Q_{\mathbb{Z}} = \frac{1}{2\pi} \sum_P [x_P], \quad \text{where} \quad (5)$$

$$[x_P] \equiv x_P - 2\pi \left\lfloor \frac{x_P + \pi}{2\pi} \right\rfloor. \quad (6)$$

¹Here, the sum \sum_P denotes the sum over all such plaquettes.

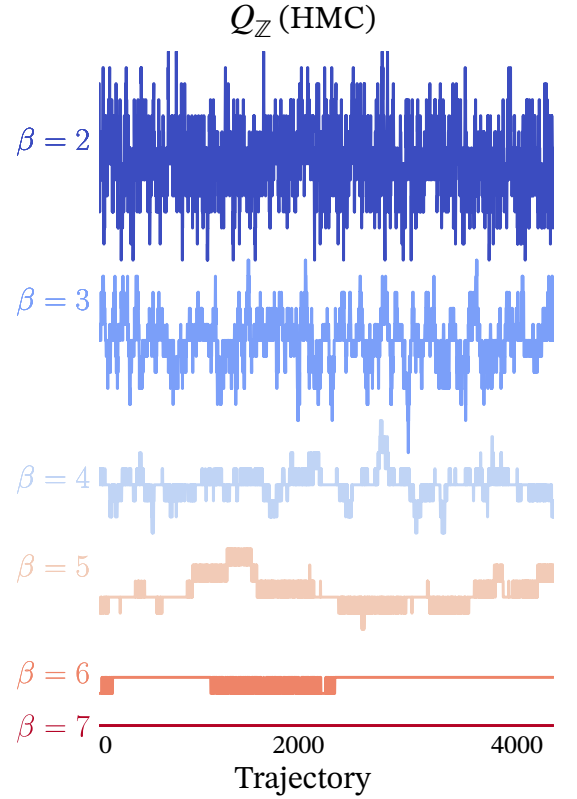


Figure 1: Illustration of the topological charge, $Q_{\mathbb{Z}}$, freezing as $\beta : 2 \rightarrow 7$ for traditional HMC.

As $\beta \rightarrow \infty$, the $Q_{\mathbb{Z}} = 0$ mode becomes dominant and we see that the value of $Q_{\mathbb{Z}}$ remains fixed for large durations of the simulation. This freezing can be quantitatively measured by introducing the *tunneling rate*²

$$\delta Q_{\mathbb{Z},i} \equiv |Q_{\mathbb{Z},i+1} - Q_{\mathbb{Z},i}| \in \mathbb{Z}, \quad (7)$$

which serves as a measure for how efficiently our chain is able to jump (tunnel) between sectors of distinct topological charge. From Figure 1, we can see that $\overline{\delta Q_{\mathbb{Z}}} \rightarrow 0$ as $\beta \rightarrow \infty$.

2. Hamiltonian Monte Carlo (HMC)

The Hamiltonian Monte Carlo algorithm begins by introducing a fictitious momentum, \mathbf{v} , for each coordinate variable, \mathbf{x} , typically taken from an independent normal distribution. This allows us to write the joint target density of the $\xi \equiv (\mathbf{x}, \mathbf{v})$ system as

$$p(\mathbf{x}, \mathbf{v}) = p(\mathbf{x})p(\mathbf{v}) = e^{-S_{\beta}(\mathbf{x})} e^{-\mathbf{v}^2/2} = e^{-\mathcal{H}(\mathbf{x}, \mathbf{v})} \quad (8)$$

where $\mathcal{H}(\xi) = \mathcal{H}(\mathbf{x}, \mathbf{v}) = S_{\beta}(\mathbf{x}) + \frac{1}{2}\mathbf{v}^2$ is the Hamiltonian of the system. We use the *leapfrog integrator* to approximately numerically integrate Hamilton's equations

$$\dot{\mathbf{x}} = \frac{\partial \mathcal{H}}{\partial \mathbf{v}}, \quad \dot{\mathbf{v}} = -\frac{\partial \mathcal{H}}{\partial \mathbf{x}} \quad (9)$$

along iso-probability contours of $\mathcal{H} = \text{const.}$ from $\xi = (\mathbf{x}, \mathbf{v}) \rightarrow (\mathbf{x}^*, \mathbf{v}^*) = \xi^*$. The error in this integration is then corrected by a Metropolis-Hastings (MH) accept/reject step.

