

# ENERGETIC VARIATIONAL NEURAL NETWORK DISCRETIZATIONS OF GRADIENT FLOWS\*

ZIQUING HU<sup>†</sup>, CHUN LIU<sup>‡</sup>, YIWEI WANG<sup>§</sup>, AND ZHILIANG XU<sup>†</sup>

**Abstract.** We present a structure-preserving Eulerian algorithm for solving  $L^2$ -gradient flows and a structure-preserving Lagrangian algorithm for solving generalized diffusions. Both algorithms employ neural networks as tools for spatial discretization. Unlike most existing methods that construct numerical discretizations based on the strong or weak form of the underlying PDE, the proposed schemes are constructed based on the energy-dissipation law directly. This guarantees the monotonic decay of the system's free energy, which avoids unphysical states of solutions and is crucial for the long-term stability of numerical computations. To address challenges arising from nonlinear neural network discretization, we perform temporal discretizations on these variational systems before spatial discretizations. This approach is computationally memory-efficient when implementing neural network-based algorithms. The proposed neural network-based schemes are mesh-free, allowing us to solve gradient flows in high dimensions. Various numerical experiments are presented to demonstrate the accuracy and energy stability of the proposed numerical schemes.

**Key words.** structure-preserving, gradient flows, neural networks, energy stability, Lagrangian schemes

**MSC codes.** 35K35, 35K55, 49J40, 65M06, 65M12

**DOI.** 10.1137/22M1529427

**1. Introduction.** Evolution equations with variational structures, often termed as gradient flows, have a wide range of applications in physics, material science, biology, and machine learning [14, 60, 78]. These systems not only possess but also are determined by an energy-dissipation law, which consists of an energy of state variables that describes the energetic coupling and competition, and a dissipative mechanism that drives the system to an equilibrium.

More precisely, beginning with a prescribed energy-dissipation law

$$(1.1) \quad \frac{d}{dt} E^{\text{total}} = -\Delta \leq 0,$$

where  $E^{\text{total}}$  is the sum of the Helmholtz free energy  $\mathcal{F}$  and the kinetic energy  $\mathcal{K}$ , and  $\Delta$  is the rate of energy dissipation, one could derive the dynamics, the corresponding partial differential equation (PDE), of the system by combining the Least Action Principle (LAP) and the Maximum Dissipation Principle (MDP). The LAP states that the equation of motion of a Hamiltonian system can be derived from the variation of the action functional  $\mathcal{A} = \int_0^T (\mathcal{K} - \mathcal{F}) dt$  with respect to state variable  $\mathbf{x}$ . This gives

\*Submitted to the journal's Numerical Algorithms for Scientific Computing section October 18, 2022; accepted for publication (in revised form) April 12, 2024; published electronically August 5, 2024.

<https://doi.org/10.1137/22M1529427>

**Funding:** The work of the second and third authors was partially supported by the National Science Foundation grants DMS-1950868 and DMS-2153029. The work of the fourth author was partially supported by the National Science Foundation grant CDS&E-MSS 1854779.

<sup>†</sup>Department of Applied and Computational Mathematics and Statistics, University of Notre Dame, Notre Dame, IN 46556, USA (zhu4@nd.edu, zxu2@nd.edu).

<sup>‡</sup>Department of Applied Mathematics, Illinois Institute of Technology, Chicago, IL 60616, USA (cliu124@iit.edu).

<sup>§</sup>Corresponding author. Department of Mathematics, University of California, Riverside, Riverside, CA 92521, USA (yiweiw@ucr.edu).

rise to a unique procedure to derive the “conservative force” in the system. The MDP, on the other hand, derives the “dissipative force” by taking the variation of the dissipation potential  $\mathcal{D}$ , which equals  $\frac{1}{2}\Delta$  in the linear response regime (near equilibrium), with respect to  $\mathbf{x}_t$ . In turn, the force balance condition leads to the PDE of the system

$$(1.2) \quad \frac{\delta \mathcal{D}}{\delta \mathbf{x}_t} = \frac{\delta \mathcal{A}}{\delta \mathbf{x}}.$$

This procedure is known as the energetic variational approach (EnVarA) [24, 48, 74], developed based on the celebrated work of Onsager [57, 58] and Rayleigh [63] in nonequilibrium thermodynamics. During the past decades, the framework of EnVarA has shown to be a powerful tool for developing thermodynamically consistent models for various complex fluid systems, including two-phase flows [80, 81, 83], liquid crystals [48], ionic solutions [22], and reactive fluids [75, 76]. In a certain sense, the energy-dissipation law (1.1), which comes from the first and second law of thermodynamics [23], provides a more intrinsic description of the underlying physics than the derived PDE (1.2), in particular for systems that possess multiple structures and involve multiscale and multiphysics coupling.

From a numerical standpoint, the energy-dissipation law (1.1) also serves as a valuable guideline for developing structure-preserving numerical schemes for these variational systems, as many straightforward PDE-based discretizations may fail to preserve the continuous variational structures, as well as the physical constraints, such as the conservation of mass or momentum, the positivity of mass density, and the dissipation of the energy. As a recent development in this field, in [50], the authors proposed a numerical framework, called a *discrete energetic variational approach*, to address these challenges. Similar methods are also used in [72, 82]. The key idea of this approach is to construct numerical discretizations directly based on the continuous energy-dissipation laws without using the underlying PDE (see section 3.1 for details). The approach has advantages in preserving a discrete counterpart of the continuous energy-dissipation law, which is crucial for avoiding unphysical solutions and the long-term stability of numerical computations. Within the framework of the *discrete energetic variational approach*, Eulerian [55], Lagrangian [50, 51], and particle methods [73] have been developed for various problems.

However, tackling high-dimensional variational problems remains a challenge. Traditional numerical discretizations, such as finite difference, finite element, and finite volume methods, suffer from the well-known curse of dimensionality (CoD) [27]. Namely, the computational complexity of the numerical algorithm increases exponentially with the dimensionality of the problem [4, 27]. Although particle methods hold promise for addressing high-dimensional problems, standard particle methods often lack accuracy and are not suitable for problems involving derivatives.

During recent years, neural networks have demonstrated remarkable success across a wide spectrum of scientific disciplines [12, 31, 41, 43, 47]. Leveraging the potent expressive power of neural network [3, 10, 20], particularly deep neural network (DNN) architectures, there exists a growing interest in developing neural network-based algorithms for PDEs, especially for those in high dimensions [13, 21, 37, 38, 42, 62, 68, 70, 77]. Examples include physics-informed neural network (PINN) [62], deep Ritz method (DRM) [21], deep Galerkin method (DGM) [68], variational PINN [37], and weak adversarial network (WAN) [84] to name a few. A key component of these aforementioned approaches is to represent solutions of PDEs via neural networks. All of these approaches determine the optimal parameters of the neural network by minimizing a

loss function, which is often derived from either the strong or weak form of the PDE (see section 2.1 for more details). By employing neural network approximations, the approximate solution belongs to a space of nonlinear functions, which may lead to a robust estimation by sparser representation and cheaper computation [37].

The goal of this paper is to combine neural network-based spatial discretization with the framework of the discrete energetic variational approach [50], to develop efficient and robust numerical schemes, termed as energetic variational neural network (EVNN) methods. To clarify the idea, we consider the following two types of gradient flows modeled by EnVarA:

- **$L^2$ -gradient flow** that satisfies an energy-dissipation law

$$(1.3) \quad \frac{d}{dt} \mathcal{F}[\varphi] = - \int_{\Omega} \eta(\varphi) |\varphi_t|^2 d\mathbf{x}.$$

- **Generalized diffusion** that satisfies an energy-dissipation law

$$(1.4) \quad \frac{d}{dt} \mathcal{F}[\rho] = - \int_{\Omega} \eta(\rho) |\mathbf{u}|^2 d\mathbf{x},$$

where  $\rho$  satisfies  $\rho_t + \nabla \cdot (\rho \mathbf{u}) = 0$ , known as the mass conservation.

Derivation of the underlying PDEs for these types of systems is described in sections 2.1.1 and 2.1.2 of this paper.

Our primary aim is to develop structure-preserving Eulerian algorithms to solve  $L^2$ -gradient flows and structure-preserving Lagrangian algorithms to solve generalized diffusions based on their energy-dissipation law by utilizing neural networks as a new tool of spatial discretization. To overcome difficulties arising from neural network discretization, we develop a discretization approach that performs temporal discretizations before spatial discretizations. This approach leads to a computer-memory-efficient implementation of neural network-based algorithms. Since neural networks are advantageous due to their ability to serve as parametric approximations for unknown functions even in high dimensions, the proposed neural-network-based algorithms are capable of solving gradient flows in high-dimension, such as these appear in machine learning [19, 32, 69, 73].

The rest of this paper is organized as follows. Section 2 reviews the EnVarA and some existing neural network-based numerical approaches for solving PDEs. Section 3 of the paper is devoted to the development of the proposed EVNN schemes for  $L^2$ -gradient flows and generalized diffusions. Various numerical experiments are presented in section 4 to demonstrate the accuracy and energy stability of the proposed numerical methods. Conclusions are drawn in section 5.

## 2. Preliminary.

**2.1. Energetic variational approach.** In this subsection, we provide a detailed derivation of underlying PDEs for  $L^2$ -gradient flows and generalized diffusions by the general framework of EnVarA. We refer interested readers to [24, 74] for a more comprehensive review of the EnVarA. In both systems, the kinetic energy  $\mathcal{K} = 0$ .

**2.1.1.  $L^2$ -gradient flow.**  $L^2$ -gradient flows are those systems satisfying the energy-dissipation law:

$$(2.1) \quad \frac{d}{dt} \mathcal{F}[\varphi] = - \int_{\Omega} \eta |\varphi_t|^2 d\mathbf{x},$$

where  $\varphi$  is the state variable,  $\mathcal{F}[\varphi]$  is the Helmholtz free energy, and  $\eta > 0$  is the dissipation rate. One can also view  $\varphi$  as the generalized coordinates of the system

[14]. By treating  $\varphi$  as  $\mathbf{x}$ , the variational procedure (1.2) leads to the following  $L^2$ -gradient flow equation:

$$\frac{\delta(\frac{1}{2} \int \eta |\varphi_t|^2 d\mathbf{x})}{\delta \varphi_t} = -\frac{\delta \mathcal{F}}{\delta \varphi} \Rightarrow \eta \varphi_t = -\frac{\delta \mathcal{F}}{\delta \varphi}.$$

Many problems in soft matter physics, material science, and machine learning can be modeled as  $L^2$ -gradient flows. Examples include Allen–Cahn equation [15], Oseen–Frank and Landau–de Gennes models of liquid crystals [11, 45], phase field crystal models of quasicrystal [34, 44], and generalized Ohta–Kawasaki model of diblock and triblock copolymers [56]. For example, by taking

$$\mathcal{F}[\varphi] = \int_{\Omega} \frac{1}{2} |\nabla \varphi|^2 + \frac{1}{4\epsilon^2} (\varphi^2 - 1)^2 d\mathbf{x},$$

which is the classical Ginzburg–Landau free energy, and letting  $\eta(\varphi) = 1$ , one gets the Allen–Cahn equation

$$\varphi_t = \Delta \varphi - \frac{1}{\epsilon^2} (\varphi^2 - 1) \varphi.$$

**2.1.2. Generalized diffusion.** Generalized diffusion describes the space-time evolution of a conserved quantity  $\rho(\mathbf{x}, t)$ . Due to the physics law of mass conservation,  $\rho(\mathbf{x}, t)$  satisfies the kinematics

$$(2.2) \quad \partial_t \rho + \nabla \cdot (\rho \mathbf{u}) = 0,$$

where  $\mathbf{u}$  is an averaged velocity in the system.

To derive the generalized diffusion equation by the EnVarA, one should introduce a Lagrangian description of the system. Given a velocity field  $\mathbf{u}(\mathbf{x}, t)$ , one can define a flow map  $\mathbf{x}(\mathbf{X}, t)$  through

$$\begin{cases} \frac{d}{dt} \mathbf{x}(\mathbf{X}, t) = \mathbf{u}(\mathbf{x}(\mathbf{X}, t), t), \\ \mathbf{x}(\mathbf{X}, 0) = \mathbf{X}, \end{cases}$$

where  $\mathbf{X} \in \Omega_0$  is the Lagrangian coordinates and  $\mathbf{x} \in \Omega_t$  is the Eulerian coordinates. Here  $\Omega_0$  is the initial configuration of the system, and  $\Omega_t$  is the configuration of the system at time  $t$ . For a fixed  $\mathbf{X}$ ,  $\mathbf{x}(\mathbf{X}, t)$  describes the trajectory of a particle (or a material point) labeled by  $\mathbf{X}$ ; while for a fixed  $t$ ,  $\mathbf{x}(\mathbf{X}, t)$  is a diffeomorphism from  $\Omega_0$  to  $\Omega_t$  (See Figure 2.1 for the illustration). The existence of the flow map  $\mathbf{x}(\mathbf{X}, t)$  requires a certain regularity of  $\mathbf{u}(\mathbf{x}, t)$ , for instance, being Lipschitz in  $\mathbf{x}$ .

In a Lagrangian picture, the evolution of the density function  $\rho(\mathbf{x}, t)$  is determined by the evolution of the  $\mathbf{x}(\mathbf{X}, t)$  through the kinematics relation (2.2) (written in the

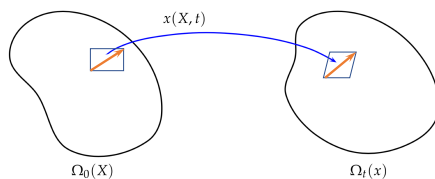


FIG. 2.1. Schematic illustration of a flow map  $\mathbf{x}(\mathbf{X}, t)$ .



Lagrangian frame of reference). More precisely, one can define the deformation tensor associated with the flow map  $\mathbf{x}(\mathbf{X}, t)$  by

$$(2.3) \quad \tilde{\mathbf{F}}(\mathbf{x}(\mathbf{X}, t), t) = \mathbf{F}(\mathbf{X}, t) = \nabla_{\mathbf{X}} \mathbf{x}(\mathbf{X}, t).$$

Without ambiguity, we do not distinguish  $\mathbf{F}$  and  $\tilde{\mathbf{F}}$  in the rest of this paper. Let  $\rho_0(\mathbf{X})$  be the initial density, then the mass conservation indicates that

$$(2.4) \quad \rho(\mathbf{x}(\mathbf{X}, t), t) = \rho_0(\mathbf{X}) / \det \mathbf{F}(\mathbf{X}, t) \quad \forall \mathbf{X} \in \Omega_0$$

in the Lagrangian frame of reference. This is equivalent to  $\rho_t + \nabla \cdot (\rho \mathbf{u}) = 0$  in the Eulerian frame.

