

# LEARNING DEFORMATION TRAJECTORIES OF BOLTZMANN DENSITIES

Bálint Máté & François Fleuret

University of Geneva

{balint.mate, francois.fleuret}@unige.ch

## ABSTRACT

We introduce a training objective for continuous normalizing flows that can be used in the absence of samples but in the presence of an energy function. Our method relies on either a prescribed or a learnt interpolation  $f_t$  of energy functions between the target energy  $f_1$  and the energy function of a generalized Gaussian  $f_0(x) = \|x/\sigma\|_p^p$ . The interpolation of energy functions induces an interpolation of Boltzmann densities  $p_t \propto e^{-f_t}$  and we aim to find a time-dependent vector field  $V_t$  that transports samples along the family  $p_t$  of densities. The condition of transporting samples along the family  $p_t$  can be translated to a PDE between  $V_t$  and  $f_t$  and we optimize  $V_t$  and  $f_t$  to satisfy this PDE.

## 1 INTRODUCTION

We consider the task of estimating the expectation value  $\mathbb{E}_{x \sim p}[\mathcal{O}(x)]$  of some observable  $\mathcal{O}$ , under a probability density  $p$  proportional to the unnormalized density  $e^{-f}$ , where  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a given energy function. In particular, we don't have access to true samples from  $p$ , all we have is the ability to evaluate  $f$  and its derivatives for any  $x \in \mathbb{R}^n$ . A popular technique (Boyda et al., 2021; Albergo et al., 2021a;b; 2022; Abbott et al., 2022; de Haan et al., 2021; Gerdes et al., 2022; Noé et al., 2018; Köhler et al., 2020; Nicoli et al., 2020; 2021) for attacking this problem is to use a normalizing flow to parametrize a variational density  $q_\theta$  and optimize the parameters  $\theta$  to minimize the reverse KL-divergence

$$KL[q_\theta, p] = \mathbb{E}_{x \sim q_\theta} [\log q_\theta(x) - \log p(x)] = \mathbb{E}_{x \sim q_\theta} [\log q_\theta(x) + f(x)] + Z. \quad (1)$$

The use of normalizing flows for this problem is particularly attractive because  $q_\theta$  can be used as a proposal for importance sampling,  $\mathbb{E}_{x \sim p}[\mathcal{O}(x)] = \mathbb{E}_{x \sim q_\theta} \left[ \frac{p(x)}{q_\theta(x)} \mathcal{O}(x) \right]$ , to account for the inaccuracies of  $q_\theta$ .

Unfortunately, the reverse KL-divergence is mode-seeking, making the training prone to mode-collapse (Fig. 1). We propose an alternative training objective based on infinitesimal deformations of Boltzmann densities (Pfau & Rezende, 2020; Máté & Fleuret, 2022). The contributions of this work can be summarised as follows.

- In §3 we describe our method which relies on either a prescribed or a learnt interpolation  $f_t$  of energy functions between the target energy  $f_1$  and the energy function of a generalized Gaussian  $f_0(x) = \|x/\sigma\|_p^p$ . Given  $f_t$  we optimize a vector field  $V_t$  to transport samples along the family  $p_t(x) \propto e^{-f_t(x)}$  of Boltzmann densities. After translating this condition to a PDE between  $V_t$  and  $f_t$  we propose to minimize the amount by which this PDE fails to hold.
- In §4 we run experiments on the Boltzmann density of a quantum particle moving in a double-well potential and report improvements in KL-divergence, effective sample size, and mode coverage.

Consider the following multimodal density

$$\frac{1}{3} \left( 2\mathcal{N}([-8 \quad -8], 1) + \mathcal{N}([4 \quad 4], 1) \right), \quad (2)$$

where  $\mathcal{N}(\mu, \sigma)$  denotes a normal density centered at  $\mu$  with covariance matrix  $\text{diag}(\sigma^2)$ . Fig. 1 shows the mode-collapse of a normalizing flow trained with the reverse KL-divergence on this target.