2.1 Leapfrog Integrator

1. Starting from \mathbf{x}_i , resample the momentum $\mathbf{v}_i \sim \mathcal{N}(0, \mathbb{I})$ and construct the state $\xi = (\mathbf{x}_i, \mathbf{v}_i)$.
2. Generate a *proposal configuration* ξ^* by integrating $\dot{\xi}$ along $\mathcal{H} = \text{const.}$ for N leapfrog steps. i.e.

$$\xi \rightarrow \xi_1 \rightarrow \dots \rightarrow \xi_N \equiv \xi^*, \quad (10)$$

where a single leapfrog step $\xi_i \rightarrow \xi_{i+1}$ above consists of:

$$\text{(a.) } \mathbf{v}' \leftarrow \mathbf{v} - \frac{\varepsilon}{2} \partial_{\mathbf{x}} S(\mathbf{x}), \quad \text{(b.) } \mathbf{x}' \leftarrow \mathbf{x} + \varepsilon \mathbf{v}', \quad \text{(c.) } \mathbf{v}'' \leftarrow \mathbf{v}' - \frac{\varepsilon}{2} \partial_{\mathbf{x}} S(\mathbf{x}). \quad (11)$$

3. At the end of the trajectory, accept or reject the proposal configuration ξ^* using the MH test

$$\mathbf{x}_{i+1} \leftarrow \begin{cases} \mathbf{x}^* & \text{with probability } A(\xi^*|\xi) \equiv \min \left\{ 1, \frac{p(\xi^*)}{p(\xi)} \right\} \\ \mathbf{x}_i & \text{with probability } 1 - A(\xi^*|\xi). \end{cases} \quad (12)$$

An illustration of this procedure can be seen in Figure 2.

²As measured between subsequent states in our chain $i, i+1$.

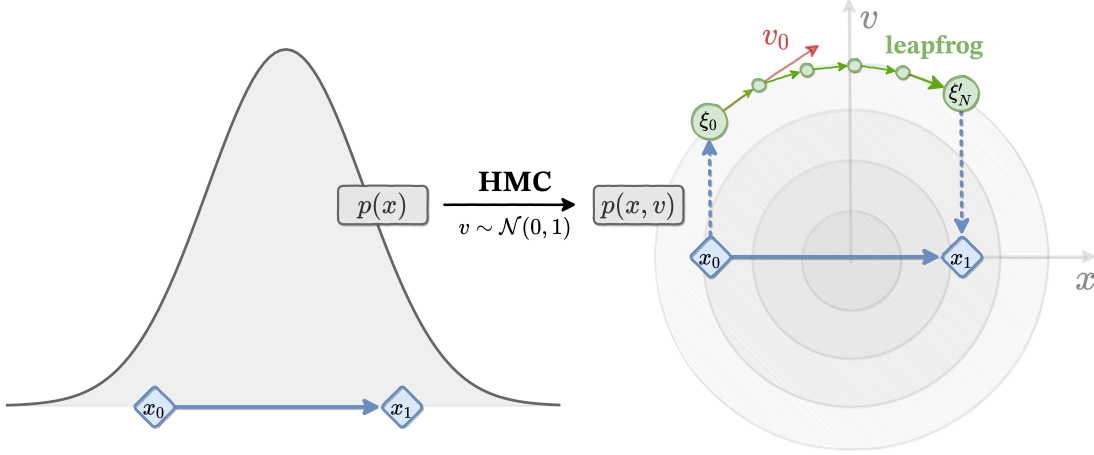


Figure 2: High-level overview of the HMC algorithm.