Within (2.4), one can rewrite the energy-dissipation law (1.4) in terms of the flow map  $\mathbf{x}(\mathbf{X}, t)$  and its velocity  $\mathbf{x}_t(\mathbf{X}, t)$  in the reference domain  $\Omega_0$ . A typical form of the free energy  $\mathcal{F}[\rho]$  is given by

$$(2.5) \quad \mathcal{F}[\rho] = \int_{\Omega} \omega(\rho) + \rho V(\mathbf{x}) + \frac{1}{2} \left( \int_{\Omega} K(\mathbf{x}, \mathbf{y}) \rho(\mathbf{x}) \rho(\mathbf{y}) d\mathbf{y} \right) d\mathbf{x},$$

where  $V(\mathbf{x})$  is a potential field, and  $K(\mathbf{x}, \mathbf{y})$  is a symmetric nonlocal interaction kernel. In Lagrangian coordinates, the free energy (2.5) and the dissipation in (1.4) become

$$\mathcal{F}[\mathbf{x}(\mathbf{X}, t)] = \int_{\Omega_0} \omega \left( \frac{\rho_0}{\det \mathbf{F}} \right) \det \mathbf{F} + V(\mathbf{x}) \rho_0(\mathbf{X}) + \frac{\rho_0(\mathbf{X})}{2} \left( \int_{\Omega_0} K(\mathbf{x}, \mathbf{y}) \rho_0(\mathbf{X}) d\mathbf{X} \right) d\mathbf{X}$$

and

$$\mathcal{D}(\mathbf{x}, \mathbf{x}_t) = \frac{1}{2} \int_{\Omega_0} \eta(\rho_0 / \det \mathbf{F}) |\mathbf{x}_t|^2 \det \mathbf{F} d\mathbf{X},$$

By a direct computation, the force balance equation (1.2) can be written as (see [24] for a detailed computation)

$$(2.6) \quad \eta(\rho) \mathbf{u} = -\rho \nabla \mu, \quad \mu = \omega'(\rho) + V(\mathbf{x}) + K * \rho.$$

In Eulerian coordinates. Combining (2.6) with the kinematics equation (2.2), one obtains a generalized diffusion equation

$$(2.7) \quad \rho_t = \nabla \cdot (m(\rho) \nabla \mu), \quad m(\rho) = \rho^2 / \eta(\rho),$$

where  $m(\rho)$  is known as the mobility in physics. Many PDEs in a wide range of applications, including the porous medium equations (PME) [71], nonlinear Fokker–Planck equations [35], Cahn–Hilliard equations [50], Keller–Segel equations [36], and Poisson–Nernst–Planck (PNP) equations [22], can be obtained by choosing  $\mathcal{F}[\rho]$  and  $\eta(\rho)$  differently.

The velocity equation (2.6) can be viewed as the equation of the flow map  $\mathbf{x}(\mathbf{X}, t)$  in Lagrangian coordinates, i.e.,

$$(2.8) \quad \eta \left( \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}(\mathbf{X})} \right) \frac{d}{dt} \mathbf{x}(\mathbf{X}, t) = - \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}(\mathbf{X})} \mathbf{F}^{-T} \nabla_{\mathbf{X}} \mu \left( \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}(\mathbf{X})} \right),$$

where  $\mathbf{F}(\mathbf{X}) = \nabla_{\mathbf{X}} \mathbf{x}(\mathbf{X}, t)$ . The flow map equation (2.8) is a highly nonlinear equation of  $\mathbf{x}(\mathbf{X}, t)$  that involves both  $\mathbf{F}$  and  $\det \mathbf{F}$ . It is rather unclear how to introduce a suitable spatial discretization to (2.8) by looking at the equation. However, a

generalized diffusion can be viewed as an  $L^2$ -gradient flow of the flow map in the space of diffeomorphism. This perspective gives rise to a natural discretization of generalized diffusions [50].

*Remark 2.1.* In the case that  $\eta(\rho) = \rho$ , a generalized diffusion can be viewed as a Wasserstein gradient flow in the space of all probability densities having finite second moments  $\mathcal{P}_2(\Omega)$  [35]. Formally, the Wasserstein gradient flow can be defined as a continuous time limit ( $\tau \rightarrow 0$ ) from a semidiscrete Jordan-Kinderlehrer-Otto (JKO) scheme,

$$(2.9) \quad \rho^{n+1} = \arg \min_{\rho \in \mathcal{P}_2(\Omega)} \frac{1}{2\tau} W_2(\rho, \rho^n)^2 + \mathcal{F}[\rho], \quad n = 0, 1, 2, \dots,$$

where  $\mathcal{P}_2(\Omega) = \{\rho: \Omega \rightarrow [0, \infty) \mid \int_{\Omega} \rho \, d\mathbf{x} = 1, \int_{\Omega} |\mathbf{x}|^2 \rho(\mathbf{x}) \, d\mathbf{x} < \infty\}$  and  $W_2(\rho, \rho^n)$  is the Wasserstein distance between  $\rho$  and  $\rho^n$ . The Wasserstein gradient flow is an Eulerian description [6]. Other choices of dissipation can define other metrics in the space of probability measures [1, 46].

**2.2. Neural-network-based numerical schemes for PDEs.** In this subsection, we briefly review some existing neural network-based algorithms for solving PDEs. We refer interested readers to [18, 52] for detailed reviews.

Considering a PDE subject to a certain boundary condition

$$(2.10) \quad \mathcal{L}\varphi(\mathbf{x}) = f(\mathbf{x}), \quad \mathbf{x} \in \Omega \subset \mathbb{R}^d, \quad \mathcal{B}\varphi(\mathbf{x}) = g(\mathbf{x}), \quad \mathbf{x} \in \partial\Omega,$$

and assuming the solution can be approximated by a neural network, denoted by  $\varphi_{\text{NN}}(\mathbf{x}; \Theta)$ , where  $\Theta$  is the set of parameters in the neural network, the PINN [62] and similar methods [13, 68, 77] find the optimal parameter  $\Theta^*$  by minimizing a loss function defined as

$$(2.11) \quad L(\Theta) = \frac{1}{N_{in}} \sum_{i=1}^{N_{in}} (\mathcal{L}\varphi_{\text{NN}}(\mathbf{x}_i, \Theta) - f(\mathbf{x}_i))^2 + \frac{\lambda}{N_b} \sum_{j=1}^{N_b} (\mathcal{B}\varphi_{\text{NN}}(\mathbf{s}_j; \Theta) - g(\mathbf{s}_j))^2.$$

Here  $\{\mathbf{x}_i\}_{i=1}^{N_{in}}$  and  $\{\mathbf{s}_j\}_{j=1}^{N_b}$  are sets of samples in  $\Omega$  and  $\partial\Omega$ , respectively, which can be drawn uniformly or by following some prescribed distributions. The parameter  $\lambda$  is used to weigh the sum in the loss function. Minimizers of the loss function (2.11) can be obtained by using some optimization methods, such as AdaGrad and Adam. Of course, the objective function of this minimization problem, in general, is not convex even when the initial problem is. Obtaining the global minimum of (2.11) is highly nontrivial.

In contrast to the PINN, which is based on the strong form of a PDE, the DRM [21] is designed for solving PDEs using their variational formulations. If (2.10) is an Euler-Lagrangian equation of some energy functional

$$(2.12) \quad I[\varphi, \nabla\varphi] = \int_{\Omega} W(\varphi(\mathbf{x}), \nabla_{\mathbf{x}}\varphi(\mathbf{x})) \, d\mathbf{x} + \lambda \int_{\partial\Omega} |\mathcal{B}\varphi - g(\mathbf{x})|^2 \, dS,$$

then the loss function can be defined directly by

$$(2.13) \quad L(\theta) = \frac{1}{N_{in}} \sum_{i=1}^{N_{in}} W(\varphi_{\text{NN}}(\mathbf{x}_i), \nabla_{\mathbf{x}}\varphi_{\text{NN}}(\mathbf{x}_i)) + \frac{\lambda}{N_b} \sum_{j=1}^{N_b} (\mathcal{B}\varphi_{\text{NN}}(\mathbf{s}_j) - g(\mathbf{s}_j))^2.$$

The last term in (2.12) is a penalty term for the boundary condition. The original Dirichlet boundary condition can be recovered with  $\lambda \rightarrow \infty$ . Again, the samples

$\{\mathbf{x}_i\}_{i=1}^{N_{in}}$  and  $\{\mathbf{s}_j\}_{j=1}^{N_b}$  can be uniformly sampled in  $\Omega$  and  $\partial\Omega$  or sampled following some other prescribed distributions, respectively. Additionally, both the variational PINN [37] and the WAN [84] utilize the Galerkin formulation to solve PDEs. The variational PINN [37] stems from the Petrov–Galerkin method, represents the solution via a DNN, and keeps test functions belonging to linear function spaces. The WAN employs the primal and adversarial networks to parameterize the weak solution and test functions, respectively [84].

The neural network-based algorithms mentioned above focus on elliptic equations. For an evolution equation of the form

$$(2.14) \quad \begin{cases} \partial_t \varphi(\mathbf{x}, t) = F(t, \mathbf{x}, \varphi), & (\mathbf{x}, t) \in \Omega \times (0, \infty), \\ \varphi(\mathbf{x}, 0) = \varphi_0(\mathbf{x}), & \mathbf{x} \in \Omega \end{cases}$$

with a suitable boundary condition, the predominant approach is to treat the time variable as an additional dimension, and the PINN type techniques can be applied. However, these approaches are often expensive and may fail to capture the dynamics of these systems. Unlike spatial variables, there exists an inherent order in time, where the solution at time  $T$  is determined by the solutions in  $t < T$ . It is difficult to generate samples in time to maintain this causality. Very recently, new efforts have been made in using neural networks to solve evolution PDEs [8, 16]. The idea of these approaches is to use neural networks with time-dependent parameters to represent the solutions. For instance, Neural Galerkin method, proposed in [8], parameterizes the solution as  $\varphi_h(\mathbf{x}; \Theta(t))$  with  $\Theta(t)$  being the neural network parameters, and defines a loss function in terms of  $\Theta$  and  $\dot{\Theta}$  through a residual function, i.e.,

$$(2.15) \quad J(\Theta, \eta) = \frac{1}{2} \int |\nabla_{\Theta} \varphi_h \cdot \eta - F(\mathbf{x}, t, \varphi_h(\mathbf{x}; \Theta))|^2 d\nu_{\Theta}(\mathbf{x}),$$

where  $\nu_{\Theta}(\mathbf{x})$  is a suitable measure which might depend on  $\Theta$ . By taking variation of  $J(\Theta, \eta)$  with respect to  $\eta$ , the neural Galerkin method arrives at an ODE of  $\Theta(t)$ :

$$(2.16) \quad \mathbf{M}(\Theta) \dot{\Theta} = F(t, \Theta), \quad \Theta(0) = \Theta_0,$$

where  $\mathbf{M}(\Theta) = \int \nabla_{\Theta} \varphi_h \otimes \nabla_{\Theta} \varphi_h d\nu_{\Theta}(\mathbf{x})$ , and  $F(t, \Theta) = \int \nabla_{\Theta} \varphi_h F(\mathbf{x}, t, \varphi_h) d\nu_{\Theta}(\mathbf{x})$ . The ODE (2.16) can be solved by standard explicit or implicit time-marching schemes.

It's important to note that, with the exception of DRM, all existing neural-network-based methods are developed based on either the strong or weak forms of PDEs. While DRM utilizes the free energy functional, it is only suitable for solving static problems, i.e., finding the equilibrium of the system. Additionally, most of these existing approaches are Eulerian methods. For certain types of PDEs, like generalized diffusions, these methods may fail to preserve physical constraints, such as the positivity and conservation of mass of a probability function.

**3. Energetic variational neural network.** In this section, we present the structure-preserving EVNN discretization for solving both  $L^2$ -gradient flows and generalized diffusions. As mentioned earlier, the goal is to construct a neural-network discretization based on the energy-dissipation law, without working on the underlying PDE. One can view our method as a generalization of the DRM to evolution equations.

**3.1. EVNN scheme for  $L^2$ -gradient flows.** Before we discuss the neural network discretization, we first briefly review the discrete energetic variational approach proposed in [50]. Given a continuous energy-dissipation law (1.1), the discrete

energetic variational approach first constructs a finite-dimensional approximation to this law by introducing a spatial discretization of the state variable  $\varphi$ , denoted by  $\varphi_h(\mathbf{x}; \boldsymbol{\Theta}(t))$ , where  $\boldsymbol{\Theta}(t) \in \mathbb{R}^K$  is the parameter to be determined. By replacing  $\varphi$  by  $\varphi_h$ , one can obtain a semidiscrete energy-dissipation law (here we assume  $\mathcal{K} = 0$  in the continuous model for simplicity) in terms of  $\boldsymbol{\Theta}$ :

$$(3.1) \quad \frac{d}{dt} \mathcal{F}_h[\boldsymbol{\Theta}(t)] = -\Delta_h[\boldsymbol{\Theta}(t), \boldsymbol{\Theta}'(t)],$$

where  $\mathcal{F}_h[\boldsymbol{\Theta}(t)] = \mathcal{F}[\varphi_h(\mathbf{x}; \boldsymbol{\Theta})]$ , and  $\Delta_h[\boldsymbol{\Theta}(t), \boldsymbol{\Theta}'(t)] = \int \eta(\varphi_h) |\nabla_{\boldsymbol{\Theta}} \varphi_h \cdot \boldsymbol{\Theta}'(t)|^2 d\mathbf{x}$ . For example, in Galerkin methods, one can let  $\varphi_h(\mathbf{x}, t)$  be  $\varphi_h(\mathbf{x}, t) = \sum_{i=1}^K \gamma_i(t) \psi_i(\mathbf{x})$  with  $\{\psi_i\}_{i=1}^K$  being the set of basis functions. Then  $\boldsymbol{\Theta}(t)$  is a vector given by  $\boldsymbol{\Theta}(t) = (\gamma_1(t), \gamma_2(t), \dots, \gamma_K(t))^T \in \mathbb{R}^K$ .

