

# Neural Canonical Transformation with Symplectic Flows


Shuo-Hui Li<sup>1,2</sup>, Chen-Xiao Dong<sup>1,2</sup>, Linfeng Zhang<sup>3,\*</sup> and Lei Wang<sup>1,4,†</sup>

<sup>1</sup>*Institute of Physics, Chinese Academy of Sciences, Beijing 100190, China*

<sup>2</sup>*University of Chinese Academy of Sciences, Beijing 100049, China*

<sup>3</sup>*Program in Applied and Computational Mathematics, Princeton University, Princeton, New Jersey 08544, USA*

<sup>4</sup>*Songshan Lake Materials Laboratory, Dongguan, Guangdong 523808, China*

 (Received 13 October 2019; revised manuscript received 22 January 2020; accepted 9 March 2020; published 28 April 2020)

Canonical transformation plays a fundamental role in simplifying and solving classical Hamiltonian systems. Intriguingly, it has a natural correspondence to normalizing flows with a symplectic constraint. Building on this key insight, we design a neural canonical transformation approach to automatically identify independent slow collective variables in general physical systems and natural datasets. We present an efficient implementation of symplectic neural coordinate transformations and two ways to train the model based either on the Hamiltonian function or phase-space samples. The learned model maps physical variables onto an independent representation where collective modes with different frequencies are separated, which can be useful for various downstream tasks such as compression, prediction, control, and sampling. We demonstrate the ability of this method first by analyzing toy problems and then by applying it to real-world problems, such as identifying and interpolating slow collective modes of the alanine dipeptide molecule and MNIST database images.

DOI: [10.1103/PhysRevX.10.021020](https://doi.org/10.1103/PhysRevX.10.021020)

Subject Areas: Computational Physics,  
Nonlinear Dynamics, Statistical Physics

## I. INTRODUCTION

The inherent symplectic structure of classical Hamiltonian mechanics has profound theoretical and practical implications [1]. For example, the symplectic symmetry underlies Liouville's theorem [2], which states that the phase-space density is incompressible under the Hamiltonian evolution. Canonical transformations which preserve the symplectic symmetry in the phase space have been a key technique for simplifying and solving Hamiltonian dynamics. Respecting the intrinsic symplectic symmetry of Hamiltonian systems is also crucial for stable and energy-conserving numerical integration schemes [3] which play central roles in the investigations of celestial mechanics and molecular dynamics.

Molecular dynamics (MD) simulation investigates the dynamical and statistical properties of matter by integrating the equations of motion of a large number of atoms. MD is a vital tool for understanding complex physical, chemical, and biological phenomena, as well as for practical

applications in material discovery and drug design. Modern MD simulation generates huge datasets, which encapsulate the full microscopic details of the molecular system [4]. However, this also poses challenges to the development of data analysis tools. In particular, one typically is interested in the emerging slow modes, which are often related to the collective property of the system. Moreover, identifying such degrees of freedom is also crucial for an enhanced sampling of molecular conformations. See Refs. [5,6] for recent reviews.

Techniques of machine learning provide promising solutions to these problems in MD. For example, the time-lagged independent component analysis [7–11] separates a linear mixture of independent time-series signals. The approach shows a close connection to the dynamic mode decomposition scheme developed in the fluid mechanics community [12,13]. Many of these linear analysis methods have nonlinear generalizations based on kernel approaches [14,15]. More recently, several approaches of deep learning have also been proposed to identify a nonlinear coordinate transformation of dynamical systems [16–20]. Parallel to these efforts, it is also an active research direction to extract slow features in general time-series data [14,21,22] within the community of machine learning.

In this paper, we develop a different approach by exploiting the inherent connection between canonical transformation and normalizing flows [23,24]. We design a class of learnable neural canonical transformations to

\*linfengz@princeton.edu

†wanglei@iphy.ac.cn

*Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.*

simplify complex Hamiltonian dynamics of the physical variables towards independent collective motions in the transformed phase space. Correspondingly, the canonical transformation also reduces complex phase-space densities towards an independent Gaussian prior distribution. After learning, one can directly control nonlinear collective variables with different frequencies by tuning independent collective variables in the latent space of the normalizing flows. We present learning algorithms and discuss applications of the neural canonical transformation on the extraction of slow collective variables of the physical and realistic dataset. We stress that most techniques that target the extraction of dynamical information require time-series data. However, in the present approach, the dynamical information is imposed on the structure of the neural canonical transformation, so that the training scheme does not necessarily follow a specific time step, and the data sample may come from other types of sampling methods, such as Monte Carlo, biased dynamics, etc.

There have been related research works exploiting the symplectic property in tasks of machine learning. Reference [25] solves Hamiltonian equations using neural networks. Reference [26] learns a Hamiltonian dynamics from observed data using neural networks. Both studies found that exploiting the symplectic structure in the learning helps in boosting the performance. More recently, there have been more preprints on related topics [27–31] which also aim at improving the performance in tasks of machine learning by imposing physics-motivated inductive biases in the design of neural network. Our work finds the closest connection to Ref. [32], which investigates classical integrable systems using symplectic neural networks. Our paper targets more general settings and aims to identify nonlinear slow collective modes of complex systems. In addition, we also note that there are efforts on learning neural networks for the force fields of molecular dynamics [33–39], wherein imposing the physical invariance is also crucial.

The organization of the paper is as follows. In Sec. II, we reveal the key connection of canonical transformation to normalizing flows with the symplectic condition. In Sec. III, we present the design and training of the symplectic neural networks for canonical transformation. We also discuss potential applications of the neural canonical transformation. In Sec. IV we demonstrate applications of the neural canonical transformation to toy problems and realistic data. Finally, we discuss possible prospects of neural canonical transformation in Sec. V.

## II. THEORETICAL BACKGROUND

We review the canonical transformation and its connection to the normalizing flow model.

### A. Canonical transformation of Hamiltonian systems

We denote the canonical variables, namely the momenta and coordinates of a Hamiltonian system, as a row vector

with  $2n$  elements  $\mathbf{x} \equiv (\mathbf{p}, \mathbf{q})$ . The Hamiltonian equation can be concisely written as  $\dot{\mathbf{x}} = \nabla_{\mathbf{x}} H(\mathbf{x}) J$ , where  $\dot{\mathbf{x}}$  denotes the time derivative of the canonical variables.  $H(\mathbf{x})$  is the Hamiltonian function and

$$J = \begin{pmatrix} & I \\ -I & \end{pmatrix}$$

is a  $2n \times 2n$  symplectic metric matrix.

The canonical transformation is a bijective mapping from the original canonical variables to a new set of canonical variables, i.e.,  $\mathcal{T} : \mathbf{x} \mapsto \mathbf{z} \equiv (\mathbf{P}, \mathbf{Q})$ , whose Jacobian matrix  $M_{ij} = \nabla_{x_j} z_i$  satisfies the symplectic condition

$$MJM^T = J. \quad (1)$$

Canonical transformation preserves the Hamiltonian equation, i.e., one has  $\dot{\mathbf{z}} = \nabla_{\mathbf{z}} K(\mathbf{z}) J$ , where  $K(\mathbf{z}) = H \circ \mathcal{T}^{-1}(\mathbf{z})$  is a transformed Hamiltonian in terms of the new phase-space variables. The canonical transformation establishes a bijective mapping between the Hamiltonian trajectories in the original and the transformed phase spaces. Thus, one can search for canonical transformations which simplify and even solve the Hamiltonian dynamics. One such searching strategy is to compose elementary canonical transformations since the symplectic condition Eq. (1) forms a group.

### B. Normalizing flow models

Generative modeling aims at modeling the joint distribution of complex high-dimensional data [40]. A generative model can capture the key variations of the dataset and draw samples efficiently from the learned probability distribution. A large class of generative models achieves these goals by learning a transformation from a simple latent distribution to a complex physical distribution. Examples well known to the machine-learning community include the generative adversarial networks (GAN) [41], the variational autoencoders (VAE) model [42], and the normalizing flow models [23,24].

The normalizing flow models, or the flow-based generative models, are particularly suitable for our purpose since they employ a bijective mapping from the latent space  $\mathbf{z}$  to the target space  $\mathbf{x}$  for the probability transformation. Essentially, the normalizing flow models perform a change of variables to induce transformations in the probability densities. Typically, the normalizing flow model transforms a simple prior distribution of the latent variables following, e.g., the Gaussian distribution  $\mathcal{N}(\mathbf{z}; \mathbf{0}, \Sigma)$  with zero mean and covariance  $\Sigma$  to the more complex distribution of the realistic data, and vice versa.