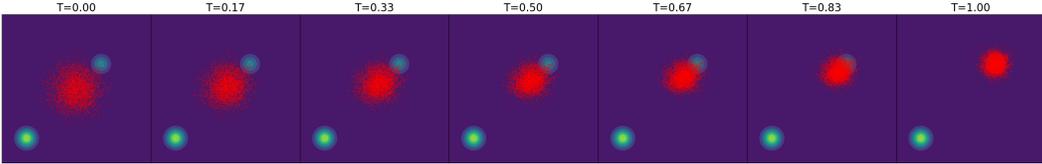


Figure 1: Mode-seeking nature of the reverse KL-divergence. The figures from left to right show how the latent gaussian is transformed by the continuous normalizing flow trained with the reverse KL objective. The green blobs denote the target density (2), a mixture of two Gaussians.

## 2 BACKGROUND

**Boltzmann densities.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be an energy function with a finite normalizing constant  $Z = \int e^{-f(x)} d^n x$ . The function  $f$  then induces a Boltzmann density over the configurations  $x \in \mathbb{R}^n$ ,  $p(x) = \frac{1}{Z} e^{-f(x)}$ . Conversely, given a probability density function  $p : \mathbb{R}^n \rightarrow \mathbb{R}_{+,0}$  the corresponding energy function can be recovered up to a constant  $f = -\log p - \log Z$ .

**The deformation equation.** A probability density  $p_0 \propto e^{-f_0}$  and a one-parameter family of vector field  $V_t$  generates a one-parameter family of probability densities  $p_t \propto e^{-f_t}$  simply by following the integral curves of  $V_t$ . Conversely, given a one-parameter family of probability densities  $p_t$ , we are interested in finding the *transport field* or *deformation field*  $V_t$ . Once  $V_t$  is found, it can be used to transform samples between different members of the family  $p_t$ . Either way, the families  $f_t$  and  $V_t$  are related by

$$\partial_t f_t + \langle \nabla f_t, V_t \rangle - \nabla \cdot V_t + C_t = 0, \tag{3}$$

where  $C_t$  is a spatially constant function, i.e. it only depends on time. We refer the reader to the works of Pfau & Rezende (2020, Eq. 6) and Máté & Fleuret (2022, Eq. 16) for details.

## 3 APPROXIMATING THE TRANSPORT FIELD

From here on, we will use the vector field  $V_t$  and the term “continuous normalizing flow” interchangeably. Our goal is to sample from a target Boltzmann density  $p_{\text{target}} \propto e^{-f_{\text{target}}}$  by

1. defining a family of energy functions  $f_t, 0 \leq t \leq 1$  interpolating between the target energy  $f_1 = f_{\text{target}}$  and the energy function of a generalized Gaussian  $f_0(x) = \|x/\sigma\|_p^p$ ,
2. finding a transport field  $V_t$  such that  $(f_t, V_t)$  “solves” the deformation equation (3).

If we succeed at both of these constructions, then we can obtain samples from  $p_{\text{target}}$  by sampling from  $p_0 \propto e^{-\|x/\sigma\|_p^p}$  and let the samples follow the integral curves of  $V_t$  from  $t = 0$  to  $t = 1$ .

Regarding the second item of the above list, an analytical expression for  $V_t$  is not easy to find if we are given a family of energy functions  $f_t$ . This would amount to solving (3), which is difficult in general. Therefore we will parametrise  $V_t$  with a neural network and train it to minimize the amount by which the pair  $(f_t, V_t)$  fails to satisfy the deformation equation (3). In what follows,  $V_t^\theta$  and  $C_t^\theta$  are parametrized by neural networks and are optimized to minimize some monotonically increasing function  $L^1$  of the pointwise *deformation error*

$$\mathcal{E}_{\theta,x,t} = |\partial_t f_t(x) + \langle \nabla f_t(x), V_t^\theta(x) \rangle - \nabla \cdot V_t^\theta(x) + C_t^\theta|. \tag{4}$$

The expression  $L(\mathcal{E}_{\theta,x,t})$  measures the incompatibility of  $f_t$  and  $V_t$  at a single  $(t, x)$  pair of coordinates, we will need to optimize some sort of integral of this pointwise error over  $t$  and  $x$ .