By employing the energetic variational approach in the semidiscrete level (3.1), one can obtain an ODE system of  $\boldsymbol{\Theta}(t)$ . Particularly, in the linear response regime,  $\mathcal{D}_h(\boldsymbol{\Theta}(t), \boldsymbol{\Theta}'(t)) = \frac{1}{2} \int \eta(\varphi_h) |\nabla_{\boldsymbol{\Theta}} \varphi_h \cdot \boldsymbol{\Theta}'(t)|^2 d\mathbf{x}$  is a quadratic function of  $\boldsymbol{\Theta}'$ . The ODE system of  $\boldsymbol{\Theta}(t)$  can then be written as

$$(3.2) \quad \mathbf{D}(\boldsymbol{\Theta}) \boldsymbol{\Theta}'(t) = -\frac{\delta \mathcal{F}_h}{\delta \boldsymbol{\Theta}},$$

where

$$(3.3) \quad \frac{\delta \mathcal{F}_h}{\delta \boldsymbol{\Theta}} = \frac{\delta \mathcal{F}}{\delta \varphi} \nabla_{\boldsymbol{\Theta}} \varphi_h, \quad \mathbf{D}(\boldsymbol{\Theta}) = \int \eta(\varphi_h) (\nabla_{\boldsymbol{\Theta}} \varphi_h \otimes \nabla_{\boldsymbol{\Theta}} \varphi_h) d\mathbf{x}.$$

The ODE (3.2) is the same as the ODE (2.16) in the Neural Galerkin method [8] for  $L^2$  gradient flows, although the derivation is different.

Since (3.2) is a finite-dimensional gradient flow, one can then construct a minimizing movement scheme for  $\boldsymbol{\Theta}$  [25]: finding  $\boldsymbol{\Theta}^{n+1}$  such that

$$(3.4) \quad \boldsymbol{\Theta}^{n+1} = \arg \min_{\boldsymbol{\Theta} \in \mathcal{S}_{ad}^h} J_n(\boldsymbol{\Theta}), \quad J_n(\boldsymbol{\Theta}) = \frac{(\mathbf{D}_*^n(\boldsymbol{\Theta} - \boldsymbol{\Theta}^n)) \cdot (\boldsymbol{\Theta} - \boldsymbol{\Theta}^n)}{2\tau} + \mathcal{F}_h(\boldsymbol{\Theta}).$$

Here  $\mathcal{S}_{ad}^h$  is the admissible set of  $\boldsymbol{\Theta}$  inherited from the admissible set of  $\varphi$ , denoted by  $\mathcal{S}$ , and  $\mathbf{D}_*^n$  is a constant matrix. A typical choice is  $\mathbf{D}_*^n = \mathbf{D}(\boldsymbol{\Theta}^n)$ . An advantage of this scheme is that

$$(3.5) \quad \mathcal{F}_h(\boldsymbol{\Theta}^{n+1}) \leq J_h^n(\boldsymbol{\Theta}^{n+1}) \leq J_h^n(\boldsymbol{\Theta}^n) = \mathcal{F}_h(\boldsymbol{\Theta}^n),$$

if  $\mathbf{D}_*^n$  is positive definite, which guarantees the energy stability with respect to the discrete free energy  $\mathcal{F}_h(\boldsymbol{\Theta})$ . Moreover, by choosing a proper optimization method, we can assure that  $\varphi_h^{n+1}$  stays in the admissible set  $\mathcal{S}$ .

Although neural networks can be used to construct  $\varphi_h(\mathbf{x}; \boldsymbol{\Theta}(t))$ , it might be expensive to compute  $\nabla_{\boldsymbol{\Theta}} \varphi_h \otimes \nabla_{\boldsymbol{\Theta}} \varphi_h$  in  $\mathbf{D}(\boldsymbol{\Theta})$ . Moreover,  $\mathbf{D}(\boldsymbol{\Theta})$  is not a sparse matrix and requires a lot of computer memory to store when a deep neural network is used.

To overcome these difficulties, we propose an alternative approach by introducing temporal discretization before spatial discretization. Let  $\tau$  be the time step size. For the  $L^2$ -gradient flow (2.1), given  $\varphi^n$ , which represents the numerical solution at  $t^n = n\tau$ , one can obtain  $\varphi^{n+1}$  by solving the following optimization problem: finding  $\varphi^{n+1}$  in some admissible set  $\mathcal{S}$  such that

$$(3.6) \quad \varphi^{n+1} = \arg \min_{\varphi \in \mathcal{S}} J^n(\varphi), \quad J^n(\varphi) = \frac{1}{2\tau} \int \eta(\varphi^n) |\varphi - \varphi^n|^2 d\mathbf{x} + \mathcal{F}[\varphi].$$

Let  $\varphi_h(\mathbf{x}; \Theta)$  be a finite-dimensional approximation of  $\varphi$  with  $\Theta \in \mathbb{R}^K$  being the parameter of the spatial discretization (e.g., weights of linear combination in a Galerkin approximation) yet to be determined, then the minimizing movement scheme (3.6) can be written in terms of  $\Theta$ : finding  $\Theta^{n+1}$  such that

$$(3.7) \quad \Theta^{n+1} = \arg \min_{\Theta \in S_h} J_h^n(\Theta), \quad J_h^n(\Theta) = \frac{1}{2\tau} \int \eta^n |\varphi_h(\mathbf{x}; \Theta) - \varphi_h(\mathbf{x}; \Theta^n)|^2 d\mathbf{x} + \mathcal{F}_h[\Theta].$$

Here  $\Theta^n$  is the value of  $\Theta$  at time  $t^n$ ,  $\mathcal{F}_h[\Theta] = \mathcal{F}[\varphi_h(\mathbf{x}; \Theta)]$ , and  $\eta^n = \eta(\varphi_h(\mathbf{x}; \Theta^n))$ .

*Remark 3.1.* The connection between the minimizing movement scheme (3.7), derived by a temporal-then-spatial approach, and the minimizing movement scheme (3.4), derived by a spatial-then-temporal approach, can be shown with a direct calculation. Indeed, according to the first-order necessary condition for optimality, an optimal solution of the minimization problem (3.7)  $\Theta^{n+1}$  satisfies

$$(3.8) \quad \begin{aligned} \frac{\delta J_h^n(\Theta)}{\delta \Theta} \Big|_{\Theta^{n+1}} &= \frac{1}{\tau} \int \eta^n (\varphi_h(\mathbf{x}; \Theta^{n+1}) - \varphi_h(\mathbf{x}; \Theta^n)) \nabla_{\Theta} \varphi_h \Big|_{\Theta^{n+1}} d\mathbf{x} + \frac{\delta \mathcal{F}_h}{\delta \Theta} \Big|_{\Theta^{n+1}} \\ &= \int \eta^n \nabla_{\Theta} \varphi_h \Big|_{\Theta^*} \otimes \nabla_{\Theta} \varphi_h \Big|_{\Theta^{n+1}} d\mathbf{x} \frac{\Theta^{n+1} - \Theta^n}{\tau} + \frac{\delta \mathcal{F}_h}{\delta \Theta} \Big|_{\Theta^{n+1}} = 0 \end{aligned}$$

for some  $\Theta^*$ , where the second equality follows the mean value theorem. In the case of Galerkin methods, as  $\nabla_{\Theta} \varphi_h$  is independent on  $\Theta$ ,  $\Theta^{n+1}$  is a solution to an implicit Euler scheme for the ODE (3.2) in which  $\eta(\varphi_h)$  is treated explicitly.

By choosing a certain neural network to approximate  $\varphi$ , denoted by  $\varphi_{\text{NN}}(\mathbf{x}; \Theta)$  ( $\Theta$  is used to denote all parameters in the neural network), we can summarize the EVNN scheme for solving  $L^2$ -gradient flows in Algorithm 3.1.

It can be noticed that both (3.9) and (3.10) involve integration in the computational domain  $\Omega$ . This integration is often computed by using a grid-based numerical quadrature or Monte-Carlo/Quasi-Monte-Carlo algorithms [65]. It is worth mentioning that due to the nonconvex nature of the optimization problem and the error in estimating the integration, it might be difficult to find an optimal  $\Theta^{n+1}$  at each step. But since we start the optimization procedure with  $\Theta^n$ , we'll always be able to get a  $\Theta^{n+1}$  that lowers the discrete free energy at least on the training set. In practice, the optimization problem can be solved by either deterministic optimization algorithms, such as L-BFGS and gradient descent with Barzilai-Borwein step-size, or stochastic gradient descent algorithms, such as AdaGrad and Adam.

---

**Algorithm 3.1** Numerical algorithm for solving  $L^2$ -gradient flows.

---

For a given initial condition  $\varphi_0(\mathbf{x})$ , compute  $\Theta^0$  by solving

$$(3.9) \quad \Theta^0 = \arg \min_{\Theta} \int_{\Omega} |\varphi_{\text{NN}}(\mathbf{x}; \Theta) - \varphi_0(\mathbf{x})|^2 d\mathbf{x}.$$

At each step, update  $\Theta^{n+1}$  by solving the optimization problem

$$(3.10) \quad \Theta^{n+1} = \arg \min_{\Theta} \left( \frac{1}{2\tau} \int_{\Omega} \eta^n |\varphi_{\text{NN}}(\mathbf{x}; \Theta) - \varphi_{\text{NN}}(\mathbf{x}; \Theta^n)|^2 d\mathbf{x} + \mathcal{F}[\varphi_{\text{NN}}(\mathbf{x}; \Theta)] \right).$$

We have  $\varphi_{\text{NN}}(\mathbf{x}; \Theta^n)$  as a numerical solution at time  $t^n = n\tau$ .

---

*Remark 3.2.* It is straightforward to incorporate other variational high-order temporal discretizations to solve the  $L^2$ -gradient flows. For example, a second-order accurate BDF2 scheme can be reformulated as an optimization problem

$$(3.11) \quad \begin{aligned} \Theta^n = \arg \min_{\Theta} & \left( \frac{\eta}{\tau} \int_{\Omega} |\varphi_{\text{NN}}(\mathbf{x}; \Theta) - \varphi_{\text{NN}}(\mathbf{x}, \Theta^{n-1})|^2 d\mathbf{x} \right. \\ & \left. - \frac{\eta}{4\tau} \int_{\Omega} |\varphi_{\text{NN}}(\mathbf{x}; \Theta) - \varphi_{\text{NN}}(\mathbf{x}, \Theta^{n-2})|^2 d\mathbf{x} + \mathcal{F}[\varphi_{\text{NN}}(\mathbf{x}; \Theta)] \right). \end{aligned}$$

A modified Crank–Nicolson time-marching scheme can be reformulated as

$$(3.12) \quad \begin{aligned} \Theta^n = \arg \min_{\Theta} & \left( \frac{\eta}{\tau} \int_{\Omega} |\varphi_{\text{NN}}(\mathbf{x}; \Theta) - \varphi_{\text{NN}}(\mathbf{x}, \Theta^{n-1})|^2 d\mathbf{x} + \mathcal{F}[\varphi_{\text{NN}}(\mathbf{x}; \Theta)] \right. \\ & \left. + \langle \nabla_{\Theta} \mathcal{F}[\varphi_{\text{NN}}(\mathbf{x}; \Theta^{n-1})], \Theta - \Theta^{n-1} \rangle \right). \end{aligned}$$

Here we assume that  $\eta$  is a constant for simplicity.

**3.2. Lagrangian EVNN scheme for generalized diffusions.** In this subsection, we show how to formulate an EVNN scheme in the Lagrangian frame of reference for generalized diffusions.

As discussed previously, a generalized diffusion can be viewed as an  $L^2$ -gradient flow of the flow map  $\mathbf{x}(\mathbf{X}, t)$  in the space of diffeomorphisms. Hence, the EVNN scheme for a generalized diffusion can be formulated in terms of a minimizing movement scheme of the flow map given by

$$(3.13) \quad \Phi^{n+1} = \arg \min_{\Phi \in \text{Diff}} \frac{1}{2\tau} \int |\Phi(\mathbf{X}) - \Phi^n(\mathbf{X})|^2 \rho_0(\mathbf{X}) d\mathbf{X} + \mathcal{F}[\Phi_{\#} \rho_0],$$

where  $\Phi^n(\mathbf{X})$  is a numerical approximation of the flow map  $\mathbf{x}(\mathbf{X}, t)$  at  $t^n = n\tau$ ,  $\rho_0(\mathbf{X})$  is the initial density,  $\mathcal{F}[\rho]$  is the free energy for the generalized diffusion defined in (2.5), and

$$(\Phi_{\#} \rho_0)(\mathbf{x}) := \frac{\rho_0(\Phi^{-1}(\mathbf{x}))}{\det F(\Phi^{-1}(\mathbf{x}))}, \quad \text{Diff} = \{\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^d \mid \Phi \text{ is a diffeomorphism}\}.$$

One can parameterize  $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$  by a suitable neural network. The remaining procedure is nearly identical to that of the previous subsection. However, it is often difficult to solve this optimization problem directly, and one might need to build a large neural network to approximate  $\Phi^{n+1}$  when  $n$  is large.

To overcome this difficulty, we propose an alternative approach. Instead of seeking an optimal map  $\Phi^{n+1}$  between  $\rho^0$  and  $\rho^{n+1}$ , we seek an optimal map  $\Psi^{n+1}$  between  $\rho^n$  and  $\rho^{n+1}$ . More precisely, we assume that

$$\Phi^{n+1} = \Psi^{n+1} \circ \Psi^n \circ \Psi^{n-1} \dots \circ \Psi^1.$$

Given  $\rho^n$ , one can compute  $\Psi^{n+1}$  by solving the following optimization problem:

$$(3.14) \quad \Psi^{n+1} = \arg \min_{\Psi \in \text{Diff}} \frac{1}{2\tau} \int |\Psi(\mathbf{x}) - \mathbf{x}|^2 \rho^n(\mathbf{x}) d\mathbf{x} + \mathcal{F}[\Psi_{\#} \rho^n].$$

The corresponding  $\rho^{n+1}$  can then be computed through  $\rho^{n+1}(\mathbf{x}) = (\Psi_{\#}^{n+1} \rho^n)(\mathbf{x})$ . An advantage of this approach is that we only need a small size neural network to approximate  $\Psi^{n+1}$  at each time step when  $\tau$  is small.