There have been a great variety of normalizing flow models [43–47] that achieved nice performance in applications of machine learning. Compared to GAN and VAEs, the normalizing flow models have appealing features such as

tractable likelihood for any given data  $\mathbf{x}$  and exact reversibility between  $\mathbf{z}$  from  $\mathbf{x}$ . These features are particularly attractive for principled and quantitative scientific applications such as learning renormalization group flow [48], holographic mapping [49], Monte Carlo sampling [50–52], molecular simulations [53], and spin glasses [54]. Interestingly, in the continuous-time limit, the normalizing flow model exhibits intriguing connections to a variety of topics including dynamical systems, ordinary differentiation equations, optimal transport theory, and fluid dynamics [55–57].

### C. Connections between canonical transformation and normalizing flow models

Since canonical transformations are changes of variables in the phase space, they can be naturally parametrized and learned as normalizing flow models with the added symplectic condition. Moreover, it is instructive to consider the probabilistic interpretation of this change of variables. Consider the phase-space density in the canonical ensemble  $\pi(\mathbf{x}) = e^{-\beta H(\mathbf{x})}/Z$ , where  $Z = \int d\mathbf{x} e^{-\beta H(\mathbf{x})}$  is the partition function,  $\beta = k_B T$  is the inverse temperature. The change of variables to a simplified Hamiltonian function  $K(\mathbf{z})$  implies reaching a simplified density in the latent phase space  $e^{-\beta K(\mathbf{z})}/Z$ . Thus, one can regard the canonical transformation as a flow-based generative model connecting the physical and latent phase spaces and the associated phase-space densities. There is, however, a crucial symplectic constraint in Eq. (1) on the transformation compared to ordinary normalizing flow models. From the generative modeling perspective, the additional symplectic condition further restricts the expressibility of the network. On the other hand, the symplectic inductive bias offers a physical guarantee and interpretability on the training results, e.g., the network will always be a canonical transformation that preserves the dynamics in the latent space. Table I summarizes the connection between canonical transformation and normalizing flow.

This key insight has several profound consequences. First, one can search for canonical transformations by learning the flow models in the phase space, which simplifies complex Hamiltonians both in the statistical and in the dynamical senses. Second, the learned latent spaces attain the physical meaning of transformed

TABLE I. The correspondence of a canonical transformation and normalizing flow. See Sec. II C for explanations of the connection.

Canonical transformation	Normalizing flow
$\mathbf{x} = (\mathbf{p}, \mathbf{q})$	Physical variables
$\mathbf{z} = (\mathbf{P}, \mathbf{Q})$	Latent variables
$\mathbf{x} \leftrightarrow \mathbf{z}$	Symplectic flow
$e^{-H(\mathbf{x})}/Z$	Physical phase-space density
$e^{-K(\mathbf{z})}/Z$	Latent phase-space density

canonical variables which can be useful for various downstream tasks. Last, the symplectic neural network necessarily has the volume-preserving property, which can be computationally efficient since the Jacobian determinant is always unity by construction.

## III. CANONICAL TRANSFORMATION USING NORMALIZING FLOW MODELS

It is usually difficult to devise useful canonical transformations for generic Hamiltonians since it typically involves solving a large set of coupled nonlinear equations. However, building on the connections of canonical transformation and normalizing flow models [23,24], we can construct a family of expressive canonical transformations with symplectic neural networks and train them with optimization techniques.

To train the model, one can follow either the variational approach or the data-driven approach. As a result, the neural canonical transformation helps simplify the dynamics and identify nonlinear slow collective variables of complex Hamiltonians.

### A. Model architectures

As a flow-based generative model, the neural canonical transformation consists of a symplectic network and a prior distribution which corresponds to the transformation and the target phase density distribution, respectively. In the most general setting, the canonical transformation can even mix the momenta and coordinates. Here, we restrict ourselves to point transformations [58] for balanced flexibility and interpretability. We list several other possible implementations of the neural canonical transformations in the Appendix A. We note that one can compose symplectic neural networks to form more expressive canonical transformations.

#### 1. Neural point transformations

In the point transformation, one performs a nonlinear transformation to the coordinates  $\mathbf{q}$  and a linear transformation to the momenta  $\mathbf{p}$  accordingly,

$$\mathbf{Q} = \mathcal{F}(\mathbf{q}), \quad (2)$$

$$\mathbf{P} = \mathbf{p}(\nabla_{\mathbf{q}}\mathbf{Q})^{-1} = \mathbf{p}\nabla_{\mathbf{Q}}\mathbf{q}. \quad (3)$$

The overall transformation in the phase space  $\mathcal{T}:\mathbf{x} = (\mathbf{p}, \mathbf{q}) \mapsto \mathbf{z} = (\mathbf{P}, \mathbf{Q})$  satisfies the symplectic condition Eq. (1) [58]. Training of point transformation will involve both momenta and coordinates in the phase space. Since the coordinate transformation Eq. (2) is independent of momenta, we can use the resulting coordinates  $\mathbf{Q}$  alone as a set of collective coordinates.

The coordinate transformation  $\mathcal{F}:\mathbf{q} \mapsto \mathbf{Q}$  in Eq. (2) can be any nonlinear bijective mapping. We implement it with a

real-valued non-volume-preserving (real NVP) network [45], which is a typical normalizing flow model [23,24]. The momenta transformation has the form of a vector-Jacobian product, which is commonly implemented for the reverse mode automatic differentiation [59]. So we can leverage the automatic differentiation mechanism for the momentum transformation. In practice, we run the coordinate transformation first forward and then reverse for  $\mathbf{q} = \mathcal{F}^{-1}(\mathbf{Q})$ . Then, we compute its inner product with the initial momenta  $\mathbf{p} \cdot \mathbf{q}$ . Finally, we compute the derivative of the scalar with respect to  $\mathbf{Q}$  to obtain the transformed momenta  $\mathbf{P}$  in Eq. (3).

## 2. Latent-space Hamiltonian and prior distribution

We assume the transformed Hamiltonian in the latent space has the simple form of an independent harmonic oscillator

$$K(\mathbf{z}) = \sum_{k=1}^n \frac{P_k^2 + \omega_k^2 Q_k^2}{2}, \quad (4)$$

where  $\omega_k$  are learnable frequencies for each pair of conjugated canonical variables. Without loss of generality, we set the inverse temperature in the latent space to be one. Therefore, in terms of canonical density distribution, the prior density in the latent space is an independent Gaussian  $\mathcal{N}(\mathbf{z}; \mathbf{0}, \Sigma) = [(\prod_{k=1}^n \omega_k) / (2\pi)^n] e^{-K(\mathbf{z})}$ , where  $\Sigma = \text{diag}(\underbrace{1, \dots, 1}_n, \underbrace{\omega_1^{-2}, \dots, \omega_n^{-2}}_n)$  is a diagonal covariance matrix. In this setting, each pair of canonical variables in the latent space corresponds to an independent collective mode with a learnable frequency. Thus, after training one can select a desired number of slow modes according to the frequencies.

## B. Training approaches

The principle for training is to match the phase-space density of the generative model  $\rho$  to the target density  $\pi$ . Depending on specific applications one may either have direct access to the Hamiltonian function or the samples from the target distribution. Thus, we devised two training schemes of the neural canonical transformation based on variational calculation and the data-driven approach, respectively.

### 1. Variational approach

We can learn the canonical transformation based on the analytical expression of the physical Hamiltonian. For this purpose, we minimize the variational free energy

$$\mathcal{L} = \int d\mathbf{x} \rho(\mathbf{x}) [\ln \rho(\mathbf{x}) + \beta H(\mathbf{x})]. \quad (5)$$

This objective function is upper bounded by the free energy since  $\mathcal{L} + \ln Z = \mathbb{K}\mathbb{L}(\rho \parallel \pi) \geq 0$ , where the

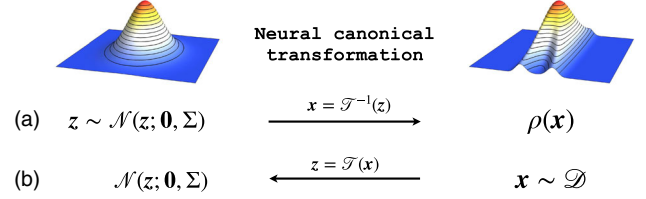


FIG. 1. A neural canonical transformation maps between the latent variables  $\mathbf{z}$  and physical variables  $\mathbf{x}$  via a symplectic neural network. The transformation preserves the Hamiltonian equation and connects the phase-space trajectories in the physical and the latent spaces. There are two ways to train the neural canonical transformation: (a) variational free energy based on the Hamiltonian [Eq. (5)]. (b) density of phase-space estimation based on data [Eq. (6)].

