

# LeapfrogLayers: A Trainable Framework for Effective Topological Sampling

---

**Sam Foreman,<sup>a,\*</sup> Xiao-Yong Jin<sup>a,b</sup> and James C. Osborn<sup>a,b</sup>**

<sup>a</sup>*Leadership Computing Facility, Argonne National Laboratory,  
Lemont, IL, USA*

<sup>b</sup>*Computational Science Division, Argonne National Laboratory,  
Lemont, IL, USA*

*E-mail:* [foremans@anl.gov](mailto:foremans@anl.gov), [xjin@anl.gov](mailto:xjin@anl.gov), [osborn@alcf.anl.gov](mailto: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.

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

---

\*Speaker

## 1. Introduction

Background

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

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

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

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

where  $\tau_O^{\text{int}}$  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

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

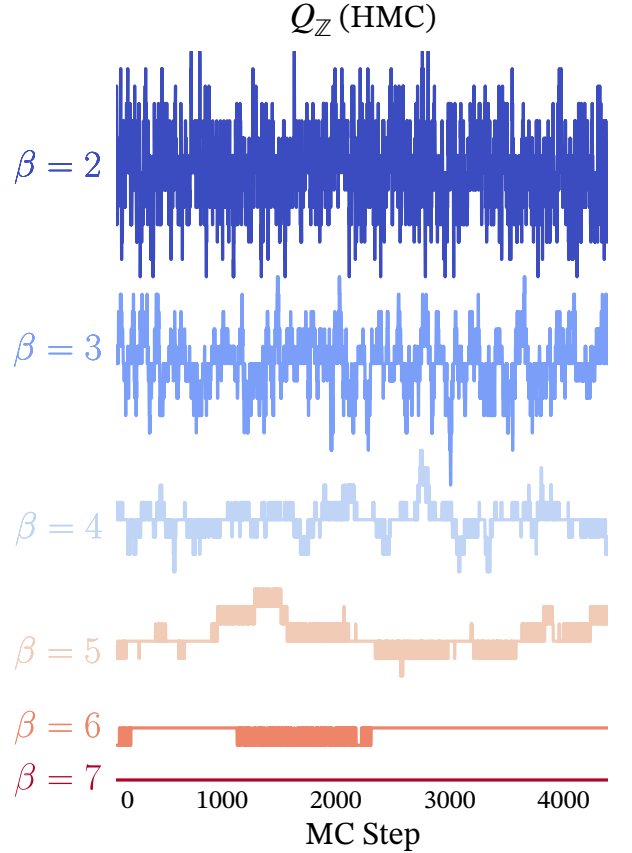
$$x_P \equiv x_\mu(n) + x_\nu(n + \hat{\mu}) - x_\mu(n + \hat{\nu}) - x_\nu(n), \quad (5)$$

defined in terms of the sum of the gauge variables around the elementary plaquette starting at site  $n$  of the lattice. For a given lattice configuration, we can calculate the topological charge  $Q \in \mathbb{Z}$  using

$$Q_{\mathbb{Z}} = \frac{1}{2\pi} \sum_P [x_P], \quad \text{where} \quad (6)$$

$$[x_P] \equiv x_P - 2\pi \left\lfloor \frac{x_P + \pi}{2\pi} \right\rfloor. \quad (7)$$

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.



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

This freezing can be quantatively measured by introducing the *tunneling rate*<sup>1</sup>

$$\delta Q_{\mathbb{Z},i} \equiv |Q_{\mathbb{Z},i+1} - Q_{\mathbb{Z},i}| \in \mathbb{Z}, \quad (8)$$

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

$$p(\mathbf{x}, \mathbf{v}) = p(\mathbf{x})p(\mathbf{v}) = e^{-S_\beta(\mathbf{x})} e^{-\mathbf{v}^2/2} = e^{-\mathcal{H}(\mathbf{x}, \mathbf{v})} \quad (9)$$

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

$$\dot{\mathbf{x}} = \frac{\partial \mathcal{H}}{\partial \mathbf{v}}, \quad \dot{\mathbf{v}} = -\frac{\partial \mathcal{H}}{\partial \mathbf{x}} \quad (10)$$

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*  $\xi^*$  by integrating  $\dot{\xi}$  along  $\mathcal{H} = \text{const.}$  for  $N$  leapfrog steps. i.e.