2.2 Issues with HMC

Re-sampling the momentum at the start of each trajectory causes the energy levels we explore to be randomly selected. This is reminiscent of the “random-walk” behavior of traditional MCMC and leads to a slow exploration of our target distribution (i.e. long autocorrelations). Additionally, the HMC sampler is known to have difficulty traversing low-density zones, resulting in poor performance for distributions which have multiple isolated modes. This is particularly relevant in the case of sampling topological quantities in lattice gauge models.

3. Generalizing HMC: LeapfrogLayers

In order to preserve the asymptotic behavior of HMC, our update must be explicitly reversible with a tractable Jacobian determinant. To simplify notation, we introduce two functions, Γ (Λ) to denote the \mathbf{v} (\mathbf{x}) updates. As in HMC, we follow the general pattern of performing alternating updates of \mathbf{v} and \mathbf{x} .

We can ensure our update is reversible by splitting the \mathbf{x} update into two parts and sequentially updating complementary subsets using a binary mask \mathbf{m} and its complement $\bar{\mathbf{m}}$. As in [5, 8], we introduce $d \sim \mathcal{U}(+, -)$, distributed independently of both \mathbf{x} and \mathbf{v} , to determine the “direction” of our update³. Here, we associate $+$ ($-$) with the forward (backward) direction and note that running sequential updates in opposite directions has the effect of inverting the update. We denote the complete state by $\xi = (\mathbf{x}, \mathbf{v}, \pm)$, with target density given by $p(\xi) = p(\mathbf{x})p(\mathbf{v})p(\pm)$.

Explicitly, we can write this series of updates as⁴

$$\begin{aligned} \text{(a.) } \mathbf{v}' &\leftarrow \Gamma^\pm[\mathbf{v}; \zeta_{\mathbf{v}}] & \text{(b.) } \mathbf{x}' &\leftarrow \mathbf{m} \odot \mathbf{x} + \bar{\mathbf{m}} \odot \Lambda^\pm[\mathbf{x}_{\bar{\mathbf{m}}}; \zeta_{\bar{\mathbf{x}}}] \\ \text{(c.) } \mathbf{x}'' &\leftarrow \bar{\mathbf{m}} \odot \mathbf{x}' + \mathbf{m} \odot \Lambda^\pm[\mathbf{x}_{\mathbf{m}}; \zeta_{\mathbf{x}'}] & \text{(d.) } \mathbf{v}'' &\leftarrow \Gamma^\pm[\mathbf{v}', \zeta_{\mathbf{v}'}] \end{aligned}$$

where $\zeta_{\bar{\mathbf{x}}} = [\bar{\mathbf{m}} \odot \mathbf{x}, \mathbf{v}]$, $\zeta_{\mathbf{x}} = [\mathbf{m} \odot \mathbf{x}, \mathbf{v}]$ ($\zeta_{\mathbf{v}} = [\mathbf{x}, \partial_{\mathbf{x}} S(\mathbf{x})]$) is independent of \mathbf{x} (\mathbf{v}) and is passed as input to the update functions Λ^\pm (Γ^\pm).

³As indicated by the superscript \pm on Γ^\pm , Λ^\pm in the update functions.

⁴Here we denote by $\mathbf{x}_{\mathbf{m}} = \mathbf{m} \odot \mathbf{x}$ and $\mathbf{x}_{\bar{\mathbf{m}}} = \bar{\mathbf{m}} \odot \mathbf{x}$ with $\mathbb{1} = \mathbf{m} + \bar{\mathbf{m}}$.