Kullback-Leibler (KL) divergence is a non-negative measure of the dissimilarity between the model and the target distributions. The equality is reached only when two distributions match each other. The objective function of this form was recently employed in the probability density distillation of generative models [60].

To evaluate Eq. (5) we first draw samples from the normal distribution and then scale them according to the frequencies in the prior distribution to obtain  $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \Sigma)$ . Next, we pass the samples through the symplectic transformation to obtain  $\mathbf{x} = \mathcal{T}^{-1}(\mathbf{z})$ , as shown in Fig. 1(a). Since the symplectic transformation is volume preserving, the probability density of the produced samples reads  $\rho(\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \Sigma)$ . The objective function is estimated on these samples as  $\mathcal{L} = \mathbb{E}_{\mathbf{x} \sim \rho(\mathbf{x})} [\ln \rho(\mathbf{x}) + \beta H(\mathbf{x})]$ . To minimize the objective function we compute the gradient over such a sampling procedure with the reparametrization trick [42], which is an unbiased and a low variance gradient estimator for the learnable parameters [61].

### 2. Maximum likelihood estimation

Alternatively, one can also learn the canonical transformation in a purely data-driven approach. Assuming one already has access to independent and identically distributed samples from the target distribution  $\pi(\mathbf{x})$ , one can learn the neural canonical transformation with the maximum likelihood estimation on the data. This amounts to performing the density estimation in the phase space with a flow-based probabilistic generative model. The goal is to minimize the negative log-likelihood (NLL) on the dataset  $\mathcal{D} = \{\mathbf{x}\}$

$$\text{NLL} = -\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\ln \rho(\mathbf{x})], \quad (6)$$

which reduces the observed phase-space density and the model density  $\mathbb{K}\mathbb{L}(\pi \parallel \rho)$  based on empirical observations. To train the network we run the transformation from the physical to latent space as shown in Fig. 1(b) and compute the model density  $\rho(\mathbf{x}) = \mathcal{N}(\mathbf{z} = \mathcal{T}(\mathbf{x}); \mathbf{0}, \Sigma)$ .



The density estimation Eq. (6) requires the phase-space data, which involves both the coordinates and the momenta information. This appears to pose difficulties for applications to MD data which typically only contain the trajectory in the coordinate space. Fortunately, the momenta and coordinate distribution are factorized for separable Hamiltonians encounters in most MD simulations. That is, the momenta follow an independent Gaussian distribution whose variances depend on the atom masses and temperature. Therefore, one can exploit this fact and augment the training dataset by sampling momenta data directly from a Gaussian distribution.

### C. Applications

Neural canonical transformation learns a latent representation with independent modes and simplified dynamics. In principle, the learned representation is useful for the prediction, control, and sampling of the original system. We list a few concrete applications below.

#### 1. Thermodynamics and excitation spectra

Since the training approach of Sec. III B 1 satisfies the variational principle, the loss function Eq. (6) provides an upper bound to the physical free energy of the system. Besides, one can also estimate entropy and free-energy differences of the Hamiltonian with different parameters. Similar variational free-energy calculation of statistical mechanics problems using deep generative models have been carried out recently in Refs. [48,49,52–54,57,62]. In particular, Ref. [53] has obtained encouraging results for sampling equilibrium molecular configurations. The present approach differs since it works in the phase space which involves both momenta and coordinates, which allows extract dynamical information in addition to statistical properties.

Since the neural canonical transformation preserves the Hamiltonian dynamics of the system, the learned frequencies in the latent space Eq. (4) reflect the intrinsic timescale of the target problem. In this way, the present approach captures coherent excitation of the system in the latent-space harmonic motion. One may also estimate the spectral density based on the learned frequencies.

#### 2. Identifying collective variables from slow modes

Since the neural canonical transformation automatically separates dynamical modes with different frequencies in the latent space, one can extract nonlinear slow modes of the original physical system by selecting the latent variables with small frequencies.

The neural canonical transformation differs fundamentally with these general time-series analysis approaches [7,10–13] which do not exploit the domain-specific symplectic symmetry of the Hamiltonian systems. Another fundamental difference is that the canonical transformation

is performed in the phase space which contains both momenta and coordinates, rather than for the time sequence of the coordinates. Since the explicit time information was never used in the present approach, there is no need to choose the time lag hyperparameter as in the time-lagged independent component analysis and related approaches. Last, variational training of the transformation also allows one to identify the canonical transformation directly from the microscopic Hamiltonian even without the time-series data.

We note that in practice, exact dynamical information is usually lost when one cares only about the structural, or static, information of a system. Sophisticated thermostating and enhanced sampling techniques are used to accelerate the sampling, but, meanwhile, the dynamics is destroyed. In this case, MD plays the role of a sampler, rather than a real-time simulator, yet dynamical information can be extracted from statistical data of Hamiltonian systems with the neural canonical transformation approach.

## IV. EXAMPLES

We demonstrate the application of the neural canonical transformation with concrete examples. We start from simple toy problems and then move on to more challenging realistic problems. In all examples, the trainable parts of the network are the coordinate transformation  $\mathcal{F}$  [Eq. (2)] and the latent-space frequencies [Eq. (4)]. The code implementation is publicly available at [63].

### A. Ringworld

First, we consider a two-dimensional toy problem with the Hamiltonian  $H = \frac{1}{2}(p_1^2 + p_2^2) + (\sqrt{q_1^2 + q_2^2} - 2)^2/0.32$  [50]. Canonical distribution of this Hamiltonian resides in a four-dimensional phase space. The canonical ensemble samples from  $\pi(\mathbf{x}) = e^{-H(\mathbf{x})}/Z$  are confined in a manifold embedded in the phase space due to the potential term. In the Euclidean space, the coordinates are correlated.

Taking the Hamiltonian and training it with the variational approach, we obtain a neural point transformation from the original variables to a new set of canonical variables. Figure 2 shows the samples projected onto the plane of latent coordinates  $Q_k$  and the polar coordinate variables  $\varphi = \arctan(q_2/q_1)$ ,  $r = \sqrt{q_1^2 + q_2^2}$ . One observes a significant correlation between the slowest variable  $Q_1$  and the polar angle  $\varphi$ , while the other transformed coordinate  $Q_2$  shows a strong correlation with the radius variable  $r$ .

This example demonstrates that, as a bottom line, the neural point transformation can automatically identify nonlinear transformation (such as polar coordinates) of the original coordinates. In the learned representation, the dynamics of each degree of freedom becomes independent.

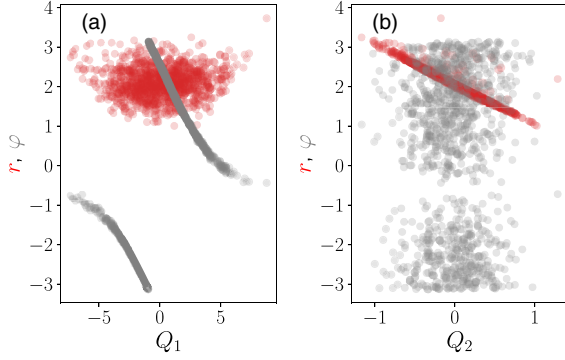


FIG. 2. The Ringworld samples of Sec. IV A projected to the plane of learned latent coordinates and the polar coordinates.

### B. Harmonic chain

Next, we consider a harmonic chain with the Hamiltonian  $H = \frac{1}{2} \sum_{i=1}^n [p_i^2 + (q_i - q_{i-1})^2]$ . We set  $q_0 = q_{n+1} = 0$  to fix both ends of the chain. The system can be readily diagonalized by finding the normal mode representation  $H = \frac{1}{2} \sum_{k=1}^n (\dot{Q}_k^2 + \omega_k^2 Q_k^2)$ , where  $\omega_k = 2 \sin[\pi k / (2n + 2)]$  is the normal mode frequency and  $Q_k = \sqrt{2 / (n + 1)} \sum_{i=1}^n q_i \sin[ik\pi / (n + 1)]$  is the normal coordinate.

We train a neural point transformation with the variational loss Eq. (5) at the inverse temperature  $\beta = 1$ . Figure 3(a) shows learned frequencies in the latent-space harmonic Hamiltonian Eq. (4) together with the analytical dispersion. The agreement is particularly good for the low frequencies which are populated by the canonical distribution at the given temperature. Moreover, we pick the two slowest coordinates  $Q_k$  and compute their Jacobians with respect to the physical variables  $q_i$  as shown in Fig. 3(b). The comparison shows that the neural canonical transformation nicely identifies slow collective modes of the system based on its Hamiltonian. On the other hand, the model is also able to learn these slow modes from data.