$$\xi \rightarrow \xi_1 \rightarrow \dots \rightarrow \xi_N \equiv \xi^*, \quad (11)$$

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

$$\text{(a.) } \mathbf{v}' \leftarrow \mathbf{v} - \frac{\varepsilon}{2} \partial_{\mathbf{x}} S(\mathbf{x}), \quad \text{(b.) } \mathbf{x}' \leftarrow \mathbf{x} + \varepsilon \mathbf{v}', \quad \text{(c.) } \mathbf{v}'' \leftarrow \mathbf{v}' - \frac{\varepsilon}{2} \partial_{\mathbf{x}} S(\mathbf{x}). \quad (12)$$

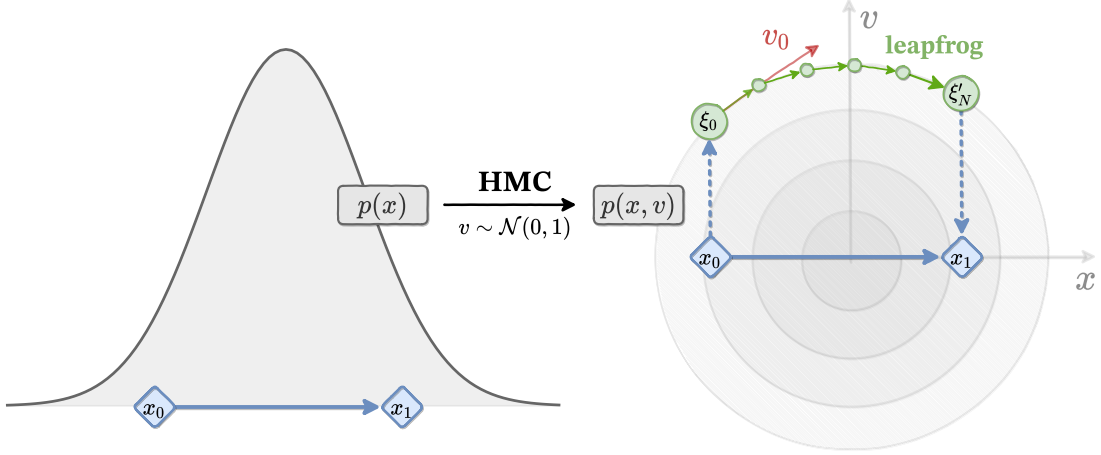
3. At the end of the trajectory, accept or reject the proposal configuration  $\xi^*$  using the MH test

$$\mathbf{x}_{i+1} \leftarrow \begin{cases} \mathbf{x}^* & \text{with probability } A(\xi^*|\xi) \equiv \min \left\{ 1, \frac{p(\xi^*)}{p(\xi)} \right\} \\ \mathbf{x}_i & \text{with probability } 1 - A(\xi^*|\xi). \end{cases} \quad (13)$$

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

---

<sup>1</sup>As 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,  $\Gamma$  ( $\Lambda$ ) 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 [4, 6], we introduce  $d \sim \mathcal{U}(+, -)$ , distributed independently of both  $\mathbf{x}$  and  $\mathbf{v}$ , to determine the “direction” of our update<sup>2</sup>. 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 as<sup>3</sup>

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

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  $\Lambda^\pm$  ( $\Gamma^\pm$ ).

<sup>2</sup>As indicated by the superscript  $\pm$  on  $\Gamma^\pm$ ,  $\Lambda^\pm$  in the update functions.

<sup>3</sup>Here 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  $\mathbb{1} = \mathbf{m} + \bar{\mathbf{m}}$ .