<sup>1</sup>In our experiments we tried  $L : \mathbb{R}_+ \rightarrow \mathbb{R}_+ \in \{|\mathcal{E}| \mapsto |\mathcal{E}|, |\mathcal{E}| \mapsto |\mathcal{E}|^2, |\mathcal{E}| \mapsto |\mathcal{E}| + |\mathcal{E}|^2, |\mathcal{E}| \mapsto \log(1 + |\mathcal{E}|)\}$ .

**The deformation loss.** Suppose that we have an interpolation of energy functions  $f_t$ . We propose to train  $V_t^\theta$  and  $C_t^\theta$  to minimize the deformation error (4) along the trajectories of  $V_t^\theta$ . Formally, let  $q_\theta$  be a parametric density parametrized by a continuous normalizing  $V_t^\theta$ . We update the parameters to minimize the integral of  $L(\mathcal{E}_{\theta,x,t})$  along the trajectories of the flow,

$$\mathcal{L}(\theta) = \mathbb{E}_{z \sim \mathcal{B}} \left[ \int_0^1 L(\mathcal{E}_{\theta, \gamma_t^\theta(z), t}) dt \right], \tag{5}$$

where  $\mathcal{B}$  denotes the base distribution and  $\gamma_t^\theta(z)$  is given by transporting  $z$  along the  $V_t^\theta$  between 0 and  $\tau$ . The standard deviation of the base is an important hyperparameter, its role can be explained as follows. Ideally, we would like to minimize the deformation error everywhere in space, but we can only evaluate it along the trajectories of  $V_t^\theta$ . To provide better coverage, we can increase the standard deviation of the base density during training.

**Parametrizing the interpolation.** We use a neural network to parametrize the interpolation  $f_t$  as

$$f_t(x) = (1-t)f_0(x) + tf_1(x) + t(1-t)f_t^\theta(x), \quad 0 \leq t \leq 1, \tag{6}$$

where  $f_t^\theta$  is parametrized by a neural network. This parametrization ensures that the boundary conditions at  $t \in \{0, 1\}$  are satisfied, and allows for flexibility for  $0 < t < 1$ . The parameters of the interpolation are trained together with the those of the flow (and those of  $C_t$ ) with the objective of minimizing the deformation loss. Fig. 2 shows that this flexibility allows a flow trained with the deformation loss to capture both modes of the distribution (2).

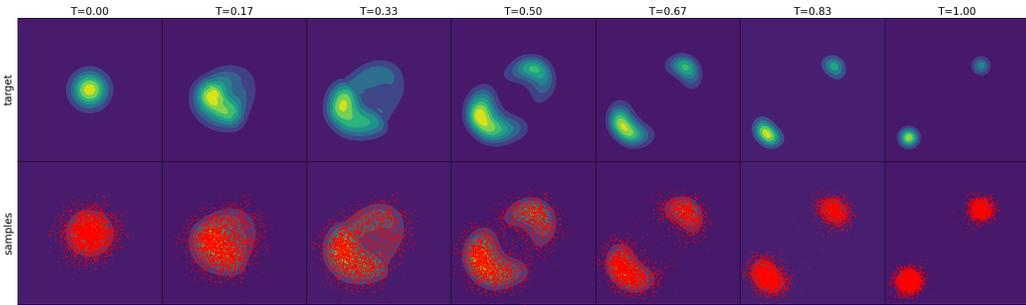


Figure 2: The same continuous normalizing flow as in Fig. 1 trained with the deformation loss using the trainable interpolation. The green blobs denote the (unnormalized) target density (2), a mixture of two Gaussians. The top row shows how the target density evolves under the learned interpolation, while the bottom row shows how the samples from  $q_\theta$  evolve along  $V_t$  as  $t$  is varied.