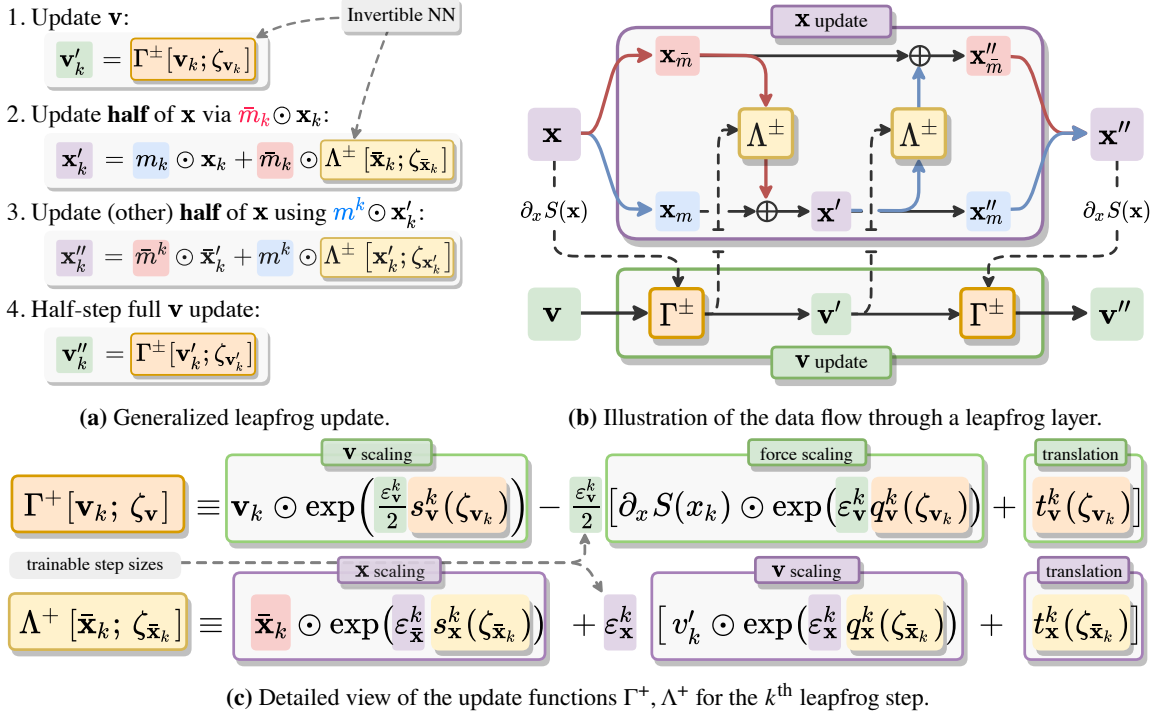


Figure 3: Illustrations of the generalized leapfrog update. In Figure 3a, 3c, we denote $\bar{\mathbf{x}}_k = \bar{\mathbf{m}} \odot \mathbf{x}_k$.

3.1 Network Details

Normalizing flows [4] are an obvious choice for the structure of the update functions. These architectures are easily invertible while maintaining a tractable Jacobian determinant, and have also been shown to be effective at approximating complicated target distributions in high-dimensional spaces [1, 2, 4, 5, 7, 8, 10, 11, 14].

We maintain separate networks Γ, Λ with identical architectures for updating \mathbf{v} and \mathbf{x} , respectively. Without loss of generality, we describe below the details of the \mathbf{x} update for the forward ($d = +$) direction, $\Lambda^+[\mathbf{x}_{\bar{\mathbf{m}}}; \zeta_{\bar{\mathbf{x}}}]$ ⁵. For simplicity, we describe the data flow through a single leapfrog layer, which takes as input $\zeta_{\bar{\mathbf{x}}} = (\mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{v})$. For the 2D $U(1)$ model, the gauge links are encoded as $[\cos(x), \sin(x)]$ for $x \in [-\pi, \pi]$. Explicitly⁶,

$$h_1 = \sigma(w_x \mathbf{x} + w_v \mathbf{v} + b_1) \in \mathbb{R}^{n_1} \quad (13)$$

$$h_2 = \sigma(w_2 h_1 + b_2) \in \mathbb{R}^{n_2} \quad (14)$$

\vdots

$$h_k = \sigma(w_k h_{k-1} + b_{k-1}) \in \mathbb{R}^{n_k} \implies \quad (15)$$

$$\text{(a.) } s_x = \lambda_s \tanh(w_s h_k + b_s); \quad \text{(b.) } t_x = w_t h_k + b_t; \quad \text{(c.) } q_x = \lambda_q \tanh(w_q h_k + b_q);$$

where the outputs s_x, t_x, q_x are of the same dimensionality as \mathbf{x} , and λ_s, λ_q are trainable parameters. These outputs are then used to update \mathbf{x} , as shown in Figure 3.