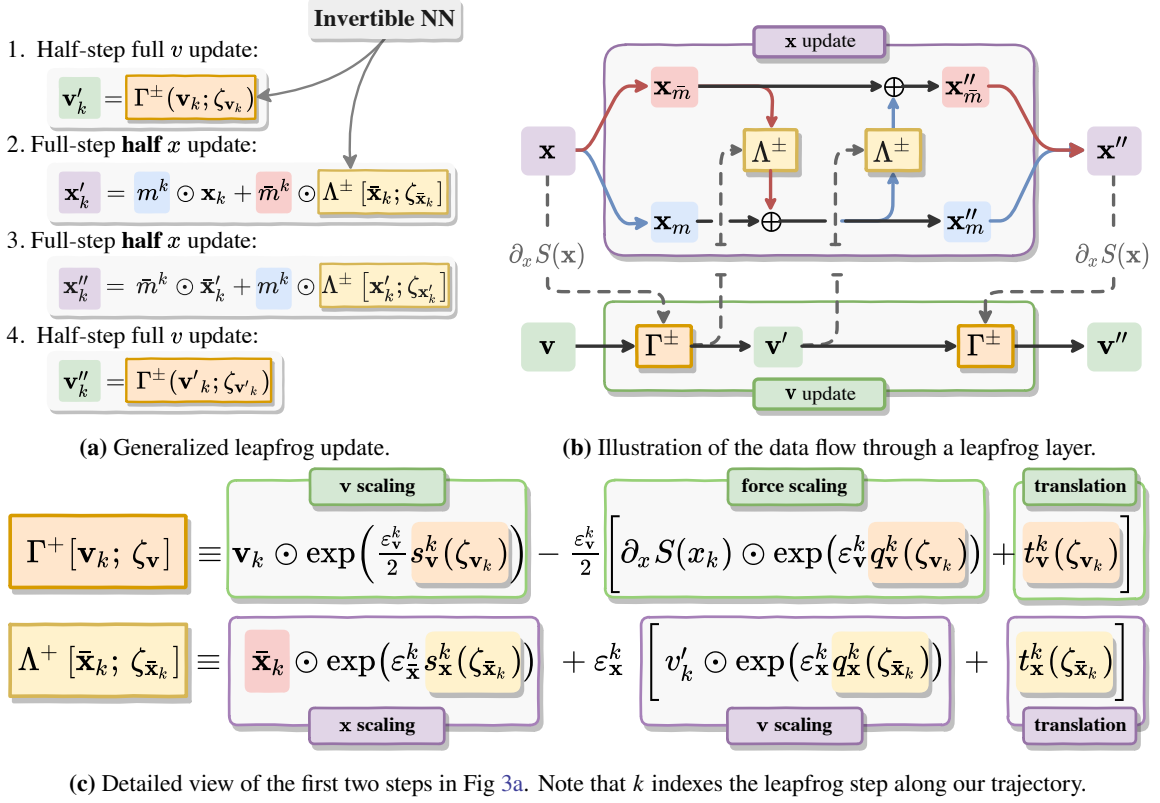


Figure 3: Various illustrations of the generalized leapfrog update.

### 3.1 Network Details

Normalizing flows [3] 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 [1–9].

Note that we maintain separate networks  $\Gamma, \Lambda$  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,  $\Lambda^+[\mathbf{x}_{\bar{m}}; \zeta_{\bar{\mathbf{x}}}]$ <sup>4</sup>. For simplicity, we describe the data flow through a single leapfrog layer, which takes as input  $\zeta_{\bar{\mathbf{x}}} = (\mathbf{x}_{\bar{m}}, \mathbf{v})$ . For the 2D  $U(1)$  model, the gauge links are encoded as  $[\cos(x), \sin(x)]$  for  $x \in [-\pi, \pi]$ . Explicitly<sup>5</sup>,

$$h_1 = \sigma(w_x \mathbf{x} + w_v \mathbf{v} + b_1) \in \mathbb{R}^{n_1} \quad (14)$$

$$h_2 = \sigma(w_2 h_1 + b_2) \in \mathbb{R}^{n_2} \quad (15)$$

$\vdots$

$$h_k = \sigma(w_k h_{k-1} + b_{k-1}) \in \mathbb{R}^{n_k} \implies \quad (16)$$

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

<sup>4</sup>To get the update for the  $d = -$  direction, we invert the update functions and run them in the opposite order.

<sup>5</sup>Here  $\sigma(\cdot)$  is a generic nonlinear activation function acting independently on the elements of the output vector.