#### 4 QUANTUM PARTICLE IN A DOUBLE-WELL POTENTIAL

In this section we consider the trajectory of a quantum mechanical particle moving in a double-well potential. The action associated to a discretised trajectory  $(\phi_1, \dots, \phi_N)$  is

$$S(\phi_1, \dots, \phi_N) = \sum_{i=1}^N \left( (\phi_i - \phi_{i+1})^2 + V(\phi_i) \right), \quad V(\phi_i) = -m\phi_i^2 + \lambda\phi_i^4 \tag{7}$$

where  $m$  and  $\lambda$  are numerical parameters and the subscript  $i + 1$  is to be understood modulo  $N$ . The Boltzmann density then reads

$$e^{-S(\phi_1, \dots, \phi_N)} = \prod_{i=1}^N e^{-(\phi_i - \phi_{i+1})^2} \prod_{i=1}^N e^{-V(\phi_i)}. \tag{8}$$

In all our experiments  $N = 16, \lambda = 1/16$  and  $m \in \{0.25, 0.50, 0.75, 1.00\}$ . In all cases, the one-dimensional Boltzmann density  $e^{-V(\phi)}$  has two modes. Their separation is controlled by  $m$ . Intuitively, the second term in (8) encourages every node to follow the unnormalized density  $e^{-V(\phi)}$  at every lattice site, while the first one penalizes neighbors that differ too much (i.e. jump between the modes of  $e^{-V(\phi)}$ ).

**Performance metrics.** To quantify the results of the experiments, for each model we report a subset of the following metrics. For all runs we report the reverse KL-divergence (minus  $\log Z$ ),  $\mathbb{E}_{x \sim q_\theta}(\log q_\theta(x) - \log p(x))$  and the effective sample size,

$$ESS_r = \frac{\left(\frac{1}{N} \sum_i p(x_i)/q_\theta(x_i)\right)^2}{\frac{1}{N} \sum_i (p(x_i)/q_\theta(x_i))^2}, \quad x_i \sim q_\theta, \quad (9)$$

where  $N$  is the number of samples. These metrics capture how good a fit  $q_\theta$  is for  $p$ , but only in those regions from which samples are available. They are therefore insensitive to mode collapse. To compensate for this, we compute the Hausdorff distance between the means of the modes  $M = \{m_1, \dots, m_k\}$  and a batch of  $N$  samples  $X = \{x_1, \dots, x_N\}$  from the model,

$$H(M, X) = \max_{m \in M} \min_{x \in X} \sqrt{\|m - x\|^2}. \quad (10)$$

The means are given  $(\phi_1, \phi_2, \dots, \phi_N) = (a, a, \dots, a)$  and  $(\phi_1, \phi_2, \dots, \phi_N) = (b, b, \dots, b)$  where  $a$  and  $b$  are the two local minima of the double-well potential  $V(\phi) = -m\phi^2 + \lambda\phi^4$ . The Hausdorff distance is a good metric for measuring mode coverage but is insensitive to the shape of the distributions.

**Experiments.** The quantitative results of our runs are summarized in Table 1. In Figure 3 we compare  $e^{-V(\phi)}$  to the histogram of flattened samples from the trained models.

	$m = 0.25$			$m = 0.50$		
	$H(M, X) \downarrow$	Rev. KL $\downarrow$	Rev. ESS $\uparrow$	$H(M, X) \downarrow$	Rev. KL $\downarrow$	Rev. ESS $\uparrow$
$KL(q_\theta, p)$	$1.074 \pm .02$	$-9.777 \pm .00$	$0.984 \pm .00$	$0.879 \pm .01$	$-18.70 \pm .01$	$0.959 \pm .00$
Def. Loss	$1.063 \pm .02$	$-9.781 \pm .00$	$0.996 \pm .00$	$0.864 \pm .01$	$-18.71 \pm .00$	$0.987 \pm .00$
	$m = 0.75$			$m = 1.00$		
	$H(M, X) \downarrow$	Rev. KL $\downarrow$	Rev. ESS $\uparrow$	$H(M, X) \downarrow$	Rev. KL $\downarrow$	Rev. ESS $\uparrow$
$KL(q_\theta, p)$	$16.59 \pm .17$	$-36.18 \pm .00$	$0.958 \pm .01$	$20.13 \pm .09$	$-62.82 \pm .00$	$0.961 \pm .01$
Def. Loss	<b><math>0.755 \pm .02</math></b>	<b><math>-36.88 \pm .00</math></b>	$0.967 \pm .02$	<b><math>0.695 \pm .02</math></b>	<b><math>-63.51 \pm .00</math></b>	$0.963 \pm .01$

Table 1: Results of training the same flow with two different objectives: the reverse KL-divergence and the deformation loss with the trainable interpolation. Mean and standard deviation values over 3 seeds are reported.