⁵To get the update for the $d = -$ direction, we invert the update functions and run them in the opposite order.

⁶Here $\sigma(\cdot)$ is a generic nonlinear activation function acting independently on the elements of the output vector.

3.2 Training Details

Our goal in training the network is to find a sampler that efficiently jumps between different topological charge sectors. This can be done by constructing a loss function that maximizes the expected squared charge difference between the initial (ξ) and proposal (ξ^*) configuration generated by the sampler. Explicitly, we define

$$\mathcal{L}_\theta(\xi^*, \xi) = A(\xi^*|\xi)(Q_\mathbb{R}^* - Q_\mathbb{R})^2 \quad (16)$$

where $Q_\mathbb{R} = \frac{1}{2\pi} \sum_P \sin x_P \in \mathbb{R}$ is a real-valued approximation of the usual (integer-valued) topological charge $Q_\mathbb{Z} \in \mathbb{Z}$. This ensures that our loss function is continuous which helps to simplify the training procedure. The acceptance probability

$$A(\xi^*|\xi) = \min \left\{ 1, \frac{p(\xi^*)}{p(\xi)} \mathcal{J}(\xi^*, \xi) \right\} \quad (17)$$

now includes an extra Jacobian factor $\mathcal{J}(\xi^*, \xi)$ which allows the inclusion of non-symplectic update steps. The Jacobian comes from the $v(x)$ scaling term in the $v(x)$ update, and is easily calculated. For completeness, the details of an individual training step are summarized in Sec 3.2.1.

3.2.1 Training Step

1. Resample $\mathbf{v} \sim \mathcal{N}(0, \mathbb{I})$, $d \sim \mathcal{U}(+, -)$, and construct initial state $\xi = (\mathbf{x}, \mathbf{v}, \pm)$
2. Generate the proposal configuration ξ^* by passing the initial state sequentially through N_{LF} *leapfrog layers*: $\xi \rightarrow \xi_1 \rightarrow \dots \rightarrow \xi_{N_{\text{LF}}} = \xi^*$
3. Compute the Metropolis-Hastings acceptance $A(\xi^*|\xi) = \min \left\{ 1, \frac{p(\xi^*)}{p(\xi)} \mathcal{J}(\xi^*, \xi) \right\}$
4. Evaluate the loss function $\mathcal{L} \leftarrow \mathcal{L}_\theta(\xi^*, \xi)$, and backpropagate gradients to update weights
5. Evaluate Metropolis-Hastings criteria and assign the next state in the chain according to

$$\mathbf{x}_{t+1} \leftarrow \begin{cases} \mathbf{x}^* & \text{with prob. } A(\xi^*|\xi) \\ \mathbf{x} & \text{with prob. } 1 - A(\xi^*|\xi). \end{cases}$$

3.3 Annealing

As an additional tool to help improve the quality of the trained sampler, we scale the action during the N_T training steps using the target distribution $p_t(x) \propto e^{-\gamma_t S(x)}$ for $t = 0, 1, \dots, N_T$. The scale factors, γ_t , monotonically increase according to an *annealing schedule* up to $\gamma_{N_T} \equiv 1$, with $|\gamma_{i+1} - \gamma_i| \ll 1$. For $\gamma_i < 1$, this helps to rescale (shrink) the energy barriers between isolated modes, allowing the training to experience sufficient tunneling even when the final distribution is difficult to sample.