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

### 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 ( $\xi$ ) and proposal ( $\xi^*$ ) configuration generated by the sampler. Explicitly, we define

$$\mathcal{L}_\theta(\xi^*, \xi) = A(\xi^*|\xi)(Q_\mathbb{R}^* - Q_\mathbb{R})^2 \quad (17)$$

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

$$A(\xi^*|\xi) = \min \left\{ 1, \frac{p(\xi^*)}{p(\xi)} \mathcal{J}(\xi^*, \xi) \right\} \quad (18)$$

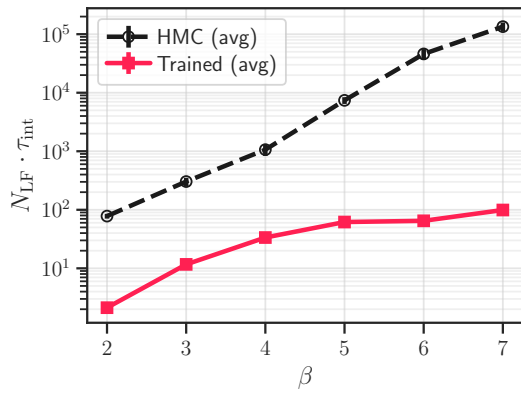
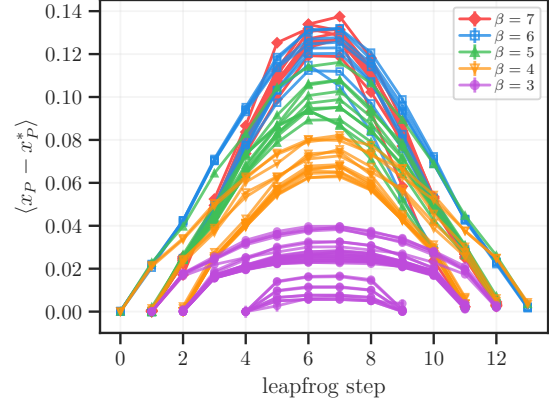
now includes an extra Jacobian factor  $\mathcal{J}(\xi^*, \xi)$  which allows the inclusion of non-symplectic update steps. The Jacobian comes from the  $v(x)$  scaling term in the  $v(x)$  update, and is easily calculated. For completeness, the details of an individual training step are summarized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

#### 3.2.1 Training Step

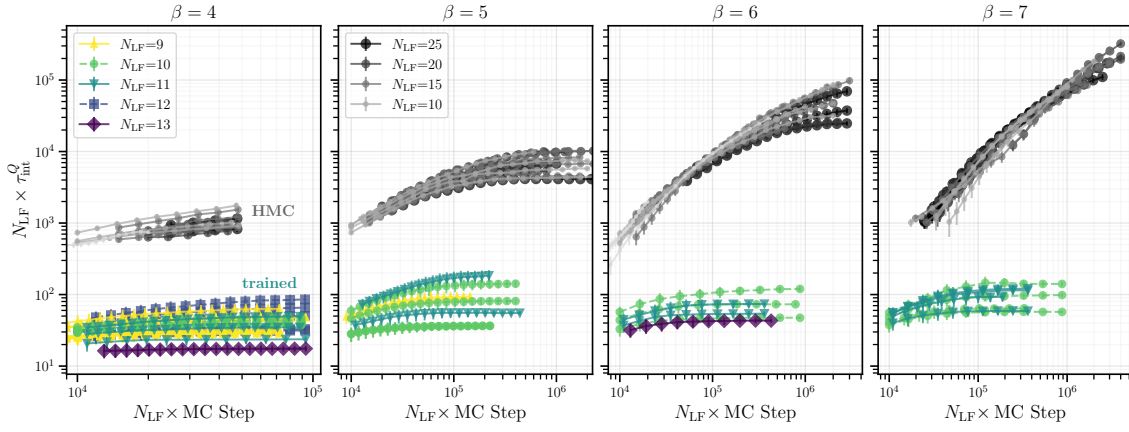
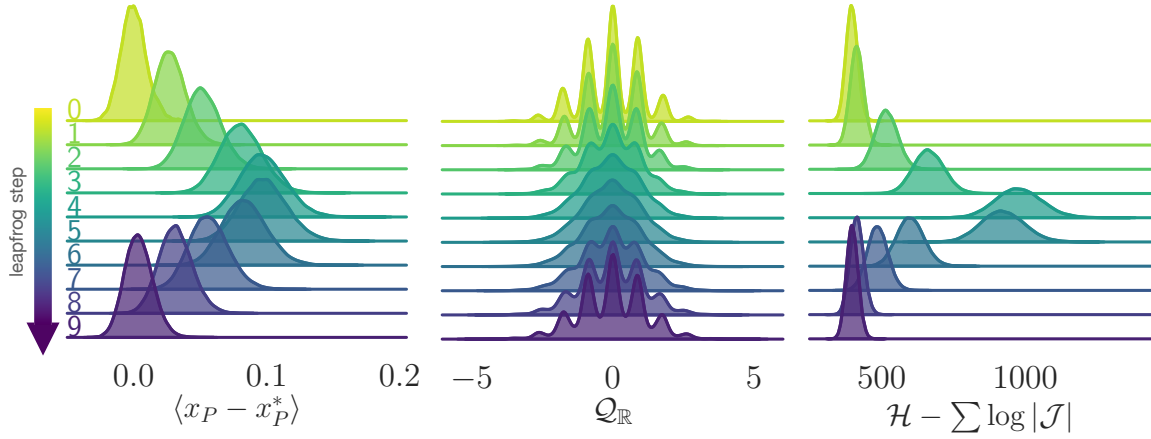
1. Resample  $\mathbf{v} \sim \mathcal{N}(0, \mathbb{I})$ ,  $d \sim \mathcal{U}(+, -)$ , and construct initial state  $\xi = (\mathbf{x}, \mathbf{v}, \pm)$
2. Generate the proposal configuration  $\xi^*$  by passing the initial state sequentially through  $N_{\text{LF}}$  *leapfrog layers* generate a proposal configuration:  $\xi \rightarrow \xi_1 \rightarrow \dots \rightarrow \xi_{N_{\text{LF}}} = \xi^*$
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
 