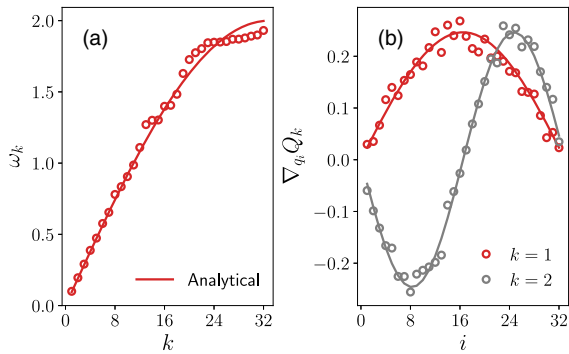


FIG. 3. (a) The learned frequencies of harmonic chain present in Sec. IV B in the prior distribution [Eq. (4)] and the analytical normal-mode frequency. (b) The Jacobian  $\nabla_{q_i} Q_k$  of the two slowest collective coordinates with respect to the original coordinates. Solid lines are the analytical solution.

this simple case, modes with small frequency correspond to latent variables with large covariance, which could also be captured by the principal component analysis [64].

Having demonstrated that the neural point transformation reproduces conventional normal-mode analysis and principal component analysis for the harmonic chain, we move on to show the major strength of the present approach in extracting nonlinear slow modes.

### C. Alanine dipeptide

Proteins show rich dynamics with multiple emergent timescales. As one of the protein's building blocks, the alanine dipeptide is a standard benchmark problem. Despite being a small organic molecule, the alanine dipeptide shows nontrivial dynamics that deserve study. The backbone of alanine dipeptide contains ten heavy atoms with the simplified molecular-input line-entry system representation CC(=O)NC(C)C(=O)NC. It is known that the two dihedral angles  $\Phi$  and  $\Psi$  which control the torsion of the molecule, as indicated in the inset of Fig. 4(a), are the key degrees of freedom which show slow dynamics.

Here, we train a neural canonical transformation to identify the slow modes of the alanine dipeptide molecule based on raw MD simulation data. For the training, we use the MD dataset [17,65] released at [66]. The dataset consists of 250 ns of Euclidean space trajectory of the ten heavy atoms in the alanine dipeptide at 300 K with an integration step of 2 fs. Since the density estimation Eq. (6) requires the phase-space data, we extend atom coordinates data to the phase space by sampling momenta from the Gaussian distribution whose variances depend on the atom masses and temperature. Note that for the phase-space density estimation we randomly shuffle the trajectory data, thus we do not use any of the time frame information in the training. We use the three MD independent trajectories at [66] for the training, validation, and testing, respectively. Each of them contains 250000 snapshots. We use the Adam

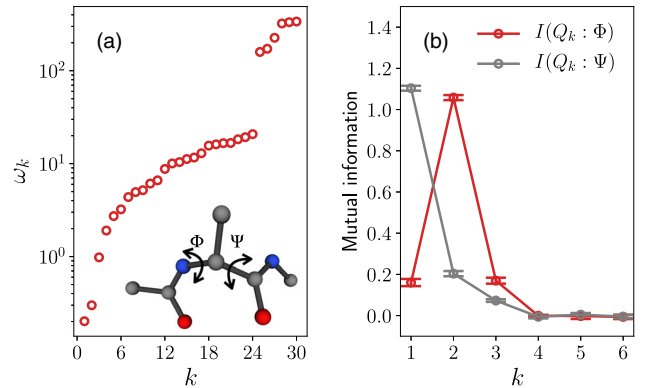


FIG. 4. (a) The learned frequencies of alanine dipeptide presented in Sec. IV C from the MD trajectories. The inset shows the molecule with slow torsion angles. (b) The mutual information between the few slowest modes and the torsion angles.

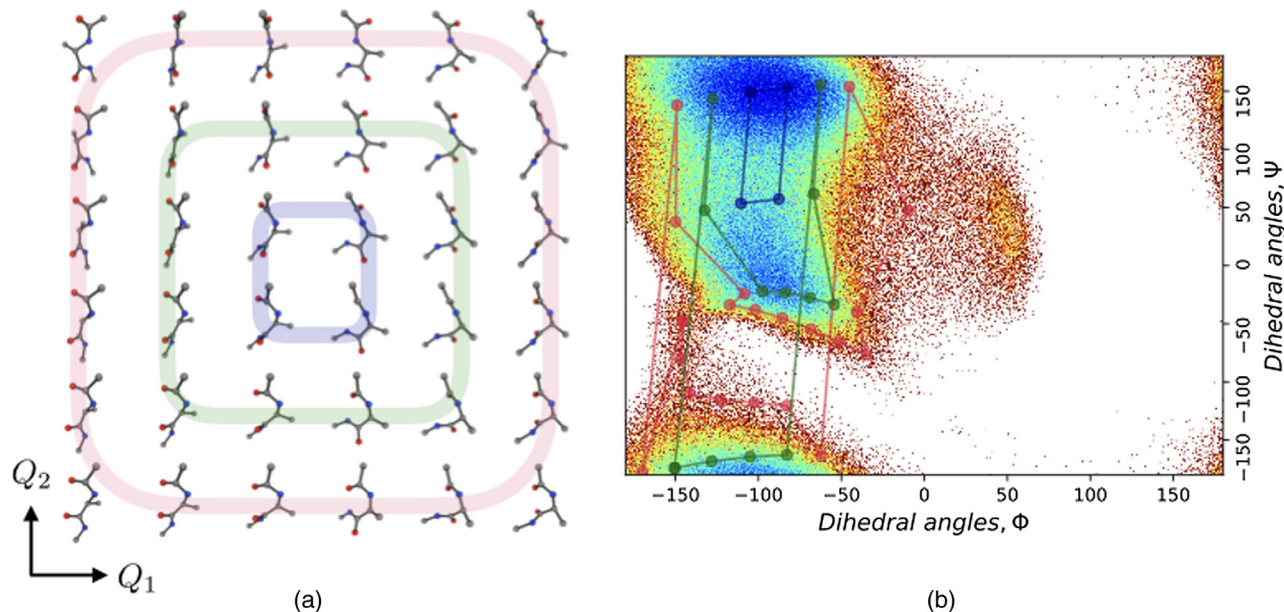


FIG. 5. (a) Latent-space interpolation in the plane of the two slowest collective coordinates ( $Q_1$ ,  $Q_2$ ) generates various molecule conformations. (b) Probability profile of the alanine dipeptide as a function of dihedral angles from generated samples of neural canonical transformation. The paths corresponding to the path with the same color in (a).

optimizer [67] with a minibatch size 200 and an initial learning rate  $10^{-3}$  for training. We reduce the learning rate by a factor of 10 if there is no improvement of the loss function on the validation set for 10 training steps.

Figure 4(a) shows the learned frequencies of the alanine dipeptide dataset, which spans a wide scale and suggests the emergence of slow collective modes in the system. The frequency of the slowest modes is smaller than the fastest mode by more than 3 orders of magnitude. Moreover, there is a notable gap in the frequencies, which suggests a separation of the fast and slow modes in the system.

We identify the latent coordinates with the smallest frequencies as the slowest nonlinear collective variables. To connect these learned collective variables to empirically known torsion angles, we estimate the mutual information [68] between them and the two torsion angles  $\Phi$  and  $\Psi$  in Fig. 4(b). One sees that the first two latent coordinates show a significant correlation with  $\Psi$  and  $\Phi$ , respectively. The mutual information between latent coordinates with larger frequencies and these two torsion angles decreases rapidly. Therefore, we conclude that the symplectic network has successfully identified the relevant slow modes which capture the low energy physics. Remarkably, this is done without having any access to the time information in the MD trajectory. This discovery highlights the usefulness of imposing the symplectic symmetry in the flow to turn statistical information into the dynamical one. We note that despite showing large mutual information, the learned two slowest coordinates do not exactly reproduce these torsion angles. The reason being that the learning objective encourages one to identify independent collective variables whose marginal distributions are independent Gaussians,

while the marginal distributions of these torsion angles are clearly not. Instead, this objective that favors independence allows us to gain better control over the identified latent variables, and thus offer advantages for practical purpose as we show next.