4. Results

We can measure the performance of our approach by comparing the integrated autocorrelation time against traditional HMC. We see in Figure 4 that our estimate of the integrated autocorrelation time is much shorter for the trained model across $\beta \in [2, 7]$. To better understand how these transformations effect the HMC update, we can look at how various quantities change over the course of a trajectory, as shown in Figure 6. We see from the plot of $\mathcal{H} - \sum \log |\mathcal{J}|$ in Figure 6c that the trained sampler artificially increases the energy towards the middle of the trajectory. This is analogous to reducing the coupling β during the trajectory, as can be seen in Figure 7b. In particular, we can see that for $\beta = 7$, the value of the plaquette at in the middle of the trajectory roughly corresponds to the expected value at $\beta = 3$, indicated by the horizontal dashed line. This effective reduction in the coupling constant allows our sampler to mix topological charge values in the middle of the trajectory, before returning to the stationary distribution at the end, as can be seen in Figure 6b. By looking at the variation in the average plaquette $\langle x_P - x_P^* \rangle$ over a trajectory for models trained at multiple values of beta we are able to better understand how our sampler behaves. This allows our trajectories to explore new regions of the target distribution which may have been previously inaccessible.

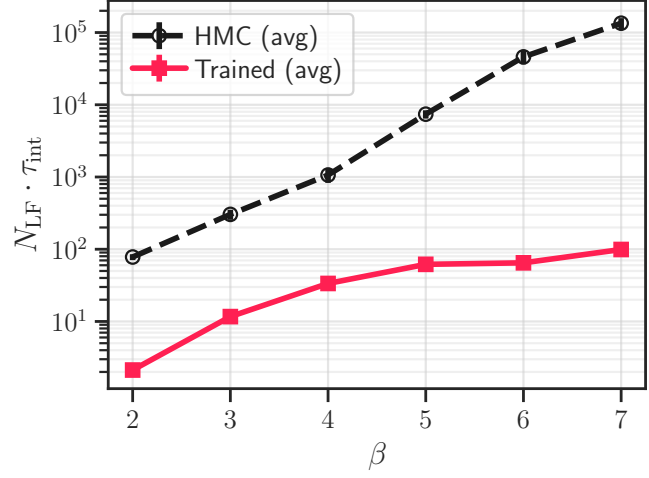


Figure 4: Estimated $\tau_{\text{int}}^{Q_R}$ vs β .

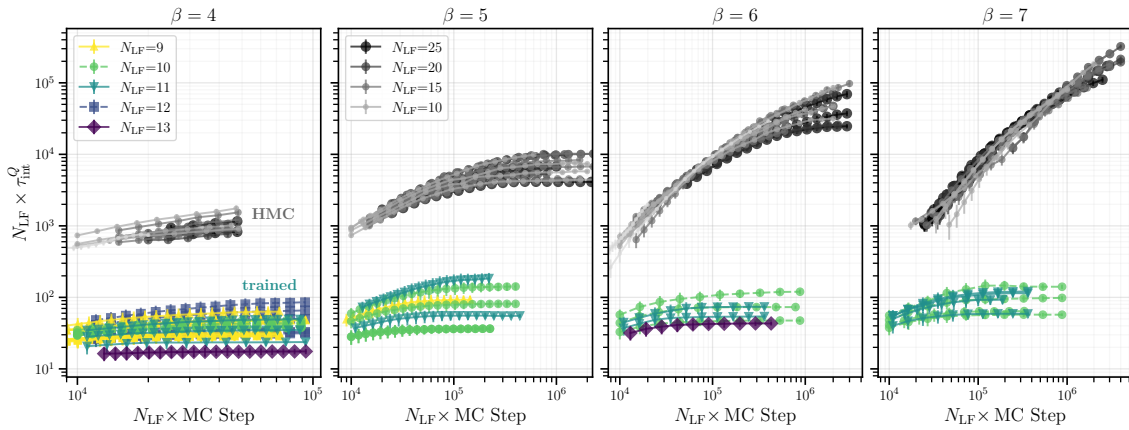


Figure 5: Comparison of the integrated autocorrelation time for trained models vs HMC with different trajectory lengths, N_{LF} , at $\beta = 4, 5, 6, 7$.