$$\mathbf{x}_{t+1} \leftarrow \begin{cases} \mathbf{x}^* & \text{with prob. } A(\xi^*|\xi) \\ \mathbf{x} & \text{with prob. } 1 - A(\xi^*|\xi). \end{cases}$$

#### 3.2.2 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,  $\gamma_t$ , 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.

(a) The integrated autocorrelation time  $\tau_{\text{int}}^Q$  vs  $\beta$ 

(b) Deviation in the average plaquette vs leapfrog step.

(c) Integrated autocorrelation time for trained models with different trajectory lengths  $N_{\text{LF}}$ , at  $\beta = 2, 3, \dots, 7$ 

(d) Illustration of how different quantities change during the generalized leapfrog update, the variation in the average plaquette (left), the real-valued topological charge (middle) and the effective energy (right).

**Figure 4:** Plots of various lattice observables.

## 4. Results

We can measure the performance of our approach by comparing the integrated autocorrelation time against traditional HMC. We see in Fig 4a 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 Fig 4d.

We see from the plot of  $\mathcal{H} - \sum \log |\mathcal{J}|$  in Fig 4d (right) that the trained sampler artificially increases the energy towards the middle of the trajectory. This is analogous to reducing the coupling  $\beta$  during the trajectory, as can be seen<sup>6</sup> in Fig 4b. By looking at the variation in the average plaquette<sup>7</sup>  $\langle x_P - x_P^* \rangle$  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.

## 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.

## 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] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. *CoRR*, abs/1605.08803, 2016.
- [4] Sam Foreman, Xiao-Yong Jin, and James C. Osborn. Deep Learning Hamiltonian Monte Carlo. In *9th International Conference on Learning Representations*, 5 2021.
- [5] 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.

<sup>6</sup>Note the trajectories have been shifted to be centered around their midpoints.

<sup>7</sup>Here  $x_P^*$  denotes the expected value in the infinite volume limit.



- [6] 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.
- [7] Kirill Neklyudov and Max Welling. Orbital MCMC. *arXiv e-prints*, page arXiv:2010.08047, October 2020.
- [8] Kirill Neklyudov, Max Welling, Evgenii Egorov, and Dmitry Vetrov. Involutive MCMC: a Unifying Framework. *arXiv e-prints*, page arXiv:2006.16653, June 2020.
- [9] Antoine Wehenkel and Gilles Louppe. You say Normalizing Flows I see Bayesian Networks. *arXiv e-prints*, page arXiv:2006.00866, June 2020.