Since the normalizing flow model is a bijective generative model, one can directly map latent variables to molecular configurations. Figure 5(a) shows the generated molecular conformation by tuning the two slowest modes in the range  $Q_k \in [-1/\omega_k, 1/\omega_k]$ . One clearly sees that the two slowest variables control the global geometry of the molecule. Figure 5(b) shows generated samples in the two-dimensional plane of torsion angles. The learned distribution is wider than the given dataset, which is a common feature of the density estimation using the normalizing flows [24]. Figure 5(b) also shows that smooth paths in the spaces of learned slow latent variables may correspond to nontrivial paths in the torsion angles plane. The neural canonical transformation has learned a compact embedding molecular conformations in the learned latent space.

Having a compact latent representation of the molecular configurations allows one to design a smooth path between stable conformations by interpolating a few latent variables. Figure 6 shows a path connecting two molecular conformations with the spherical linear interpolation of the two slowest variables [69]. The interpolation gives a path along the geodesic curve of the Gaussian distributed latent variables and thus avoids unlikely molecular conformations. As the figure shows, the interpolation yields a curved path in the torsion angle plane which avoids the low-density region that may occur with a naive interpolation in terms of atom coordinates or torsion angles.



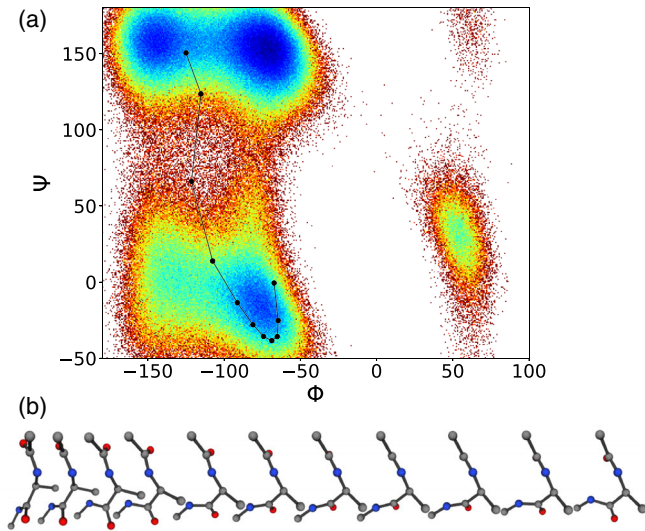


FIG. 6. (a) A path from the molecular configuration at  $(-125^\circ, 150^\circ)$  to the configuration at  $(-75^\circ, 0^\circ)$  obtained by spherical linear interpolation of the two slowest latent variables. The background is the probability profile of the alanine dipeptide dataset on the plane of the dihedral angles. (b) The molecular conformation along the interpolation path.

Moreover, thanks to the tractability of the normalizing flow model one has an exact likelihood on the path for the generated molecular unlike in the case of generating molecular conformation VAE or GAN. This quantitative access to the likelihood allows unbiased sampling of the configuration space with rejection sampling [70,71] or latent-space Monte Carlo updates [48,53].

#### D. MNIST handwritten digits

Finally, we apply the neural canonical transformation to the problems of machine learning. We consider the MNIST handwritten digits dataset, which contains 50 000 grayscale images of  $28 \times 28$  pixels. These images are divided into ten-digit classes. Treating the pixel values as coordinate variables, we can view the digit classes as stable conformations of a physical system [72]. Similar to the dipeptide studied in the Sec. IV C, one conjectures that the transition between conformations are slow, while the variations within the digits classes are the fast degrees of freedom. We assume each pixel has unity mass and augment the dataset with momenta. Then, we perform the density estimation in the phase space to train a neural canonical transformation.

Figure 7(a) shows the dispersion of the MNIST dataset, which contains a small portion of slow frequencies over all the variables. To show that these slow modes contain the relevant information of the digits classes, we pass only these slow modes to a multilayer perceptron classifier and perform supervised training. The classifier contains a single hidden layer of neurons with rectified linear units. The learned neural canonical transformation has its parameter

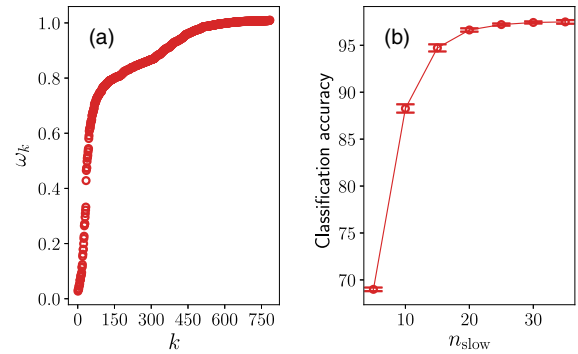


FIG. 7. (a) The learned frequencies of the MNIST dataset. (b) Classification accuracy on the test dataset based on  $n_{\text{slow}}$  slowest modes.

fixed and works as a feature extractor. By varying the number of kept slow modes from 5 up to 35 out of the total 784 dimensions, one sees that the classification accuracy on the test dataset quickly increases a plateau around 97% as shown in Fig. 7(b). Reaching high classification accuracy with only a few of the slow collective variables shows that they indeed capture digit class information.

We perform an additional experiment to directly show that the learned slow modes indeed capture the salient features of the MNIST images, i.e., the digit classes. As shown in the top panel of Fig. 8, we take a pair of

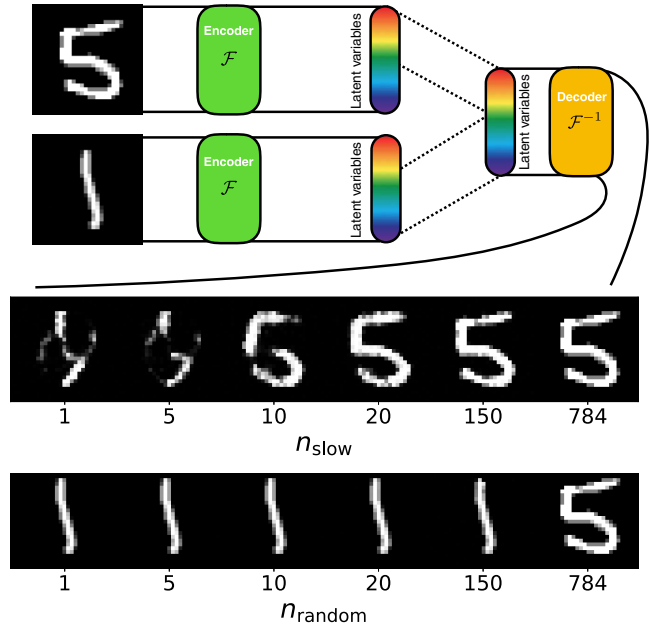


FIG. 8. The coordinate transformation Eq. (2) maps the MNIST images to a latent representation with independent frequency modes. We concatenate slow latent variables of an image together with the fast modes of another image and then map the combined latent vector back to the image space with the inverse coordinates transformation. The bottom panel shows the same experiment by concatenating latent vectors at random.



images from MNIST and map both to latent space. Then, we take  $n_{\text{slow}}$  slowest latent vector from one image and take the remaining ones from another image and concatenate them together to form a latent vector. After mapping the concatenated latent vector back to the image space, we see that even using 20 slowest modes one can already change the digit class. In comparison, if we perform the same experiment with randomly selected  $n_{\text{random}}$  modes without considering the frequency order, one sees that one needs to use a much larger number of latent variables to make the transition of digit classes.

In Appendix B we perform conceptual compression [45,73] using the slow modes learned by the neural canonical transformation.

## V. DISCUSSIONS

Neural canonical transformation extends a long-standing theoretical tool with symplectic normalizing flows. It provides a systematic way to simplify Hamiltonian dynamics and extract independent nonlinear slow modes of Hamiltonian systems and natural dataset.

Besides data analysis, the neural canonical transformation may open the door to address the sampling problem in MD simulations. As a bottom line, one is ready to employ the existing enhanced sampling approaches [74–76] with the learned slow modes as the collective variables. Since these collective variables are differentiable and exhibit slow dynamics, it fulfills the typical requirement of collective variables. Moreover, one can already sample feasible molecular configurations directly by exploiting the learned canonical transformation as a flow-based generative model. These samples would approximate the target distribution well if the generative model is trained well. One can further correct the sampling bias by using the generative model as proposals in the Markov chain Monte Carlo [70,71]. Last, one may also perform Monte Carlo sampling in the latent space and then map the latent vectors to the physical samples [48,53]. These latent-space Monte Carlo updates extend the enhanced sampling approach based on variable transformations [77] to an adaptive setting.