*Remark 3.3.* The scheme (3.14) can be viewed as a Lagrangian realization of the JKO scheme (2.9) for the Wasserstein gradient flow, although it is developed based on the  $L^2$ -gradient flow structure in the space of diffeomorphism. According to the Benamou–Brenier formulation [5], the Wasserstein distance between two probability densities  $\rho_1$  and  $\rho_2$  can be computed by solving the optimization problem

$$(3.15) \quad W_2(\rho_1, \rho_2)^2 = \min_{(\rho, \mathbf{u}) \in \mathcal{S}} \int_0^1 \int \rho |\mathbf{u}|^2 d\mathbf{x} dt,$$

where the admissible set of  $(\rho, \mathbf{u})$  is given by

$$(3.16) \quad \mathcal{S} = \{(\rho, \mathbf{u}) \mid \rho_t + \nabla \cdot (\rho \mathbf{u}) = 0, \quad \rho(\mathbf{x}, 0) = \rho_1, \quad \rho(\mathbf{x}, 1) = \rho_2\}.$$

Hence, one can solve the JKO scheme by solving

$$(3.17) \quad \begin{aligned} (\mathbf{u}^*, \hat{\rho}^*) &= \arg \min_{(\mathbf{u}, \hat{\rho})} \frac{1}{2} \int_0^\tau \int_\Omega \hat{\rho} |\mathbf{u}|^2 d\mathbf{x} dt + \mathcal{F}[\hat{\rho}(\tau)] \\ \text{s.t.} \quad &\partial_t \hat{\rho} + \nabla \cdot (\hat{\rho} \mathbf{u}) = 0, \quad \hat{\rho}(0) = \rho^n. \end{aligned}$$

and letting  $\rho^{n+1} = \hat{\rho}^*(\tau)$ . In Lagrangian coordinates, (3.17) is equivalent to

$$(3.18) \quad \mathbf{x}^*(\mathbf{X}, t) = \arg \min_{\mathbf{x}(\mathbf{X}, t)} \frac{1}{2} \int_0^\tau \int_\Omega \rho_0(\mathbf{X}) |\mathbf{x}_t(\mathbf{X}, t)|^2 d\mathbf{X} dt + \mathcal{F}(\hat{\rho}(\mathbf{x}, \tau)),$$

where  $\rho_0 = \rho^n$  and  $\hat{\rho}(\mathbf{x}(\mathbf{X}, \tau), \tau) = \rho_0(\mathbf{X}) / \det \mathbf{F}(\mathbf{X}, \tau)$ . By taking variation of (3.18) with respect to  $\mathbf{x}(\mathbf{X}, t)$ , one can show that the optimal condition is  $\mathbf{x}_{tt}(\mathbf{X}, t) = 0$  for  $t \in (0, \tau)$ , which indicates that  $\mathbf{x}(\mathbf{X}, t) = t(\Psi(\mathbf{X}) - \mathbf{X})/\tau + \mathbf{X}$  if  $\mathbf{x}(\mathbf{X}, \tau) = \Psi(\mathbf{X})$ . Hence, if  $\Psi^*$  is the optimal solution of (3.14), then  $\mathbf{x}^*(\mathbf{X}, t) = t(\Psi^*(\mathbf{X}) - \mathbf{X})/\tau + \mathbf{X}$  is the optimal solution of (3.18).

*Remark 3.4.* If  $\eta(\rho) \neq \rho$ , we can formulate the optimization problem (3.14) as

$$(3.19) \quad \Psi^{n+1} = \arg \min_{\Psi \in \text{Diff}} \frac{1}{2\tau} \int |\Psi(\mathbf{x}) - \mathbf{x}|^2 \eta(\rho^n(\mathbf{x})) d\mathbf{x} + \mathcal{F}[\Psi_{\#} \rho^n].$$

by treating  $\eta$  explicitly. A subtle fact is that  $\det \mathbf{F}^n = 1$  since we always start with an identity map.

The numerical algorithm for solving generalized diffusions is summarized in Algorithm 3.2.

---

**Algorithm 3.2** Numerical algorithm for solving generalized diffusions.

---

- Given  $\{\mathbf{x}_i^n\}_{i=1}^N$  and the densities  $\rho_i^n$  at  $t = n\tau$  and  $\mathbf{x}_i^n$ .
- Find  $\Psi^{n+1}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , by solving the optimization problem (3.14). To guarantee energy stability, we should take  $\Psi$  as an approximation to an identity map initially when solving the optimization problem (3.14).
- After obtaining  $\Psi^{n+1}$ , update  $\{\mathbf{x}_i^{n+1}\}_{i=1}^N$  and  $\rho^{n+1}$  by

$$(3.20) \quad \mathbf{x}_i^{n+1} = \Psi^{n+1}(\mathbf{x}_i^n), \quad \rho_i^{n+1} = \frac{\rho_i^n}{\det(\nabla \Psi^{n+1}(\mathbf{x}_i^n))}.$$


---

The next question is how to accurately evaluate the numerical integration in (3.14). Let  $\mathbf{x}_i^n = \Phi^n(\mathbf{x}_i^0)$ , for the general free energy (2.5). One way to evaluate the integrations in (3.14) is using

$$(3.21) \quad \begin{aligned} J(\Psi) = & \frac{1}{2\tau} \sum_{i=1}^N \rho_i^n \|\Psi(\mathbf{x}_i^n) - \mathbf{x}_i^n\|^2 |\Omega_i^n| + \sum_{i=1}^N \left( \omega \left( \frac{\rho_i^n}{\det(\nabla \Psi(\mathbf{x}_i^n))} \right) \det(\nabla \Psi(\mathbf{x}_i^n)) \right) |\Omega_i^n| \\ & + \sum_{i=1}^N V(\Psi(\mathbf{x}_i^n)) \rho_i^n |\Omega_i^n| + \frac{1}{2} \sum_{i,j=1}^N K(\Psi(\mathbf{x}_i^n), \Psi(\mathbf{x}_j^n)) \rho_i^n \rho_j^n |\Omega_j^n| |\Omega_i^n|, \end{aligned}$$

where  $|\Omega_i^n|$  is the volume of the Voronoi cells associated with the set of points  $\{\mathbf{x}_i^n\}$ ,  $\rho_i^n$  stands for  $\rho^n(\mathbf{x}_i^n)$ . Here the numerical integration is computed through a piecewisely constant reconstruction of  $\rho^n$  based on its values at  $\{\Phi^n(\mathbf{x}_i^0)\}_{i=1}^N$ . However, it is not straightforward to compute  $|\Omega_i^n|$ , particularly for high-dimensional cases. In the current study, we assume that the initial samples are drawn from  $\rho_0$ , one can roughly assume that  $\{\mathbf{x}_i^n\}$  follows the distribution  $\rho^n$ , then according to the Monte-Carlo approach, the numerical integration can be evaluated as

$$(3.22) \quad \begin{aligned} J(\Psi) = & \frac{1}{2\tau} \left( \frac{1}{N} \sum_{i=1}^N \|\Psi(\mathbf{x}_i^n) - \mathbf{x}_i^n\|^2 \right) \\ & + \frac{1}{N} \sum_{i=1}^N \left( f_\omega \left( \frac{\rho_i^n}{\det(\nabla \Psi(\mathbf{x}_i^n))} \right) + V(\Psi(\mathbf{x}_i^n)) \right) + \frac{1}{2N^2} \sum_{i,j=1}^N K(\Psi(\mathbf{x}_i^n), \Psi(\mathbf{x}_j^n)), \end{aligned}$$

where  $f_\omega(\rho) = \omega(\rho)/\rho$ . The proposed numerical method can be further improved if one can evaluate the integration more accurately, i.e., have an efficient way to estimate  $|\Omega_i^n|$ . Alternatively, as an advantage of the neural network-based algorithm, we can get  $\rho^n(\mathbf{x}) = \rho_0 \circ (\Phi^n)^{-1}(\mathbf{x}) = (\rho_0 \circ (\Phi^n)^{-1}(\mathbf{x})) / (\det(\nabla \Phi^n)|_{(\Phi^n)^{-1}(\mathbf{x})}) \forall \mathbf{x}$ , which enables us to estimate  $J(\Psi)$  by resampling. More precisely, assume  $\mu$  is a distribution that is easy to sample

$$(3.23) \quad \begin{aligned} J(\Psi) = & \frac{1}{2\tau} \left( \frac{1}{N} \sum_{i=1}^N \|\Psi(\mathbf{x}_i^n) - \mathbf{x}_i^n\|^2 \right) \frac{\rho_i^n}{\mu_i^n} + \frac{1}{N} \sum_{i=1}^N \left( f_\omega \left( \frac{\rho_i^n}{\det(\nabla \Psi(\mathbf{x}_i^n))} \right) \right) \frac{\rho_i^n}{\mu_i^n} \\ & + V(\Psi(\mathbf{x}_i^n)) \frac{\rho_i^n}{\mu_i^n} + \frac{1}{2N^2} \sum_{i,j=1}^N K(\Psi(\mathbf{x}_i^n), \Psi(\mathbf{x}_j^n)) \frac{\rho_i^n \rho_j^n}{\mu_i^n \mu_j^n}, \end{aligned}$$

where  $\mathbf{x}_i \sim \mu$ ,  $\rho_i^n = \rho^n(\mathbf{x}_i)$ . We will explore this resampling approach in future work.

The remaining question is how to parameterize a diffeomorphism using neural networks. This is discussed in detail in the next subsection.

**3.3. Neural network architectures.** In principle, the proposed numerical framework is independent of the choice of neural network architectures. However, different neural network architectures may lead to different numerical performances, arising from a balance of approximation (representation power), optimization, and generalization. In this subsection, we briefly discuss several neural network architectures that we use in the numerical experiments.

**3.3.1. Neural network architectures for Eulerian methods.** For Eulerian methods, one can construct a neural network to approximate an unknown function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ . Shallow neural networks (two-layer neural networks) approximate  $f$  by functions of the form



$$(3.24) \quad f(\mathbf{x}; \Theta) = \sum_{i=1}^N \alpha_i \sigma(\boldsymbol{\omega}_i \cdot \mathbf{x} + b_i) + \alpha_0 = \boldsymbol{\alpha} \cdot \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) + \alpha_0,$$

where  $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is a fixed nonlinear activation function,  $N$  is the number of hidden nodes (neurons), and  $\Theta = \{\alpha_i, \boldsymbol{\omega}_i, b_i\}$  are the neural network parameters to be identified. Typical choices of activation functions include the ReLU  $\sigma(x) = \max(x, 0)$ , the sigmoid  $\sigma(x) = 1/(1 + e^{-2x})$  and the hyperbolic tangent function  $\tanh(x)$ . A DNN can be viewed as a network composed of many hidden layers. More precisely, a DNN with  $L$  hidden layers represents a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  by [67]

$$(3.25) \quad f(\mathbf{x}; \Theta) = g \circ T^L \circ T^{L-1} \circ \dots \circ T^{(1)}(\mathbf{x}),$$

where  $g(\mathbf{z}) = \sum_{i=1}^{N_L} \gamma_i z_i + \gamma_0$  is a linear map from  $\mathbb{R}^{N_L}$  to  $\mathbb{R}$ ,  $T^{(l)}$  is a nonlinear map from  $\mathbb{R}^{N_{l-1}}$  to  $\mathbb{R}^{N_l}$ , and  $N_0 = d$ . The nonlinear map  $T^l$  takes the form

$$(3.26) \quad T^{(l)}(\mathbf{x}_{l-1}) = \sigma(\mathbf{W}_l \mathbf{x}_{l-1} + \mathbf{b}_l),$$

where  $\mathbf{W}_l \in \mathbb{R}^{N_{l-1} \times N_l}$ ,  $\mathbf{b}_l \in \mathbb{R}^{N_l}$ , and  $\sigma(\cdot)$  is a nonlinear activation function that acts componentwisely on vector-valued inputs.

Another widely used class of DNN model is residual neural network (ResNet). A typical  $K$ -block ResNet approximates an unknown function  $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$(3.27) \quad f_K(\mathbf{x}; \Theta) = g(\mathbf{z}_K(\mathbf{x})),$$

where  $g(\mathbf{z}) = \sum_{i=1}^N \gamma_i z_i + \gamma_0$  is a linear map from  $\mathbb{R}^N \rightarrow \mathbb{R}$  and  $\mathbf{z}_K(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^N$  is a nonlinear map defined through

$$(3.28) \quad \mathbf{z}_0 = \mathbf{V}\mathbf{x}, \quad \mathbf{z}_k = \sigma_2 \left( \boldsymbol{\alpha}_k T_k^{L_k} \circ T_k^{L_k-1} \circ \dots \circ T_k^{(1)}(\mathbf{z}_{k-1}) + \mathbf{z}_{k-1} \right), \quad k = 1, 2, \dots, K.$$

Here  $L_i$  is the number of fully connected layer in  $i$ th block,  $T_k^{(l)}$  is the same nonlinear map defined in (3.26) with  $\mathbf{W}_l^k \in \mathbb{R}^{M \times N}$ ,  $\mathbf{b}_l^k \in \mathbb{R}^M$ ,  $\mathbf{V} \in \mathbb{R}^{N \times d}$ ,  $\boldsymbol{\alpha}_i \in \mathbb{R}^{N \times M}$  and  $\sigma_i(\cdot)$  is an elementwise activation function. The model parameters are  $\Theta = \{\gamma, \boldsymbol{\alpha}_i, \mathbf{W}_l^k, \mathbf{b}_l^k, \mathbf{V}\}$ . The original ResNet [29] takes  $\sigma_2$  as a nonlinear activation function such as ReLU. Later studies indicate that one can also take  $\sigma_2$  as the identity function [18, 30]. Then at an infinite length, i.e.,  $L \rightarrow \infty$ , (3.28) corresponds to the ODE

$$(3.29) \quad \frac{d\mathbf{z}}{dt} = f(\mathbf{z}), \quad \mathbf{z}_0 = \mathbf{x}.$$

Compared with fully connected DNN models, which may suffer from numerical instabilities in the form of exploding or vanishing gradients [19, 28, 40], very deep ResNet can be constructed to avoid these issues.