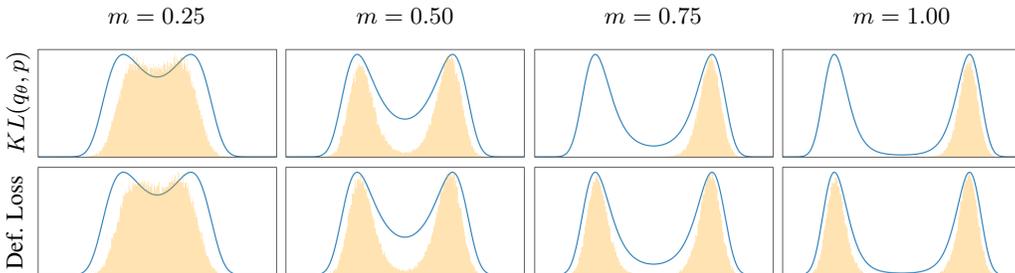


Figure 3: Mode collapse of the reverse KL-divergence for higher values of  $m$ . The unnormalized density  $p(\phi) \propto e^{-V(\phi)}$  (blue curve), compared to flattened samples  $\phi_i$  (orange histogram). Note that these two densities were only supposed to perfectly match, if the  $\phi_i$  at different time-steps were independent of each other, i.e. if (8) didn't involve the term comparing consecutive time-steps. In our setup, these plots can only be used to detect mode-collapse of the flow.

## 5 CONCLUSION

We introduced an alternative training objective of continuous normalizing flows that uses an interpolation of energy functions. We've demonstrated empirically that the proposed objective outperforms the reverse KL-divergence when the target density has multiple modes.

## 6 ACKNOWLEDGEMENT

The authors acknowledge support from the Swiss National Science Foundation under grant number CRSII5\_193716 - “Robust Deep Density Models for High-Energy Particle Physics and Solar Flare Analysis (RODEM)”. We further thank Samuel Klein, Jonas Köhler and Eloi Alonso for discussions.