It is also instructive to put the present approach in the contexts of probabilistic generative models [43–47]. Conventional generative models are concerned with the statistical properties of data in the configuration space with coordinates only, while neural canonical transformations deal with phase-space density. Therefore, they allow access to dynamical information by exploiting the symplectic structure in the transformation. Since the phase-space density is factorized to the momenta and coordinate parts for separable Hamiltonians, the KL divergence can be written as the sum of two terms for the KL divergences of the momenta and coordinates marginal distributions, respectively. The phase-space KL divergence is lower bounded by the one in the configuration space [78]. In this sense, one can

view the momenta as auxiliary variables that regularize the training of the neural coordinate transformations.

A pressing issue with the latent variable generative models is how to select a handful of most relevant latent variables after training. There have been various attempts to design hierarchical generative models [45,79–82] to capture global information of data with a few latent variables. The neural canonical transformation differs by selecting the collective variables according to the learned frequencies in the latent space. Therefore, one can assign a dynamical interpretation to the statistically learned latent representation. On the other hand, currently, we use a generic real NVP model to perform the coordinate transformation. Additional symmetry constraints like the invariance under translation, rotation, and permutation among identical particles, are useful for future applications [30,83]. In this case, one may devise an equivariant transformation by leveraging a symmetry-preserving energy model and using it to drive a gradient flow [37,57].

We remark that reaching a perfect harmonic Hamiltonian Eq. (4) in the latent space is generally not possible because it requires the original system to be integrable. In fact, a perfectly trained neural canonical transformation would reveal the invariant torus of integrable systems as shown in Ref. [32]. So we do not expect the periodic dynamics of the hidden variables under the prior harmonic Hamiltonian to replace the dynamics of the physical variables. But we do expect that the slow modes extracted from the procedure would be meaningful and useful for downstream tasks. More generally, the Kolmogorov-Arnold-Moser theory [84] shows that the phase-space trajectory of nearly integrable systems would only be deformed from quasiperiodic motions. Thus, we expect the neural canonical transformation would work well for systems with coherent collective motion. Along this line, the present approach may also be useful to study synchronization phenomena [85], where a collective motion emerges out of complex dynamical systems. Finally, extending the present work to more general time-dependent canonical transformations may give an even more powerful tool to study many-particle dynamical systems.

## ACKNOWLEDGMENTS

We thank Austen Lamacraft, Masatoshi Imada, Yuan Wan, Pan Zhang, Yi-Zhuang You, Zi Cai, Lei-Han Tang, Yueshui Zhang, Dian Wu, Yantao Wu, and Yixiao Chen for discussions. The work is supported by the Ministry of Science and Technology of China under Grants No. 2016YFA0300603 and No. 2016YFA0302400, the National Natural Science Foundation of China under Grant No. 11774398, and the Strategic Priority Research Program of Chinese Academy of Sciences Grant No. XDB28000000, and the Computational Chemical Center: Chemistry in Solution and at Interfaces funded by the DOE under Award No. DE-SC0019394.

## APPENDIX A: SYMPLECTIC FLOWS

We list several other forms of neural symplectic transformation and discuss their relation to known constructions in the literature.

### 1. Linear symplectic transformation

The simplest canonical transformation is a linear transformation to the input variables. We parametrize the linear symplectic transformation using the exponential map of its Lie algebra [86],

$$z = xe^Y \quad \text{with} \quad Y = \begin{pmatrix} A & B \\ C & -A^T \end{pmatrix}, \quad (\text{A1})$$

where  $B$ ,  $C$  are real symmetric matrices and  $A$  is an arbitrary real matrix. One can implement Eq. (A1) via efficient vector-matrix exponential multiplication. Since the symplectic group is connected [86], the exponential map covers all linear symplectic transformations. Moreover, one can obtain the reverse of the transformation by acting  $e^{-Y}$  instead. Accurate and efficient differentiation through the matrix exponential is discussed in Ref. [87].

In the special case of  $B = C = 0$  and  $A$  is a skew-symmetric matrix, i.e.,  $A = -A^T$ , the linear symplectic transformation reduces to the orthogonal transformation of both momenta and coordinates, which corresponds to the normal-mode transformation.

### 2. Continuous symplectic flow

In general, one can parametrize the canonical transformation using a scalar generating function  $G(\lambda)$ . Integrating the ordinary differential equation (ODE)

$$\dot{\lambda} = \nabla_{\lambda} G(\lambda) J \quad (\text{A2})$$

from time 0 to  $\tau$ , one can transform the original variables from  $\lambda(t=0) = \mathbf{x}$  to  $\lambda(t=\tau) = \mathbf{z}$ . As a consequence, the Hamiltonian evolution is a special form of symplectic flow in the phase space with the generating function being the Hamiltonian [1].

Equation (A2) corresponds to the infinitesimal canonical transformation [1], which covers a broad family of symplectic transformations discussed so far. For example, if the generating function is a linear function of the momenta, we will arrive at the neural point transformation equations (2) and (3) introduced in the main texts. While if the generating function is a quadratic function of  $\lambda$  we obtain the linear symplectic transformation equations (A1).

The continuous symplectic flow falls into the framework of Monge-Ampère flow in the optimal transport theory [57], where the transportation is induced by a gradient flow under a scalar potential function. The symplectic structure in Eq. (A2) simplifies the computation due to the volume-preserving property. In practice, the continuous

transformation equation (A2) can be implemented via the neural ODE [56]. Since the symplectic symmetry is crucial for the canonical transformation, it is crucial to employ symplectic integrators [3] in the neural ODE implementation. In particular, if one employs a symplectic leapfrog discretization of Eq. (A2) for a separable generating function, one will arrive at the transformations discussed in [32,51].

## APPENDIX B: CONCEPTUAL COMPRESSION OF THE MNIST DATASET

The extraction of the salient features as the slow modes is useful for compressions. For example, conceptual compression is a lossy compression scheme that aims at capturing the global information of the input data [45,73]. The conventional approaches make use of the VAEs or the neural networks with a hierarchical structure. However, we perform the compression based on the learned frequencies since the symplectic network naturally separates fast and slow degrees of freedom of the dataset. The top of Fig. 9 shows the setup of conceptual compression with learned neural canonical transformation. First, we use the learned nonlinear coordinate transformation equation (2) to map the data to the latent space. Then, we pass only a few slow modes to a decoder network. The decoder restores the image from the latent space by running the inverse transformation as the

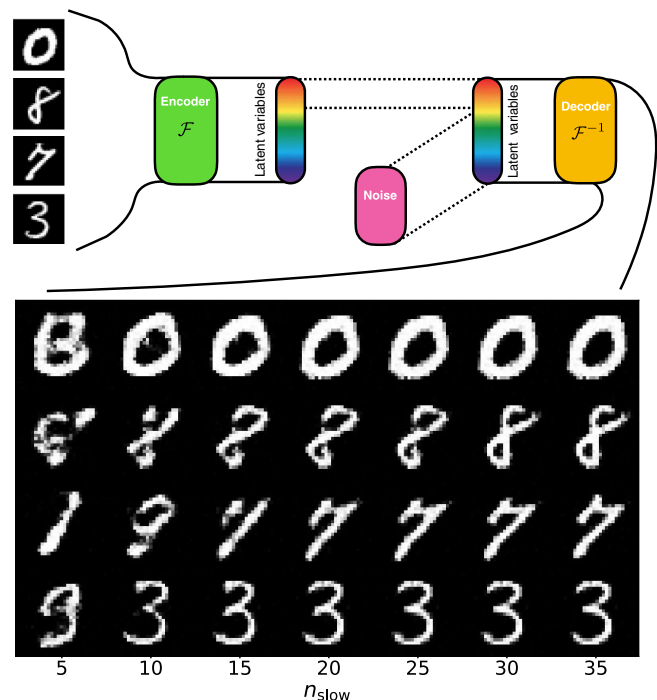


FIG. 9. Conceptual compression using the learned collective variables. One performs the coordinate transformation Eq. (2) to the input data, and restores the data based on  $n_{\text{slow}}$  slowest modes. The remaining fast modes are thrown away and resampled from the prior distribution.

encoder network. To make up the missing information, we simply sample the high-frequency modes from the prior distribution and feed them into the decoder. The bottom of Fig. 9 shows the results of the conceptual compression, that from left to right we keep 5, 10, 15, 20, 25, 30, 35 of the slowest collective variables. The conceptual compression experiments show that the symplectic transformation captures the global information in the slow modes in the latent space since one can restore the image with only a small number of the slowest variables.