Given that the proposed numerical scheme employs neural networks with time-dependent parameters to approximate the solution of a gradient flow, there is no need to employ a deep neural network. In all numerical experiments for  $L^2$ -gradient flows, we utilize ResNet with only a few blocks. The detailed settings for each numerical experiment will be described in the next section. We'll compare the numerical performance of different neural network architectures in future work.

**3.3.2. Neural network architectures for Lagrangian methods.** The proposed Lagrangian method seeks for a neural network to approximate a diffeomorphism from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ . This task involves two main challenges: one is ensuring that the map is a diffeomorphism, and the other is efficiently and robustly computing the deformation tensor  $\mathbf{F}$  and its determinant  $\det \mathbf{F}$ . Fortunately, various neural network architectures have been proposed to approximate a transport map. Examples include planar flow [64], auto-regressive flow [39], continuous-time diffusive flow [69], neural spline flow [17], and convex potential flow [32].

One way to construct a neural network  $\mathbb{R}^d \rightarrow \mathbb{R}^d$  for approximating a diffeomorphism is to use a planar flow. A  $K$ -layer planar flow is defined by  $T = T_K \circ \cdots \circ T_1 \circ T_0$ , where  $T_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is given by

$$(3.30) \quad \mathbf{x}^{k+1} = T_k(\mathbf{x}^k) = \mathbf{x}^k + \mathbf{u}_k h(\mathbf{w}_k^T \mathbf{x}^k + b_k).$$

Here,  $\mathbf{w}_k, \mathbf{u}_k \in \mathbb{R}^d$ ,  $b_k \in \mathbb{R}$ , and  $h$  is a smooth, elementwise, nonlinear activation function such as  $\tanh$ . Direct computation shows that

$$(3.31) \quad J_k = \det(\nabla T_k) = 1 + h'(\mathbf{w}_k^T \mathbf{x}^k + b_k) \mathbf{u}_k^T \mathbf{w}_k.$$

Clearly,  $T_k$  is a diffeomorphism if  $h'(\mathbf{w}_k^T \mathbf{x}^k + b_k) \mathbf{u}_k^T \mathbf{w}_k < 1 \ \forall \mathbf{x}_k$ . The determinant of the transport map can be computed as  $\det(\nabla T) = J_K J_{K-1} \cdots J_0$ , and we have  $\varphi(T(\mathbf{x})) = \frac{\varphi^n(\mathbf{x})}{\det(\nabla T)}$ .

One limitation of the planar flow is its potential lack of expressive power. Another commonly employed neural network architecture for approximating a map is the convex potential flow [32], which defines a diffeomorphism via the gradient of a strictly convex function that is parameterized by an Input Convex Neural Network (ICNN) [2]. A fully connected,  $K$ -layer ICNN can be written as

$$(3.32) \quad \mathbf{z}_{l+1} = \sigma_l(\mathbf{W}_l \mathbf{z}_l + \mathbf{A}_l \mathbf{X} + \mathbf{b}_l), \quad l = 0, 1, \dots, K-1,$$

where  $\mathbf{z}_0 = 0$ ,  $\mathbf{W}_0 = 0$ ,  $\Theta = \{\mathbf{W}_l, \mathbf{A}_l, \mathbf{b}_l\}$  are the parameters to be determined, and  $\sigma_l$  is the nonlinear activation function. As proved in [2], if all entries  $\mathbf{W}_l$  are nonnegative, and all  $\sigma_l$  are convex and nondecreasing, then  $f(\mathbf{X}; \Theta)$  is a convex function with respect to  $\mathbf{X}$ . Hence  $\nabla_{\mathbf{X}} f(\mathbf{X}; \Theta)$  provides a parametric approximation to a flow map. In the current study, we adopt the concept of the convex potential flow to develop the Lagrangian EVNN method. We'll explore other types of neural network architectures in future work.

*Remark 3.5.* It is worth mentioning that the gradient of a convex function  $f$  only defines a subspace of diffeomorphism, and it is unclear whether the optimal solution of (3.14) belongs to this subspace.

**4. Numerical experiments.** In this section, we test the proposed EVNN methods for various  $L^2$ -gradient flows and generalized diffusions. To evaluate the accuracy of different methods, we define the  $l^2$ -error

$$\left( \frac{1}{N} \sum_{i=1}^N \left| \varphi_{\text{NN}}(\mathbf{x}_i) - \varphi_{\text{ref}}(\mathbf{x}_i) \right|^2 \right)^{1/2}$$

and the relative  $l_2$ -error

$$\left( \sum_{i=1}^N \left| \varphi_{\text{NN}}(\mathbf{x}_i) - \varphi_{\text{ref}}(\mathbf{x}_i) \right|^2 / \sum_{i=1}^N \left| \varphi_{\text{ref}}(\mathbf{x}_i) \right|^2 \right)^{1/2}$$

between the neural network solution  $\varphi_{\text{NN}}$  and the corresponding reference solution  $\varphi_{\text{ref}}$  on a set of test samples  $\{\mathbf{x}_i\}_{i=1}^N$ . The reference solution  $\varphi_{\text{ref}}(\mathbf{x})$  is either an analytic solution or a numerical solution obtained by a traditional numerical method. In some numerical examples, we also plot pointwise absolute errors  $|\varphi_{\text{ref}} - \varphi_{\text{NN}}|$  at the grid points.

**4.1. Poisson equations.** Although the proposed method is designed for evolutionary equations, it can also be used to compute equilibrium solutions for elliptic problems. In this subsection, we compare the performance of the EVNN method with two classical neural network-based algorithms, PINN and DRM, in the context of solving Poisson equations.

We first consider a two-dimensional (2D) Poisson equation with a Dirichlet boundary condition

$$(4.1) \quad \begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} \in \Omega \subset \mathbb{R}^d, \\ u(\mathbf{x}) = g(\mathbf{x}), & \mathbf{x} \in \partial\Omega. \end{cases}$$

Since EVNN is developed for evolution equations, we solve the following  $L^2$ -gradient flow

$$(4.2) \quad \frac{d}{dt} \left( \int_{\Omega} \frac{1}{2} |\nabla u|^2 - f(\mathbf{x})u(\mathbf{x}) d\mathbf{x} + \lambda \int_{\partial\Omega} |u(\mathbf{x}) - g(\mathbf{x})|^2 dS \right) = - \int_{\Omega} |u_t|^2 d\mathbf{x}$$

to get a solution of the Poisson equation (4.1). Here,  $\lambda \int_{\partial\Omega} |u(\mathbf{x}) - g(\mathbf{x})|^2 dS$  is the surface energy that enforces the Dirichlet boundary condition. The corresponding gradient flow equation is given by

$$(4.3) \quad u_t = \Delta u(\mathbf{x}) + f(\mathbf{x}),$$

subject to a Robin boundary condition  $u(\mathbf{x}) - g(\mathbf{x}) = \frac{1}{2\lambda} \frac{\partial u}{\partial \mathbf{n}}$ , where  $\mathbf{n}$  is the outer normal of  $\partial\Omega$ . One can recover the original Dirichlet boundary condition by letting  $\lambda \rightarrow \infty$ . Such a penalty approach is also used in PINN and DRM, and we take  $\lambda = 500$  in all numerical experiments below.

We consider the following two cases:

- Case 1:  $\Omega = (0, \pi) \times (-\pi/2, \pi/2)$ ,  $f(\mathbf{x}) = 2 \sin x \cos y$  and  $g(\mathbf{x}) = 0$ . The exact solution is  $u(\mathbf{x}) = \sin x \cos y$ .
- Case 2:  $\Omega = \{(x, y) \mid |x| \leq 1\}$ ,  $f(\mathbf{x}) = \frac{\pi^2}{4} \sin\left(\frac{\pi}{2}(1 - |\mathbf{x}|)\right) + \frac{\pi}{2|\mathbf{x}|} \cos\left(\frac{\pi}{2}(1 - |\mathbf{x}|)\right)$  and  $g(\mathbf{x}) = 0$ . The exact solution is  $u(\mathbf{x}) = \sin\left(\frac{\pi}{2}(1 - |\mathbf{x}|)\right)$ .

We employ a 1-block ResNet with 20 hidden nodes and one hidden layer for all cases. The total number of parameters is 501. We apply Xavier Initialization [26] to initialize the weights of neural networks in all cases. To evaluate the integration in all methods, we generate 2500 samples in  $\Omega$  using a Latin hypercube sampling (LHS) and 50 samples on each boundary of  $\partial\Omega$  using a one-dimensional (1D) uniform sampling for case 1. For case 2, we generate 2500 samples in  $(-1, 1)^2$  using LHS, but only use the samples satisfies  $x^2 + y^2 < 1$  as training samples. Additionally, we generate 200 training samples  $(\cos \theta_i, \sin \theta_i)_{i=1}^{200}$  on the boundary, with  $\{\theta_i\}_{i=1}^{200}$  being generated by a uniform distribution on  $(0, 2\pi)$ . For PINN and DRM, we use Adam to minimize the loss function and use a different set of samples at each iteration. For EVNN, we use a different set of samples at each time step and employ an L-BFGS to solve the optimization problem (3.10).

Due to randomness arising from the stochastic initialization of the neural network and the sampling process, we repeated each method 10 times and plotted the mean and standard error of the relative  $l_2$ -errors. Figure 4.1(a) shows the relative  $l_2$ -errors

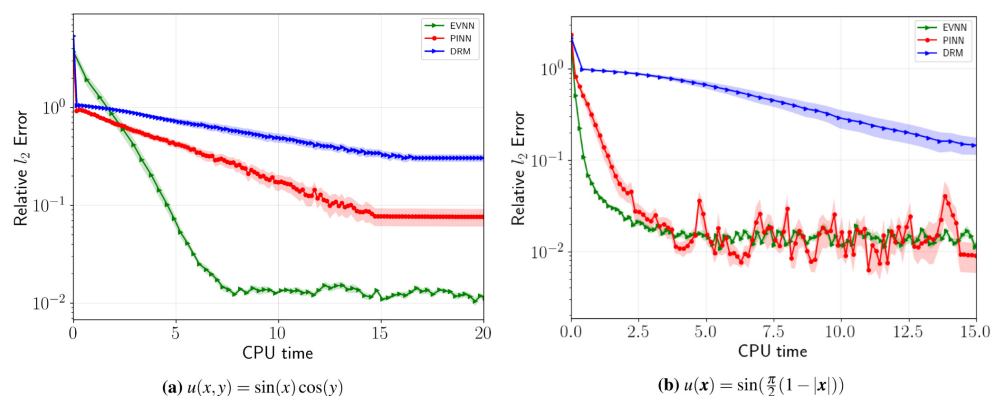


FIG. 4.1. Relative  $l^2$  error with respect to CPU time for different methods for 2D Poisson equations.

of different methods with respect to the CPU time for both cases. The training was executed on a MacBook Pro equipped with an M2 chip. In case 1, the relative  $l^2$  error of each method is evaluated on a test set comprising a uniform grid of  $101 \times 101$  points within the domain  $(0, 1)^2$ . For case 2, the test set is defined as the collection of points  $\{(x_i, y_j) | x_i^2 + y_j^2 < 1\}$ , where  $(x_i, y_j)$  forms a uniform  $201 \times 201$  grid on  $(-1, 1)^2$ . The CPU time is computed by averaging the time of 10 trials at  $n$ th iteration. The numerical results clearly show that the proposed method is more efficient than both PINN and DRM. It significantly enhances the efficiency and test accuracy of DRM through the introduction of time relaxation. While PINN may achieve better test accuracy in case 2, it exhibits a larger standard error in both cases.

Next, we consider a high-dimensional Poisson equation

$$(4.4) \quad -\Delta u = f(\mathbf{x}), \quad \mathbf{x} \in \Omega = (-1, 1)^d,$$

with a homogeneous Neumann boundary condition  $\frac{\partial u}{\partial n}|_{\partial\Omega} = 0$ . We take  $f(\mathbf{x}) = \pi^2 \sum_{k=1}^d \cos(\pi x_k)$ . Similar numerical examples are tested in [21, 53]. The exact solution to this high-dimensional Poisson equation is  $u(\mathbf{x}) = \sum_{k=1}^d \cos(\pi x_k)$ . Following [53], we solve an  $L^2$ -gradient flow associated with the free energy

$$(4.5) \quad \mathcal{F}[u] = \int_{(-1,1)^d} \frac{1}{2} |\nabla u|^2 - fu \, d\mathbf{x} + \lambda \left( \int_{(-1,1)^d} u \, d\mathbf{x} \right)^2,$$

where the last term enforces  $\int_{[-1,1]^d} u \, d\mathbf{x} = 0$ . We take  $\lambda = 0.5$  as in [53] in all numerical tests.

Since it is difficult to get an accurate estimation of high-dimensional integration, we choose Adam to minimize (3.10) at each time step and use a different set of samples at each Adam step. This approach allows us to explore the high-dimensional space with a low computational cost. The samples are drawn using LHS. Table 4.1 shows the final relative  $l^2$ -error with different neural network settings in different spatial dimensions. The relative  $l^2$ -error is evaluated on a test set that comprises 40000 samples generated by LHS. The model is trained with 200 outer iterations ( $\tau = 0.01$ ) and at most 200 iterations for inner optimization (3.10). An early stop criterion was applied if the  $l_2$ -norm of the model parameters between two consecutive epochs is less than  $10^{-6}$ . The first column of Table 4.1 specifies the numerical setting as follows: the

TABLE 4.1

Relative  $l^2$ -error of high-dimensional Poisson equation with different settings in different dimensions. The first column shows the number of residual blocks, the number of nodes in each fully connected layer, and the number of samples.

Setting	$d = 4$	$d = 8$	$d = 16$	$d = 32$
(2, 10, 1000)	0.024	0.046	0.623	0.867
(3, 60, 1000)	0.042	0.048	0.077	0.117
(3, 60, 10000)	0.018	0.022	0.036	0.077

entry (2, 10) means a ResNet with two residual blocks, each with two fully connected layers having 10 nodes is used, and 1000 represents the number of samples drawn in each epoch. As we can see, the proposed method can achieve comparable results with results reported in similar examples in previous work by DRM [21, 53]. It can be observed that increasing the width of the neural network improves test accuracy in high dimensions, as it enhances the network's expressive power. Furthermore, increasing the number of training samples also significantly improves test accuracy in high dimensions. We'll explore the effects of  $\tau$ , the number of samples, and neural network architecture for the high-dimensional Poisson equation in future work.