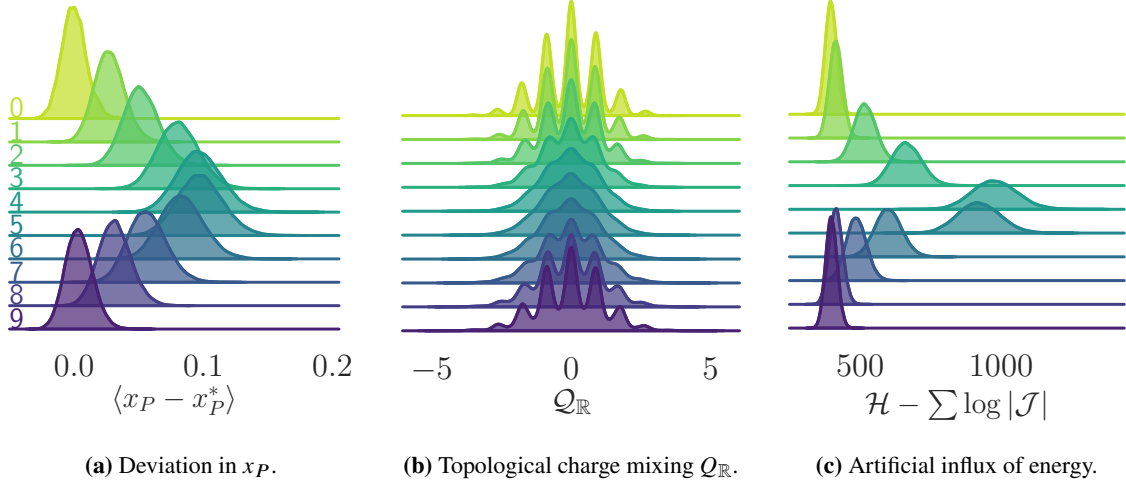


Figure 6: Illustration of how different quantities evolve over a single trajectory.

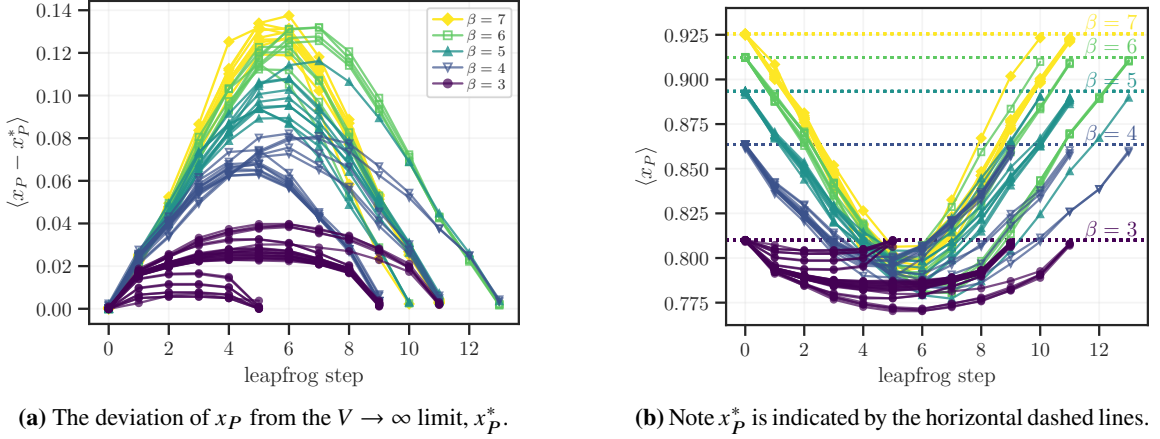


Figure 7: Plots showing how the average plaquette varies with over a single trajectory for different models trained at $\beta = 3, 4, \dots, 7$.