- 
- [1] V. I. Arnold, *Mathematical Methods of Classical Mechanics* (Springer, New York, NY, 1989).
- [2] J. Liouville, *Note sur la Théorie de la Variation des Constantes Arbitraires*, *J. Math. Pures Appl.* **3**, 342 (1838).
- [3] K. Feng and M. Qin, *Symplectic Geometric Algorithms for Hamiltonian Systems* (Springer, New York, NY, 2011).
- [4] D. E. Shaw, P. Maragakis, K. Lindorff-Larsen, S. Piana, R. O. Dror, M. P. Eastwood, J. A. Bank, J. M. Jumper, J. K. Salmon, Y. Shan, and W. Wriggers, *Atomic-Level Characterization of the Structural Dynamics of Proteins*, *Science* **330**, 341 (2010).
- [5] F. Noé, *Machine Learning for Molecular Dynamics on Long Timescales*, [arXiv:1812.07669](https://arxiv.org/abs/1812.07669).
- [6] Y. Wang, J. M. L. Ribeiro, and P. Tiwary, *Machine Learning Approaches for Analyzing and Enhancing Molecular Dynamics Simulations*, [arXiv:1909.11748](https://arxiv.org/abs/1909.11748).
- [7] L. Molgedey and H. G. Schuster, *Separation of a Mixture of Independent Signals Using Time Delayed Correlations*, *Phys. Rev. Lett.* **72**, 3634 (1994).
- [8] A. Belouchrani, K. Abed-Meraim, J.-F. Cardoso, and E. Moulines, *A Blind Source Separation Technique Using Second-Order Statistics*, *IEEE Trans. Signal Process.* **45**, 434 (1997).
- [9] A. Ziehe and K.-R. Müller, in *Proceedings of ICANN 98*, edited by L. Niklasson, M. Bodén, and T. Ziemke (Springer London, London, 1998), pp. 675–680.
- [10] G. Pérez-Hernández, F. Paul, T. Giorgino, G. De Fabritiis, and F. Noé, *Identification of Slow Molecular Order Parameters for Markov Model Construction*, *J. Chem. Phys.* **139**, 015102 (2013).
- [11] C. R. Schwantes and V. S. Pande, *Improvements in Markov State Model Construction Reveal Many Non-Native Interactions in the Folding of NTL9*, *J. Chem. Theory Comput.* **9**, 2000 (2013).
- [12] P. J. Schmid, *Dynamic Mode Decomposition of Numerical and Experimental Data*, *J. Fluid Mech.* **656**, 5 (2010).
- [13] S. Klus, F. Nüske, P. Koltai, H. Wu, I. Kevrekidis, C. Schütte, and F. Noé, *Data-Driven Model Reduction and Transfer Operator Approximation*, *J. Nonlinear Sci.* **28**, 985 (2018).
- [14] B. Schölkopf, A. Smola, and K.-R. Müller, *Nonlinear Component Analysis as a Kernel Eigenvalue Problem*, *Neural Comput.* **10**, 1299 (1998).
- [15] S. Harmeling, A. Ziehe, M. Kawanabe, and K.-R. Müller, *Kernel-Based Nonlinear Blind Source Separation*, *Neural Comput.* **15**, 1089 (2003).
- [16] A. Mardt, L. Pasquali, H. Wu, and F. Noé, *VAMPnets for Deep Learning of Molecular Kinetics*, *Nat. Commun.* **9**, 5 (2018).
- [17] C. Wehmeyer and F. Noé, *Time-Lagged Autoencoders: Deep Learning of Slow Collective Variables for Molecular Kinetics*, *J. Chem. Phys.* **148**, 241703 (2018).
- [18] C. X. Hernández, H. K. Wayment-Steele, M. M. Sultan, B. E. Husic, and V. S. Pande, *Variational Encoding of Complex Dynamics*, *Phys. Rev. E* **97**, 062412 (2018).
- [19] M. M. Sultan and V. S. Pande, *Decision Functions from Supervised Machine Learning Algorithms as Collective Variables for Accelerating Molecular Simulations*, [arXiv:1802.10510](https://arxiv.org/abs/1802.10510).
- [20] B. Lusch, J. N. Kutz, and S. L. Brunton, *Deep Learning for Universal Linear Embeddings of Nonlinear Dynamics*, *Nat. Commun.* **9**, 4950 (2018).
- [21] L. Wiskott and T. J. Sejnowski, *Slow Feature Analysis: Unsupervised Learning of Invariances*, *Neural Comput.* **14**, 715 (2002).
- [22] D. Pfau, S. Petersen, A. Agarwal, D. G. T. Barrett, and K. L. Stachenfeld, *Spectral Inference Networks: Unifying Deep and Spectral Learning*, [arXiv:1806.02215](https://arxiv.org/abs/1806.02215).
- [23] I. Kobyzev, S. Prince, and M. A. Brubaker, *Normalizing Flows: Introduction and Ideas*, [arXiv:1908.09257](https://arxiv.org/abs/1908.09257).
- [24] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, *Normalizing Flows for Probabilistic Modeling and Inference*, [arXiv:1912.02762](https://arxiv.org/abs/1912.02762).
- [25] M. Mattheakis, P. Protopapas, D. Sondak, M. Di Giovanni, and E. Kaxiras, *Physical Symmetries Embedded in Neural Networks*, [arXiv:1904.08991v1](https://arxiv.org/abs/1904.08991v1).
- [26] S. Greydanus, M. Dzamba, and J. Yosinski, *Hamiltonian Neural Networks*, [arXiv:1906.01563](https://arxiv.org/abs/1906.01563).
- [27] A. Sanchez-Gonzalez, V. Bapst, K. Cranmer, and P. Battaglia, *Hamiltonian Graph Networks with ODE Integrators*, [arXiv:1909.12790](https://arxiv.org/abs/1909.12790).
- [28] Y. D. Zhong, B. Dey, and A. Chakraborty, *Symplectic ODE-Net: Learning Hamiltonian Dynamics with Control*, [arXiv:1909.12077](https://arxiv.org/abs/1909.12077).
- [29] Z. Chen, J. Zhang, M. Arjovsky, and L. Bottou, *Symplectic Recurrent Neural Networks*, [arXiv:1909.13334](https://arxiv.org/abs/1909.13334).
- [30] D. J. Rezende, S. Racanière, I. Higgins, and P. Toth, *Equivariant Hamiltonian Flows*, [arXiv:1909.13739](https://arxiv.org/abs/1909.13739).
- [31] P. Toth, D. J. Rezende, A. Jaegle, S. Racanière, A. Botev, and I. Higgins, *Hamiltonian Generative Networks*, [arXiv:1909.13789](https://arxiv.org/abs/1909.13789).
- [32] R. Bondesan and A. Lamacraft, *Learning Symmetries of Classical Integrable Systems*, [arXiv:1906.04645](https://arxiv.org/abs/1906.04645).
- [33] J. Behler and M. Parrinello, *Generalized Neural-Network Representation of High-Dimensional Potential-Energy Surfaces*, *Phys. Rev. Lett.* **98**, 146401 (2007).
- [34] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko, *Quantum-Chemical Insights from Deep Tensor Neural Networks*, *Nat. Commun.* **8**, 13890 (2017).
- [35] J. Han, L. Zhang, R. Car, and Weinan. E, *Deep Potential: A General Representation of a Many-Body Potential Energy Surface*, *Commun. Comput. Phys.* **23**, 629 (2018).