**4.2.  $L^2$ -gradient flow.** In this subsection, we apply the EVNN scheme to two  $L^2$ -gradient flows to demonstrate its energy stability and numerical accuracy.

**4.2.1. Heat equation.** We first consider an initial-boundary value problem of a heat equation:

$$(4.6) \quad \begin{cases} u_t = \Delta u(\mathbf{x}), & \mathbf{x} \in \Omega = (0, 2)^2, \quad t \in (0, T], \\ u(\mathbf{x}, t) = 0, & \mathbf{x} \in \partial\Omega, \quad t \in (0, T], \\ u(\mathbf{x}, 0) = \sin\left(\frac{\pi}{2}x_1\right) \sin\left(\frac{\pi}{2}x_2\right), & \mathbf{x} \in \Omega = (0, 2)^2, \end{cases}$$

which can be interpreted as an  $L^2$ -gradient flow satisfying the energy-dissipation law

$$(4.7) \quad \frac{d}{dt} \left( \int_{\Omega} \frac{1}{2} |\nabla u|^2 d\mathbf{x} + \lambda \int_{\partial\Omega} |u|^2 dS \right) = - \int_{\Omega} |u_t|^2 d\mathbf{x}.$$

Again, a surface energy term is added to enforce the Dirichlet boundary condition.

In the numerical simulation, we take  $\tau = 0.01$  and use a 1-block ResNet with 20 nodes in each layer. The nonlinear activation function is chosen to be tanh. The total number of parameters is 501. To achieve energy stability, we fix the training samples to eliminate the numerical fluctuation arising from estimating the integration in (3.10). The set of training samples comprises a  $301 \times 301$  uniform grid on  $(0, 2)^2$ , and an additional 1000 uniformly spaced grid points on each edge of the boundary. To test the accuracy of the solution over time, we compare the EVNN solution with a numerical solution obtained by a finite difference method (FDM). We apply the standard central finite difference and an implicit Euler method to obtain the FDM solution. The numerical scheme can be reformulated as an optimization problem:

$$(4.8) \quad \phi_h^{n+1} = \arg \min_{u_h \in \mathcal{C}} \frac{1}{2\tau} \|u_h - u_h^n\|_h^2 + \mathcal{F}_h(u_h), \quad \mathcal{F}_h(u_h) = \frac{1}{2} \langle \nabla_h u_h, \nabla_h u_h \rangle_*.$$

Here,  $\mathcal{C} = \{u_{i,j} \mid 0 < i < I, \quad 0 < j < J\}$  is the set of grid functions defined on a uniform mesh,  $\nabla_h$  is the discrete gradient,  $\mathcal{F}_h$  is the discrete free energy,  $\|\cdot\|_h$  is the discrete  $L_2$ -norm induced by the discrete  $L_2$ -inner product  $\langle f, g \rangle = h^2 \sum_{i=1}^{N-1} f_{i,j} g_{i,j}$ , and  $\langle \cdot, \cdot \rangle_*$

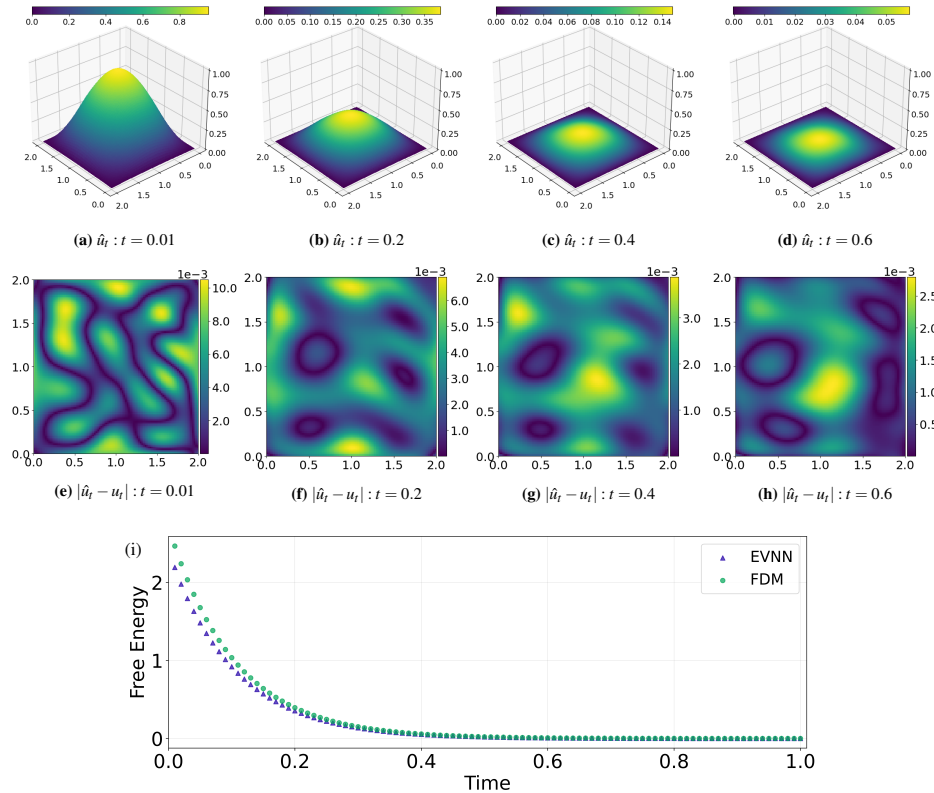


FIG. 4.2. Numerical results for the heat equation. (a)–(d) are the EVNN solutions at  $t = 0.01, 0.2, 0.4,$  and  $0.6$ . (e)–(h) are the absolute differences between the EVNN solutions and the FDM solutions. (i) The evolution of discrete free energy with respect to time for both the EVNN and FDM solutions.

is the discrete inner product defined on the staggered mesh points [49, 66]. We use a  $101 \times 101$  uniform grid and take  $\tau = 0.01$  in the FDM simulation. Figure 4.2(a)–(d) shows the EVNN solution at  $t = 0.01, 0.2, 0.4,$  and  $0.6$ . The corresponding pointwise absolute errors are shown in Figure 4.2(e)–(h). The evolution of discrete free energy (evaluated on  $101 \times 101$  uniform grid) for both methods is shown in Figure 4.2(i). Clearly, the neural-network-based algorithm achieves energy stability and the result is consistent with the FDM. It is worth mentioning that the size of the optimization problem in the neural network-based algorithm is much smaller than that in the FDM, although (4.8) is a quadratic optimization problem and can be solved directly by solving a linear equation of  $u_{i,j}$ .

**4.2.2. Allen–Cahn equation.** Next, we consider an Allen–Cahn type equation, which is extensively utilized in phase-field modeling, and becomes a versatile technique to solve interface problems arising from different disciplines [15]. In particular, we focus on the following Allen–Cahn equation on  $\Omega = (-1, 1)^2$ :

$$(4.9) \quad \begin{cases} \frac{\partial \varphi}{\partial t}(\mathbf{x}, t) = - \left( \frac{F'(\varphi)}{\varepsilon^2} - \Delta \varphi(\mathbf{x}, t) + 2W \left( \int \varphi d\mathbf{x} - A \right) \right), & \mathbf{x} \in \Omega, \quad t > 0, \\ \varphi(\mathbf{x}, t) = -1, & \mathbf{x} \in \partial\Omega, \quad t > 0, \\ \varphi(\mathbf{x}, 0) = -\tanh(10(\sqrt{4x_1^2 + x_2^2} - 0.5)), & \mathbf{x} \in \Omega. \end{cases}$$

where  $\varphi(\mathbf{x}, t)$  is the phase-field variable, constants  $\epsilon$ ,  $W$ , and  $A$  are model parameters. The Allen–Cahn equation (4.9) can be regarded as an  $L^2$ -gradient flow, derived from an energy-dissipation law

(4.10)

$$\frac{d}{dt}\mathcal{F}[\varphi] = - \int_{\Omega} |\varphi_t|^2 d\mathbf{x}, \quad \mathcal{F}[\varphi] = \int_{\Omega} \frac{1}{2} |\nabla \varphi|^2 + \frac{1}{4\epsilon^2} (\varphi^2 - 1)^2 d\mathbf{x} + W \left( \int_{\Omega} \varphi d\mathbf{x} - A \right)^2.$$

Here, the last term in the energy is a penalty term for the volume constraint  $\int \varphi d\mathbf{x} = A$ , as the standard Allen–Cahn equation does not preserve the volume fraction  $\int \varphi d\mathbf{x}$  over time. In the numerical simulation, we take  $A = -(4 - \pi r^2) + \pi r^2$ ,  $r = 0.5$ ,  $\frac{1}{\epsilon^2} = 100$ , and  $W = 1000$ .

Due to the complexity of the initial condition, we utilized a larger neural network compared to the one employed in the previous subsection. Our choice was a 1-block ResNet, which has two fully connected layers, each containing 20 nodes. The total number of parameters is 921. The nonlinear activation function is chosen to be tanh. The set of training samples comprises a  $301 \times 301$  uniform grid on  $(-1, 1)^2$ , and an additional 1000 uniformly spaced grid points on each edge of the boundary. To test the numerical accuracy of the EVNN scheme for this problem, we also solve (4.9) by a finite element method, which approximates the phase-field variable  $\varphi$  by a piecewise linear function  $\varphi_h(\mathbf{x}, t) = \sum_{i=1}^N \gamma_i(t) \psi_i(\mathbf{x})$ , where  $\psi_i(X)$  are hat functions supported on the computational mesh. Inserting  $\varphi_h$  into the continuous energy–dissipation law (4.10), we get a discrete energy–dissipation law with the discrete energy and dissipation given by

$$\mathcal{F}_h^{\text{FEM}}(\gamma) = \sum_{e=1}^{N_e} \int_{\tau_e} \frac{1}{2} \left| \sum_{i=1}^N \gamma_i \nabla \psi_i(\mathbf{x}) \right|^2 + \frac{1}{\epsilon^2} \sum_{i=1}^N (\gamma_i^2 - 1)^2 \psi_i(\mathbf{x}) d\mathbf{x},$$

and  $\mathcal{D}_h^{\text{FEM}}(\gamma, \gamma') = \sum_{e=1}^{N_e} \int_{\tau_e} \left| \sum_{i=1}^N \gamma'_i(t) \psi_i(\mathbf{x}) \right|^2 d\mathbf{x}$ , respectively. Here  $\tau_e$  is used to denote a finite element cell, and  $N_e$  is the number of cells. This form of discretization was used in [79], which constructs a piecewise linear approximation to the nonlinear term in the discrete energy. We can update  $\gamma_i$  at each time step by solving the optimization problem (3.7), i.e.,

$$(4.11) \quad \gamma^{n+1} = \arg \min_{\gamma \in \mathbb{R}^N} \frac{1}{2\tau} D(\gamma - \gamma^n) \cdot (\gamma - \gamma^n) + \mathcal{F}_h(\gamma),$$

where  $D_{ij} = \int_{\Omega} \psi_i \psi_j d\mathbf{x}$  is the mass matrix. The optimization problem (4.11) is solved by L-BFGS in our numerical implementation.

The simulation results are summarized in Figure 4.3. It clearly shows that our method can achieve comparable results with the FEM. Numerical simulation of phase-field type models is often challenging. To capture the dynamics of the evolution of the thin diffuse interface, the mesh size should be much smaller than  $\epsilon$ , the width of the diffuse interface. Traditional numerical methods often use an adaptive or moving mesh approach to overcome this difficulty. In contrast, a neural network-based numerical scheme has a mesh-free feature. The number of parameters of the neural network can be much smaller than the number of samples needed to resolve the diffuse interface. Consequently, the dimension of the optimization problem in the neural network-based scheme is much smaller than that in the FEM scheme without using adaptive or moving meshes.

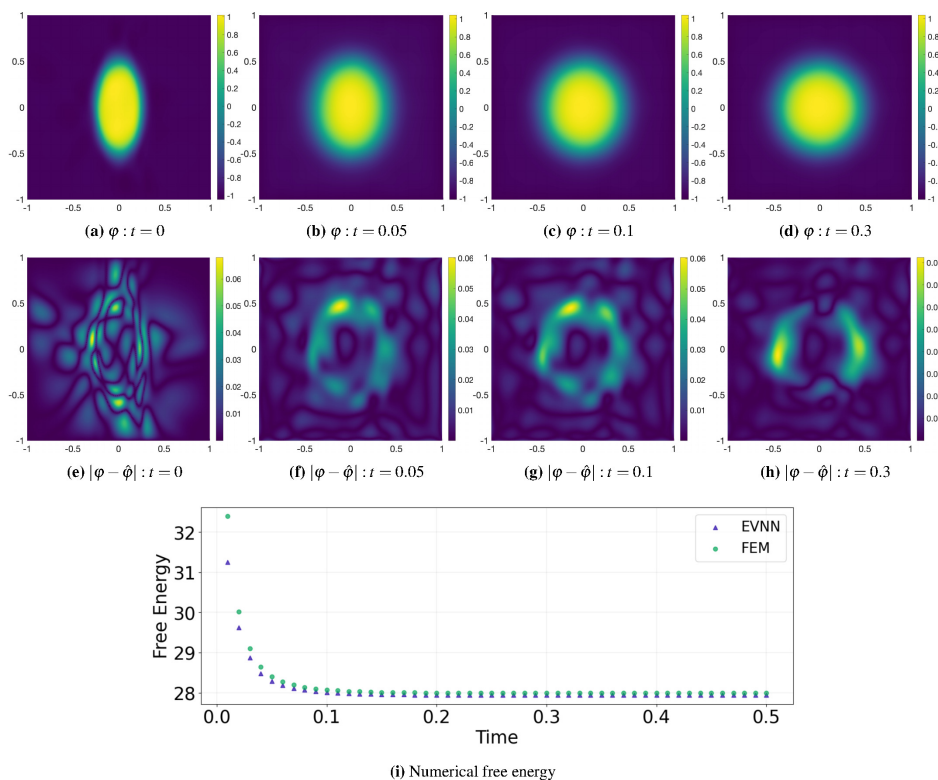


FIG. 4.3. Numerical results for the Allen–Cahn equation with a volume constraint. (a)–(d) EVNN solutions at  $t = 0, 0.05, 0.1,$  and  $0.3$ , respectively. (e)–(h) Absolute differences between the EVNN solutions and the FEM solutions at  $t = 0, 0.05, 0.1,$  and  $0.3$ , respectively. (i) Evolution of numerical free energies with respect to time for both the FEM and EVNN solutions.