## REFERENCES

- Ryan Abbott, Michael S. Albergó, Denis Boyda, Kyle Cranmer, Daniel C. Hackett, Gurtej Kanwar, Sébastien Racanière, Danilo J. Rezende, Fernando Romero-López, Phiala E. Shanahan, Betsy Tian, and Julian M. Urban. Gauge-equivariant flow models for sampling in lattice field theories with pseudofermions, 2022. URL <https://arxiv.org/abs/2207.08945>.
- Michael S. Albergó, 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, 2021a. URL <https://arxiv.org/abs/2101.08176>.
- Michael S. Albergó, Gurtej Kanwar, Sébastien Racanière, Danilo J. Rezende, Julian M. Urban, Denis Boyda, Kyle Cranmer, Daniel C. Hackett, and Phiala E. Shanahan. Flow-based sampling for fermionic lattice field theories. *Physical Review D*, 104(11), dec 2021b. doi: 10.1103/physrevd.104.114507. URL <https://doi.org/10.1103/physrevd.104.114507>.
- Michael S. Albergó, Denis Boyda, Kyle Cranmer, Daniel C. Hackett, Gurtej Kanwar, Sébastien Racanière, Danilo J. Rezende, Fernando Romero-López, Phiala E. Shanahan, and Julian M. Urban. Flow-based sampling in the lattice schwinger model at criticality, 2022. URL <https://arxiv.org/abs/2202.11712>.
- Denis Boyda, Gurtej Kanwar, Sébastien Racanière, Danilo Jimenez Rezende, Michael S. Albergó, Kyle Cranmer, Daniel C. Hackett, and Phiala E. Shanahan. Sampling using  $SU(n)$  gauge equivariant flows. *Phys. Rev. D*, 103:074504, Apr 2021. doi: 10.1103/PhysRevD.103.074504. URL <https://link.aps.org/doi/10.1103/PhysRevD.103.074504>.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Pim de Haan, Corrado Rainone, Miranda C. N. Cheng, and Roberto Bondesan. Scaling up machine learning for quantum field theory with equivariant continuous flows, 2021. URL <https://arxiv.org/abs/2110.02673>.
- Mathis Gerdes, Pim de Haan, Corrado Rainone, Roberto Bondesan, and Miranda C. N. Cheng. Learning lattice quantum field theories with equivariant continuous flows, 2022. URL <https://arxiv.org/abs/2207.00283>.
- Jonas Köhler, Leon Klein, and Frank Noé. Equivariant flows: Exact likelihood generative learning for symmetric densities, 2020. URL <https://arxiv.org/abs/2006.02425>.
- Bálint Máté and François Fleuret. Deformations of boltzmann distributions. *arXiv preprint arXiv:2210.13772*, 2022.
- Kim A. Nicoli, Shinichi Nakajima, Nils Strodthoff, Wojciech Samek, Klaus-Robert Müller, and Pan Kessel. Asymptotically unbiased estimation of physical observables with neural samplers. *Physical Review E*, 101(2), feb 2020. doi: 10.1103/physreve.101.023304. URL <https://doi.org/10.1103/physreve.101.023304>.
- Kim A Nicoli, Christopher J Anders, Lena Funcke, Tobias Hartung, Karl Jansen, Pan Kessel, Shinichi Nakajima, and Paolo Stornati. Estimation of thermodynamic observables in lattice field theories with deep generative models. *Physical review letters*, 126(3):032001, 2021.
- Frank Noé, Simon Olsson, Jonas Köhler, and Hao Wu. Boltzmann generators – sampling equilibrium states of many-body systems with deep learning, 2018. URL <https://arxiv.org/abs/1812.01729>.

David Pfau and Danilo Rezende. Integrable nonparametric flows. *arXiv preprint arXiv:2012.02035*, 2020.

## A IMPLEMENTATION AND TRAINING DETAILS

Every trainable object in our experiments is parametrized by a weighted sum of MLPs. The weighting is done by evenly spaced RBF time-kernels, one for each model. Importantly, our architecture is completely oblivious to the  $\mathbb{Z}_2 \times C_n$ -symmetry of the  $\phi^4$  theory and computes the divergence numerically. We leave the exploitation of symmetries and the use of architectures with analytic expressions for the divergence of  $V_t$  (Köhler et al., 2020; Gerdes et al., 2022), as well as for  $\nabla f_t$  and  $\partial_t f_t$ , for future work. The choices of hyperparameters are given in Table 2. Everything was implemented in JAX (Bradbury et al., 2018) and executed on one of eight A100 GPUs.

number of RBF-kernels in time	8
hidden layers per model	3
neurons per hidden layer	128
activation function	swish
base distribution	$\propto e^{-(x/2)^4}$
bath size during training	256
batch size during evaluation	4096
number of train steps	$10^4$
initial learning rate	$10^{-3}$
number of integration steps	50
deformation loss	$ \mathcal{E}_{def}  +  \mathcal{E}_{def} ^2$

Table 2: Hyperparameters and desing choices for our experiments. The learning rate was initialized to the value shown in the table and annealed to 0 following a cosine schedule.

**Computational costs** To optimize the “baseline” reverse KL objective one needs a parametrization of  $V_t$  and to integrate  $\nabla \cdot V_t$  along the trajectories. In addition to this, the deformation loss also needs a parametrization for  $C_t$  and one for  $f_t$  and to integrate an expression depending on  $\nabla \cdot V_t, \partial_t f, \nabla f$  and  $C_t$ . This means that the same number of training steps are more expensive computationally when using the deformation loss. The experiments lasted  $\sim 1.85$ -times longer when using the deformation loss instead of  $KL$ -divergence.