- [36] L. Zhang, J. Han, H. Wang, R. Car, and Weinan. E, *Deep Potential Molecular Dynamics: A Scalable Model with the Accuracy of Quantum Mechanics*, *Phys. Rev. Lett.* **120**, 143001 (2018).
- [37] L. Zhang, J. Han, H. Wang, W. A. Saidi, R. Car, and Weinan. E, *End-to-End Symmetry Preserving Inter-Atomic Potential Energy Model for Finite and Extended Systems*, in *Advances of the Neural Information Processing Systems (NIPS)*, Montreal, Canada (Curran Associates Inc., Red Hook, NY, 2018).
- [38] K. Gregor, F. Besse, D. J. Rezende, I. Danihelka, and D. Wierstra, *Towards Conceptual Compression*, [arXiv:1604.08772](https://arxiv.org/abs/1604.08772).
- [39] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, *Machine Learning of Accurate Energy-Conserving Molecular Force Fields*, *Sci. Adv.* **3**, e1603015 (2017).
- [40] I. Goodfellow, *NIPS 2016 Tutorial: Generative Adversarial Networks*, [arXiv:1701.00160](https://arxiv.org/abs/1701.00160).
- [41] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, *Generative Adversarial Networks*, [arXiv:1406.2661](https://arxiv.org/abs/1406.2661).
- [42] D. P. Kingma and M. Welling, *Auto-Encoding Variational Bayes*, [arXiv:1312.6114](https://arxiv.org/abs/1312.6114).
- [43] D. J. Rezende and S. Mohamed, *Variational Inference with Normalizing Flows*, [arXiv:1505.05770](https://arxiv.org/abs/1505.05770).
- [44] L. Dinh, D. Krueger, and Y. Bengio, *NICE: Non-Linear Independent Components Estimation*, [arXiv:1410.8516](https://arxiv.org/abs/1410.8516).
- [45] L. Dinh, J. Sohl-Dickstein, and S. Bengio, *Density Estimation Using Real NVP*, [arXiv:1605.08803](https://arxiv.org/abs/1605.08803).
- [46] D. P. Kingma and P. Dhariwal, *Glow: Generative Flow with Invertible  $1 \times 1$  Convolutions*, [arXiv:1807.03039](https://arxiv.org/abs/1807.03039).
- [47] W. Grathwohl, R. T. Q. Chen, J. Bettencourt, I. Sutskever, and D. Duvenaud, *FFJORD: Free-Form Continuous Dynamics for Scalable Reversible Generative Models*, [arXiv:1810.01367](https://arxiv.org/abs/1810.01367).
- [48] S.-H. Li and L. Wang, *Neural Network Renormalization Group*, *Phys. Rev. Lett.* **121**, 260601 (2018).
- [49] H.-Y. Hu, S.-H. Li, L. Wang, and Y.-Z. You, *Machine Learning Holographic Mapping by Neural Network Renormalization Group*, [arXiv:1903.00804](https://arxiv.org/abs/1903.00804).
- [50] J. Song, S. Zhao, and S. Ermon, *A-NICE-MC: Adversarial Training for MCMC*, [arXiv:1706.07561](https://arxiv.org/abs/1706.07561).
- [51] D. Levy, M. D. Hoffman, and J. Sohl-Dickstein, *Generalizing Hamiltonian Monte Carlo with Neural Networks*, [arXiv:1711.09268](https://arxiv.org/abs/1711.09268).
- [52] M. S. Albergo, G. Kanwar, and P. E. Shanahan, *Flow-Based Generative Models for Markov Chain Monte Carlo in Lattice Field Theory*, *Phys. Rev. D* **100**, 034515 (2019).
- [53] F. Noé, S. Olsson, J. Köhler, and H. Wu, *Boltzmann Generators—Sampling Equilibrium States of Many-Body Systems with Deep Learning*, *Science* **365** (2019).
- [54] G. S. Hartnett and M. Mohseni, *Self-Supervised Learning of Generative Spin-Glasses with Normalizing Flows*, [arXiv:2001.00585](https://arxiv.org/abs/2001.00585).
- [55] Weinan. E, *A Proposal on Machine Learning via Dynamical Systems*, *Commun. Math. Stat.* **5**, 1 (2017).
- [56] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, *Neural Ordinary Differential Equations*, [arXiv:1806.07366](https://arxiv.org/abs/1806.07366).
- [57] L. Zhang, Weinan. E, and Lei Wang, *Monge-Ampère Flow for Generative Modeling*, [arXiv:1809.10188](https://arxiv.org/abs/1809.10188).
- [58] G. J. Sussman and J. Wisdom, *Structure and Interpretation of Classical Mechanics*, 2nd ed. (The MIT Press, Cambridge, MA, 2014).
- [59] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind, *Automatic Differentiation in Machine Learning: A Survey*, *J. Machine Learning* **18**, 5595 (2015).
- [60] A. van den Oord *et al.*, *Parallel WaveNet: Fast High-Fidelity Speech Synthesis*, [arXiv:1711.10433](https://arxiv.org/abs/1711.10433).
- [61] S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih, *Monte Carlo Gradient Estimation in Machine Learning*, [arXiv:1906.10652](https://arxiv.org/abs/1906.10652).
- [62] D. Wu, L. Wang, and P. Zhang, *Solving Statistical Mechanics Using Variational Autoregressive Networks*, *Phys. Rev. Lett.* **122**, 080602 (2019).
- [63] See <https://github.com/li012589/neuralCT> for code implementation in PyTorch.
- [64] K. Pearson, *LIII. On Lines and Planes of Closest Fit to Systems of Points in Space*, *Philos. Mag.* **2**, 559 (1901).
- [65] F. Nüske, H. Wu, J. H. Prinz, C. Wehmeyer, C. Clementi, and F. Noé, *Markov State Models from Short Non-Equilibrium Simulations—Analysis and Correction of Estimation Bias*, *J. Chem. Phys.* **146**, 094104 (2017).
- [66] <https://markovmodel.github.io/mdshare/ALA2/#alanine-dipeptide>.
- [67] D. Kingma and J. Ba, *Adam: A Method for Stochastic Optimization*, in *Proceedings of the International Conference on Learning Representations (ICLR)* (2015).
- [68] A. Kraskov, H. Stögbauer, and P. Grassberger, *Estimating Mutual Information*, *Phys. Rev. E* **69**, 066138 (2004).
- [69] T. White, *Sampling Generative Networks*, [arXiv:1609.04468](https://arxiv.org/abs/1609.04468).
- [70] L. Huang and L. Wang, *Accelerated Monte Carlo Simulations with Restricted Boltzmann Machines*, *Phys. Rev. B* **95**, 035105 (2017).
- [71] J. Liu, Y. Qi, Z. Y. Meng, and L. Fu, *Self-Learning Monte Carlo Method*, *Phys. Rev. B* **95**, 041101(R) (2017).
- [72] See Ref. [57] for the preprocessing steps to map the MNIST dataset to continuous variables.
- [73] K. Gregor, F. Besse, D. J. Rezende, I. Danihelka, and D. Wierstra, *Towards Conceptual Compression*, [arXiv:1604.08772](https://arxiv.org/abs/1604.08772).
- [74] A. Laio and M. Parrinello, *Escaping Free-Energy Minima*, *Proc. Natl. Acad. Sci. U.S.A.* **99**, 12562 (2002).
- [75] O. Valsson and M. Parrinello, *Variational Approach to Enhanced Sampling and Free Energy Calculations*, *Phys. Rev. Lett.* **113**, 090601 (2014).
- [76] L. Zhang, H. Wang, and Weinan. E, *Reinforced Dynamics for Enhanced Sampling in Large Atomic and Molecular Systems*, *J. Chem. Phys.* **148**, 124113 (2018).
- [77] Z. Zhu, M. E. Tuckerman, S. O. Samuelson, and G. J. Martyna, *Using Novel Variable Transformations to Enhance Conformational Sampling in Molecular Dynamics*, *Phys. Rev. Lett.* **88**, 100201 (2002).
- [78] T. Salimans, D. P. Kingma, and M. Welling, *Markov Chain Monte Carlo and Variational Inference: Bridging the Gap*, [arXiv:1410.6460](https://arxiv.org/abs/1410.6460).

- [79] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel, *InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets*, [arXiv:1606.03657](https://arxiv.org/abs/1606.03657).
- [80] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, *beta-vae: Learning Basic Visual Concepts with a Constrained Variational Framework*, in *International Conference on Learning Representations, Toulon, France (2017)*, Vol. 3, <https://openreview.net/>.
- [81] T. Karras, S. Laine, and T. Aila, *A Style-Based Generator Architecture for Generative Adversarial Networks*, [arXiv:1812.04948](https://arxiv.org/abs/1812.04948).
- [82] H.P. Das, P. Abbeel, and C.J. Spanos, *Dimensionality Reduction Flows*, [arXiv:1908.01686](https://arxiv.org/abs/1908.01686).
- [83] J. Köhler, L. Klein, and F. Noé, *Equivariant Flows: Sampling Configurations for Multi-Body Systems with Symmetric Energies*, [arXiv:1910.00753](https://arxiv.org/abs/1910.00753).
- [84] H. S. Dumas, *The KAM Story: A Friendly Introduction to the Content, History, and Significance of Classical Kolmogorov-Arnold-Moser Theory* (World Scientific Publishing Company, Singapore, 2014).
- [85] J. A. Acebrón, L. L. Bonilla, C. J. P. Vicente, F. Ritort, and R. Spigler, *The Kuramoto Model: A Simple Paradigm for Synchronization Phenomena*, *Rev. Mod. Phys.* **77**, 137 (2005).
- [86] B. Hall, *Lie Groups, Lie Algebras and Representations*, 2nd ed. (Springer, New York, NY, 2015).
- [87] M. Lezcano-Casado, *Trivializations for Gradient-Based Optimization on Manifolds*, [arXiv:1909.09501](https://arxiv.org/abs/1909.09501).