**4.3. Generalized diffusions.** In this subsection, we apply the proposed Lagrangian EVNN scheme to solving a Fokker–Planck equation and a porous medium equation.

**4.3.1. Fokker–Planck equation.** We first consider a Fokker–Planck equation

$$(4.12) \quad \begin{cases} \rho_t = \nabla \cdot (\nabla \rho + \rho \nabla V), & \mathbf{x} \in \mathbb{R}^d, \quad t \in (0, T], \\ \rho(\mathbf{x}, 0) = \rho_0(\mathbf{x}), & \mathbf{x} \in \mathbb{R}^d, \end{cases}$$

where  $V(\mathbf{x})$  is a prescribed potential energy. The Fokker–Planck equation can be viewed as a generalized diffusion satisfying the energy-dissipation law

$$(4.13) \quad \frac{d}{dt} \int_{\mathbb{R}^d} \rho \ln \rho + \rho V(\mathbf{x}) d\mathbf{x} = - \int_{\mathbb{R}^d} \rho |\mathbf{u}|^2 d\mathbf{x},$$

where the probability density  $\rho$  satisfies the kinematics  $\rho_t + \nabla \cdot (\rho \mathbf{u}) = 0$ .

We test our numerical scheme for solving the Fokker–Planck equation in two dimensions and four dimensions, respectively. For both numerical experiments, we adopt the augmented-ICNN proposed in [32]. We use a 1 ICNN block of six fully connected layers with 32 hidden nodes in each layer. The activation function is chosen as the Gaussian-Softplus function  $\sigma(x) = \sqrt{\frac{2}{\pi}} (\sqrt{2}x \int_0^{x/\sqrt{2}} e^{-t^2} dt + e^{-\frac{x^2}{2}} + \sqrt{\frac{\pi}{2}}x)$ . In addition, we use L-BFGS to solve the optimization problem at each time step.



In the 2D case, we consider  $V(\mathbf{x})$  as follows:

$$V = \frac{1}{2}(\mathbf{x} - \mu_{\text{target}})^T \Sigma_{\text{target}}^{-1}(\mathbf{x} - \mu_{\text{target}}) \quad \text{with} \quad \mu_{\text{target}} = \left(\frac{1}{3}, \frac{1}{3}\right),$$

and  $\Sigma_{\text{target}} = \begin{bmatrix} \frac{5}{8} & -\frac{3}{8} \\ -\frac{3}{8} & \frac{5}{8} \end{bmatrix}$ . The initial condition  $\rho_0(\mathbf{x})$  is set to be the 2D standard Gaussian  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . The exact solution of (4.12) takes the following analytical form [33]:

$$(4.14) \quad \rho(\mathbf{x}, t) \sim \mathcal{N}(\mu(t), \Sigma(t)),$$

where  $\mu(t) = (1 - e^{-4t})\mu_{\text{target}}$  and  $\Sigma(t) = \begin{bmatrix} \frac{5}{8} + \frac{3}{8} \times e^{-8t} & -\frac{3}{8} + \frac{3}{8} \times e^{-8t} \\ -\frac{3}{8} + \frac{3}{8} \times e^{-8t} & \frac{5}{8} + \frac{3}{8} \times e^{-8t} \end{bmatrix}$ . We draw 10000 samples from  $\rho_0(\mathbf{x})$  as the training set. As an advantage of the neural network-based algorithm, we can compute  $\rho^n(\mathbf{x})$  point-wisely through a function composition  $\rho^n(\mathbf{x}) = \rho_0 \tilde{\circ} (\Phi^n)^{-1}(\mathbf{x})$ , where  $(\Phi^n)^{-1}(\mathbf{x})$  can be computed by solving a convex optimization problem, i.e.,  $\Phi^{-1}(\mathbf{x}) = \arg \min_{\mathbf{y}} \Phi(\mathbf{y}) - \mathbf{x}^T \mathbf{y}$ . Figure 4.4(a) - (d) shows the predicted density on a  $301 \times 301$  uniform grid on  $(-3, 3)^2$ . The absolute and relative  $l_2$ -errors on the grid are shown in Figures 4.4(e)-(f). It can be noticed that the relative  $l^2$ -error for the predicted solution is around  $10^{-2}$ , which is quite accurate given the limited number of samples used.

In the 4D case, we take  $V(\mathbf{x})$  to be

$$V = \frac{1}{2}(\mathbf{x} - \mu_{\text{target}})^T \Sigma_{\text{target}}^{-1}(\mathbf{x} - \mu_{\text{target}}) \quad \text{with} \quad \mu_{\text{target}} = \left(\frac{1}{3}, \frac{1}{3}, 0, 0\right),$$

and  $\Sigma_{\text{target}} = \begin{bmatrix} \frac{5}{8} & -\frac{3}{8} \\ -\frac{3}{8} & \frac{5}{8} \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \oplus \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ . The initial condition  $\rho_0(\mathbf{x})$  is set to be a 4D standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . The exact solution of (4.12) follows a normal distribution with  $\mu(t) = ((1 - e^{-4t})\frac{1}{3}, (1 - e^{-4t})\frac{1}{3}, 0, 0)$ , and  $\Sigma(t) = \begin{bmatrix} \frac{5}{8} + \frac{3}{8} \times e^{-8t} & -\frac{3}{8} + \frac{3}{8} \times e^{-8t} \\ -\frac{3}{8} + \frac{3}{8} \times e^{-8t} & \frac{5}{8} + \frac{3}{8} \times e^{-8t} \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \oplus$

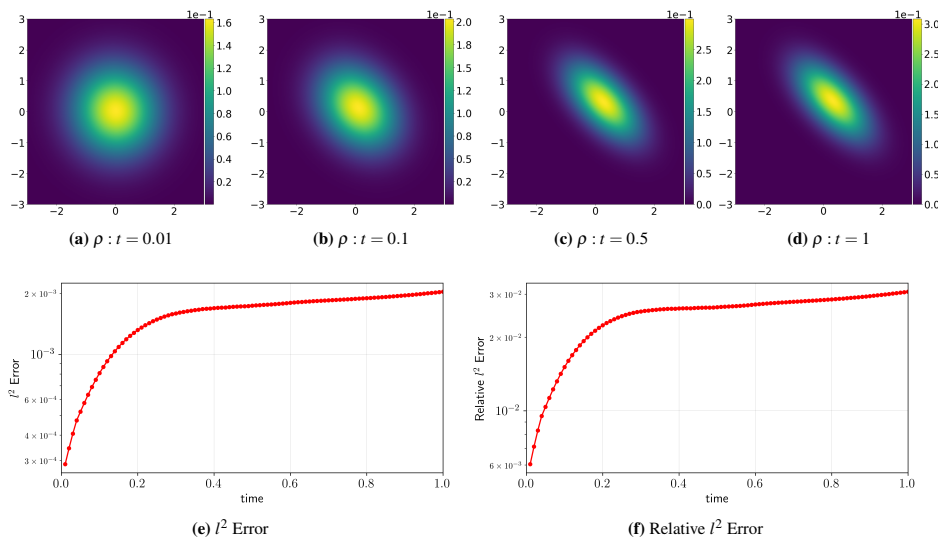


FIG. 4.4. Numerical results for the 2D Fokker-Planck equation. (a)–(d) Predicted solution on a  $301 \times 301$  uniform mesh on  $(-3, 3)^2$  at  $t = 0, 0.1, 0.5$ , and  $1$ , respectively. (e)–(f) The absolute and relative  $l_2$ -errors of the solution over time, evaluated on the uniform mesh.

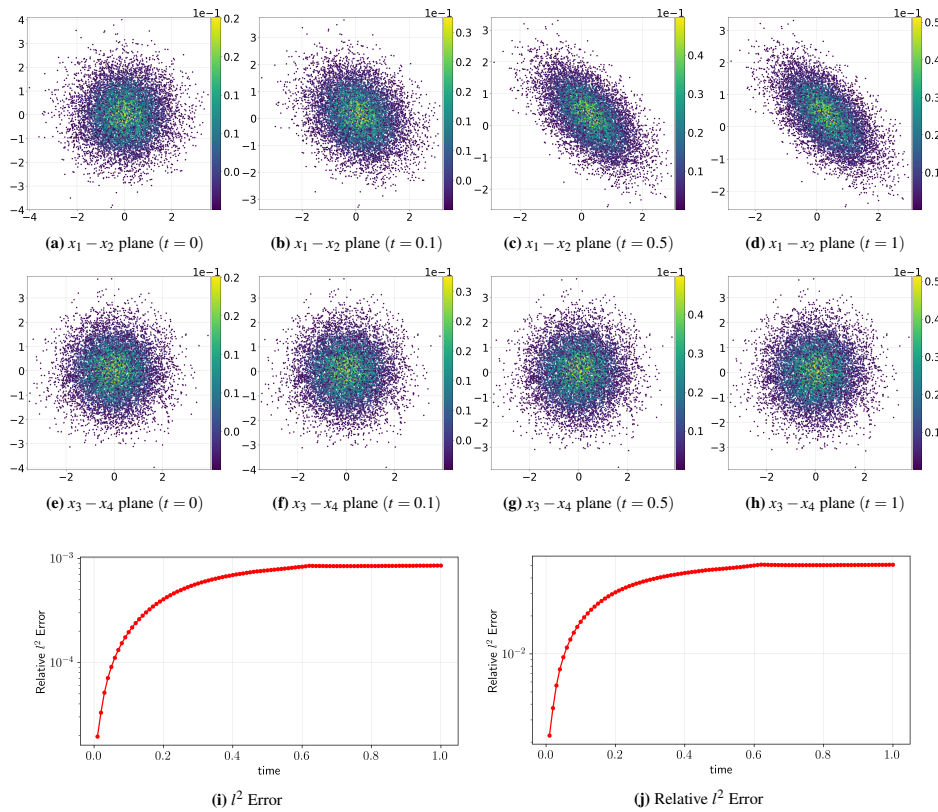


FIG. 4.5. Numerical results for the four-dimensional (4D) Fokker-Planck equation. (a)–(d) Distribution of samples, projected in  $x_1$ - $x_2$  plane, as well as its weight at  $t = 0, 0.1, 0.5$ , and  $1$ , respectively. (e)–(h) Distribution of samples, projected in  $x_3$ - $x_4$  plane, as well as its weight at  $t = 0, 0.1, 0.5$ , and  $1$ , respectively. (i)–(j) The  $l^2$ -errors and the relative  $l^2$ -errors of the solution over time, evaluated in the training set.

$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ . Here  $\oplus$  is the direct sum. Same as the 2D case, we draw 10000 initial samples from  $\rho_0(\mathbf{x})$ . Figure 4.5 shows the evolution of samples as well the values of the numerical solution on individual samples over space and time. Interestingly, the relative  $l^2$ -error in 4D case is similar to the 2D case. The dimension-independent error bound suggests the potential of the current method for handling higher-dimensional problems.

Fokker-Planck type equations have a wide range of application in machine learning. One of the fundamental tasks in modern statistics and machine learning is to estimate or generate samples from a target distribution  $\rho^*(\mathbf{x})$ , which might be completely known, partially known up to a normalizing constant, or empirically given by samples. Examples include Bayesian inference [7], numerical integration [54], space-filling design [61], density estimation [69], and generative learning [59]. These problems can be transformed as an optimization problem, which is to seek for a  $\rho^{\text{opt}} \in \mathcal{Q}$  by solving an optimization problem

$$(4.15) \quad \rho^{\text{opt}} = \arg \min_{\rho \in \mathcal{Q}} D(\rho || \rho^*),$$

where  $\mathcal{Q}$  is the admissible set,  $D(p||q)$  is a dissimilarity function that assesses the differences between two probability measures  $p$  and  $q$ . The classical dissimilarities

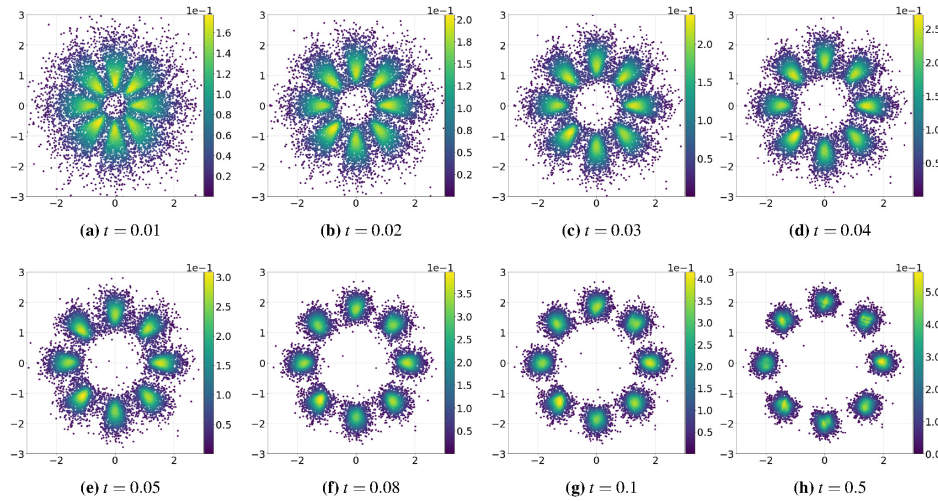


FIG. 4.6. The distribution of samples and their weight at different times for the Fokker–Planck equation with a complicated  $V(\mathbf{x})$ , which corresponds to an eight-component mixture Gaussian distribution. The initial 10,000 samples are drawn from a standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .

include the Kullback–Leibler (KL) divergence and the Maximum Mean Discrepancy (MMD). The optimal solution of the optimization problem can be obtained by solving a Fokker–Planck type equation, given by

$$(4.16) \quad \rho_t = \nabla \cdot (\rho \nabla \mu), \quad \mu = \frac{\delta D(\rho || \rho^*)}{\delta \rho}.$$

The developed numerical approach has potential applications in these machine learning problems.