5. Conclusion

In this work we have proposed a generalized sampling procedure for HMC that can be used for generating gauge configurations in the 2D $U(1)$ lattice model. We showed that our trained model is capable of outperforming traditional sampling techniques across a range of inverse coupling constants, as measured by the integrated autocorrelation time of the topological charge. By looking at the evolution of different quantities over the generalized leapfrog trajectory, we are able to gain insight into the mechanism driving this improved behavior.

6. Acknowledgments

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security

Administration. This research was performed using resources of the Argonne Leadership Computing Facility (ALCF), which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357. This work describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the work do not necessarily represent the views of the U.S. DOE or the United States Government. Results presented in this research were obtained using the Python [13], programming language and its many data science libraries [3, 6, 9, 12]

References

- [1] Michael S. Albergo, Denis Boyda, Daniel C. Hackett, Gurtej Kanwar, Kyle Cranmer, Sébastien Racanière, Danilo Jimenez Rezende, and Phiala E. Shanahan. Introduction to Normalizing Flows for Lattice Field Theory. *arXiv e-prints*, 1 2021.
- [2] Denis Boyda, Gurtej Kanwar, Sébastien Racanière, Danilo Jimenez Rezende, Michael S. Albergo, Kyle Cranmer, Daniel C. Hackett, and Phiala E. Shanahan. Sampling using $SU(N)$ gauge equivariant flows. *Phys. Rev. D*, 103(7):074504, 2021.
- [3] Thomas A Caswell, Michael Droettboom, John Hunter, Eric Firing, Antony Lee, Jody Klymak, David Stansby, Elliott Sales De Andrade, Jens Hedegaard Nielsen, Nelle Varoquaux, Benjamin Root, Tim Hoffmann, Phil Elson, Ryan May, Darren Dale, Jae-Joon Lee, Jouni K. Seppänen, Damon McDougall, Andrew Straw, Paul Hobson, Christoph Gohlke, Tony S Yu, Eric Ma, Adrien F. Vincent, Steven Silvester, Charlie Moad, Jan Katins, Nikita Kniazev, Federico Ariza, and Elan Ernest. matplotlib/matplotlib v3.1.0, May 2019.
- [4] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. *CoRR*, abs/1605.08803, 2016.
- [5] Sam Foreman, Xiao-Yong Jin, and James C. Osborn. Deep Learning Hamiltonian Monte Carlo. In *9th International Conference on Learning Representations*, 5 2021.
- [6] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020.
- [7] Gurtej Kanwar, Michael S. Albergo, Denis Boyda, Kyle Cranmer, Daniel C. Hackett, Sébastien Racanière, Danilo Jimenez Rezende, and Phiala E. Shanahan. Equivariant flow-based sampling for lattice gauge theory. *Phys. Rev. Lett.*, 125(12):121601, 2020.
- [8] Daniel Levy, Matthew D. Hoffman, and Jascha Sohl-Dickstein. Generalizing Hamiltonian Monte Carlo with Neural Networks. *arXiv e-prints*, page arXiv:1711.09268, November 2017.
- [9] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,

Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.

- [10] Kirill Neklyudov and Max Welling. Orbital MCMC. *arXiv e-prints*, page arXiv:2010.08047, October 2020.
- [11] Kirill Neklyudov, Max Welling, Evgenii Egorov, and Dmitry Vetrov. Involutive MCMC: a Unifying Framework. *arXiv e-prints*, page arXiv:2006.16653, June 2020.
- [12] Fernando Perez and Brian E. Granger. IPython: A System for Interactive Scientific Computing. *Computing in Science and Engineering*, 9(3):21–29, January 2007.
- [13] Guido Van Rossum and Fred L Drake Jr. *Python Tutorial*. Centrum voor Wiskunde en Informatica Amsterdam, The Netherlands, 1995.
- [14] Antoine Wehenkel and Gilles Louppe. You say Normalizing Flows I see Bayesian Networks. *arXiv e-prints*, page arXiv:2006.00866, June 2020.