To illustrate this point, we consider a toy problem that is widely used in the machine learning literature [32]. The goal is to sample from a target distribution  $\rho^*$ , which is known up to a normalizing constant. We take the dissimilarity function as the KL divergence  $\text{KL}(\rho || \rho^*) = \int_{\Omega} \rho(\mathbf{x}) \ln(\frac{\rho}{\rho^*}) d\mathbf{x}$ , then (4.16) is reduced to the Fokker–Planck equation (4.12) with  $V(\mathbf{x}) = -\ln \rho^*$ . In the numerical experiment, we take the target distribution  $\rho^*(\mathbf{x}) = \frac{1}{8} \sum_{i=1}^8 \mathcal{N}(\mathbf{x} | \mu_i, \Sigma)$ , an eight component mixture Gaussian distribution, where  $\mu_1 = (0, 4)$ ,  $\mu_2 = (2.8, 2.8)$ ,  $\mu_3 = (4, 0)$ ,  $\mu_4 = (-2.8, 2.8)$ ,  $\mu_5 = (-4, 0)$ ,  $\mu_6 = (-2.8, -2.8)$ ,  $\mu_7 = (0, -4)$ ,  $\mu_8 = (2.8, -2.8)$ , and  $\Sigma = \text{diag}(0.2, 0.2)$ . The simulation results are summarized in Figure 4.6, in which we first draw 10000 samples from a standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and show the distribution of samples as well as their weights at different time. It can be noticed that the proposed neural network-based algorithm can generate a weighted sample from a complicated target distribution.

**4.3.2. Porous medium equation.** Next, we consider a porous medium equation (PME),  $\rho_t = \Delta \rho^\alpha$ , where  $\alpha > 1$  is a constant. The PME is a typical example of nonlinear diffusion equations. One important feature of the PME is that the solution of the PME has a compact support at any time  $t > 0$  if the initial data has a compact support. The free boundary of the compact support moves outward with a finite speed, known as the property of finite speed propagation [71]. As a consequence, numerical simulations of the PME are often difficult to conduct using Eulerian methods, which may fail to capture the movement of the free boundary and suffer from

numerical oscillations [50]. In a recent work [50], the authors developed a variational Lagrangian scheme using a finite element method. Here we show the ability of the Lagrangian EVNN scheme to solve the PME with a free boundary. Following [50], we employ the energy-dissipation law

$$(4.17) \quad \frac{d}{dt} \int_{\mathbb{R}^d} \frac{\alpha}{(\alpha-1)(\alpha-2)} \rho^{\alpha-1} d\mathbf{x} = - \int_{\mathbb{R}^d} |\mathbf{u}|^2 d\mathbf{x}$$

to develop the EVNN scheme.

We take  $\alpha = 4$  in the simulation. To test the numerical accuracy of the EVNN scheme for solving the PME, we consider a 2D Barenblatt–Pattle solution of the PME. The Barenblatt–Pattle solution in  $d$ -dimensional space is given by

$$(4.18) \quad B_\alpha(\mathbf{x}, t) = t^{-k} \left[ \left( C_0 - \frac{k(\alpha-1)}{2d\alpha} \frac{|\mathbf{x}|^2}{t^{2k/d}} \right)_+ \right]^{1/(\alpha-1)}, \quad \mathbf{x} \in \mathbb{R}^d,$$

where  $k = (\alpha - 1 + 2/d)^{-1}$ ,  $u_+ := \max(u, 0)$ , and  $C_0$  is a constant that is related to the initial mass. This solution is radially symmetric, self-similar, and has compact support  $|\mathbf{x}| \leq \xi_\alpha(t)$  for any finite time, where  $\xi_\alpha(t) = \sqrt{\frac{2d\alpha C_0}{k(\alpha-1)}} t^{k/d}$ . We take the Barenblatt solution (4.18) with  $C_0 = 0.1$  at  $t = 0.1$  as the initial condition, and compare the numerical solution at  $T$  with the exact solution  $B_\alpha(\mathbf{x}, T + 0.1)$ . We choose  $\{x_i^0\}_{i=1}^N$  as a set of uniform grid points inside a disk  $\{(x, y) \mid \sqrt{x^2 + y^2} \leq \xi_\alpha(0.1)\}$ . To visualize the numerical results, we also draw 500 points distributed with equal arc length on the initial free boundary and evolve them using the trained ICNN. The numerical results with  $\tau = 0.005$  are shown in Figure 4.7. It can be noticed that the proposed neural network-based algorithm can well approximate the Barenblatt–Pattle solution and capture the movement of the free boundary.

We note that numerical methods to solve generalized diffusions (1.4) can also be formulated in the Eulerian frame of reference based on the notion of Wasserstein gradient flow [35]. However, it is often challenging to compute the Wasserstein distance

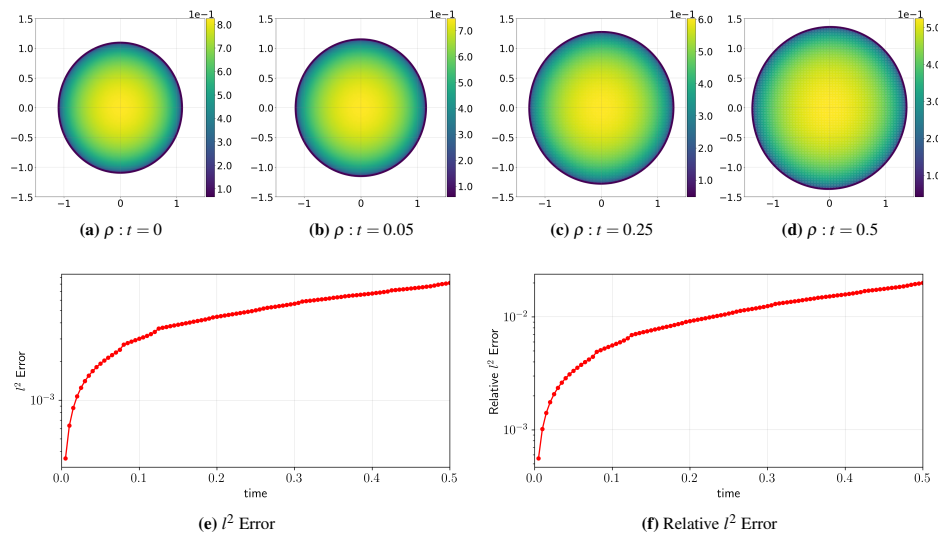


FIG. 4.7. Numerical results for the porous medium equation. (a)–(d) The numerical solution at  $t = 0, 0.05, 0.25$ , and  $0.5$ , respectively. (e)–(f) The  $l_2$ -errors of the solution.

between two probability densities efficiently with high accuracy, which often requires solving an additional min-max problem or minimization problem [5, 6, 9, 33]. In a recent study [33], the authors initially employ a fully connected neural network to approximate  $\rho$  and subsequently utilize an additional ICNN to calculate the Wasserstein distance between two probability densities. Compared with their approach, our method is indeed more efficient and accurate.

**5. Conclusion.** In this paper, we develop structure-preserving EVNN schemes for simulating  $L^2$ -gradient flows and generalized diffusions by utilizing neural networks as a tool for spatial discretization, within the framework of the discrete energetic variational approach. These numerical schemes are directly constructed based on a prescribed continuous energy-dissipation law for the system, without the need for the underlying PDE. The incorporation of mesh-free neural network discretization opens up exciting possibilities for tackling high-dimensional gradient flows arising in different applications in the future.

Various numerical experiments are presented to demonstrate the accuracy and energy stability of the proposed numerical schemes. In our future work, we will explore the effects of different neural network architectures, sampling strategies, and optimization methods, followed by a detailed numerical analysis. Additionally, we intend to employ EVNN schemes to investigate other complex fluid models, including the Cahn–Hilliard equation and Cahn–Hilliard–Navier–Stokes equations, as well as solve machine learning problems such as generative modeling and density estimation.

**Appendix A.** In this appendix, we provide the detailed computation for the EnVarA derivation of a generalized diffusion. Due to the kinematics relation (2.4), the action functional  $\mathcal{A}[\mathbf{x}]$  can be written as

$$\begin{aligned} \mathcal{A}[\mathbf{x}] = & - \int_0^T \int_{\Omega_0} \omega \left( \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}} \right) \det \mathbf{F} + V(\mathbf{x}(\mathbf{X}, t)) \rho_0(\mathbf{X}) \\ & + \frac{1}{2} \rho_0(\mathbf{X}) \left( \int_{\Omega_0} K(\mathbf{x}(\mathbf{X}, t), \mathbf{y}(\mathbf{X}, t)) \rho_0(\mathbf{X}) d\mathbf{X} \right) d\mathbf{X} dt \end{aligned}$$

in Lagrangian coordinates. A direct computation leads to (see [24] for more detailed computations)

$$\begin{aligned} \delta \mathcal{A} = & - \int_0^T \int_{\Omega_0} \left[ -\omega' \left( \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}} \right) \cdot \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}} + \omega \left( \frac{\rho_0(\mathbf{X})}{\det \mathbf{F}} \right) \right] \times (\mathbf{F}^{-T} : \nabla_{\mathbf{X}} \delta \mathbf{x}) \det \mathbf{F} \\ & + \rho_0(\mathbf{X}) \nabla V(\mathbf{x}) \cdot \delta \mathbf{x} + \rho_0(\mathbf{X}) \int \rho_0(\mathbf{X}) \nabla K(\mathbf{x}, \mathbf{y}) \cdot \delta \mathbf{x} d\mathbf{X} \quad d\mathbf{X} dt, \end{aligned}$$

where  $\delta \mathbf{x}(\mathbf{X}, t)$  is the test function satisfying  $\tilde{\delta} \mathbf{x} \cdot \mathbf{n} = 0$  with  $\mathbf{n}$  being the unit outward normal of  $\partial \Omega$  and  $\tilde{\delta} \mathbf{x}(\mathbf{x}(\mathbf{X}, t), t) = \delta \mathbf{x}(\mathbf{X}, t)$ . We omit the tilde when there is no ambiguity. Pushing forward to the Eulerian frame, we have

$$\begin{aligned} \delta \mathcal{A} = & - \int_0^T \int_{\Omega} (-\omega'(\rho)\rho + \omega) \nabla \cdot (\delta \mathbf{x}) + (\rho \nabla V + \rho(\nabla_{\mathbf{x}} K * \rho)) \cdot \delta \mathbf{x} \quad d\mathbf{x} dt \\ = & - \int_0^T \int_{\Omega} (\nabla(\omega'(\rho)\rho - \omega) + \rho \nabla V + \rho(\nabla_{\mathbf{x}} K * \rho)) \cdot \delta \mathbf{x} \quad d\mathbf{x} dt \\ = & - \int_0^T \int_{\Omega} \rho \nabla(\omega'(\rho) + V + (K * \rho)) \cdot \delta \mathbf{x} \quad d\mathbf{x} dt. \end{aligned}$$

Hence,

$$\frac{\delta \mathcal{A}}{\delta \mathbf{x}} = -\rho \nabla \mu, \quad \mu = \frac{\delta \mathcal{F}}{\delta \rho} = \omega'(\rho) + V(\mathbf{x}) + K * \rho,$$

where  $\mu$  is called the chemical potential. Since  $\mathcal{D} = \frac{1}{2} \int \eta(\rho) |\mathbf{u}|^2 d\mathbf{x}$ , it is easy to compute that  $\frac{\delta \mathcal{D}}{\delta \mathbf{u}} = \eta(\rho) \mathbf{u}$ . So the force balance condition leads to

$$\eta(\rho) \mathbf{u} = -\rho \nabla \mu, \quad \mu = \omega'(\rho) + V(\mathbf{x}) + K * \rho.$$

#### REFERENCES

- [1] S. ADAMS, N. DIRR, M. PELETIER, AND J. ZIMMER, *Large deviations and gradient flows*, Philos. Trans. Roy. Soc. A, 371 (2013), 20120341.
- [2] B. AMOS, L. XU, AND J. Z. KOLTER, *Input convex neural networks*, in Proceedings of the International Conference on Machine Learning, PMLR, 2017, pp. 146–155.
- [3] A. BARON, *Universal approximation bounds for superposition of a sigmoid function*, IEEE Trans. Inform. Theory, 39 (1993), pp. 930–945.
- [4] R. BELLMAN, *Dynamic Programming*, Dover Publications, Mineola, NY, 1957.
- [5] J.-D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem*, Numer. Math., 84 (2000), pp. 375–393.
- [6] J.-D. BENAMOU, G. CARLIER, AND M. LABORDE, *An augmented lagrangian approach to wasserstein gradient flows and applications*, ESAIM Proc. Surveys, 54 (2016), pp. 1–17.
- [7] D. M. BLEI, A. KUCUKELBIR, AND J. D. McAULIFFE, *Variational inference: A review for statisticians*, J. Amer. Statist. Assoc., 112 (2017), pp. 859–877.
- [8] J. BRUNA, B. PEHERSTORFER, AND E. VANDEN-ELJNDEN, *Neural Galerkin schemes with active learning for high-dimensional evolution equations*, J. Comput. Phys., 496 (2024), 112588.
- [9] J. A. CARRILLO, K. CRAIG, L. WANG, AND C. WEI, *Primal dual methods for Wasserstein gradient flows*, Found. Comput. Math., 22 (2022), pp. 389–443.
- [10] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, Math. Control Signals Systems, 2 (1989), pp. 303–314.
- [11] P.-G. DE GENNES AND J. PROST, *The Physics of Liquid Crystals*, Internat. Ser. Monogr. Phys. 83, Oxford University Press, 1993.
- [12] J. DEVLIN, M.-W. CHANG, K. LEE, AND K. TOUTANOVA, *Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding*, preprint, <https://arxiv.org/abs/1810.04805>, 2018.
- [13] T. DOCKHORN, *A Discussion on Solving Partial Differential Equations Using Neural Networks*, preprint, <https://arxiv.org/abs/1904.07200>, 2019.
- [14] M. DOI, *Onsager's variational principle in soft matter*, J. Phys. Condens. Matter, 23 (2011), 284118.
- [15] Q. DU AND X. FENG, *The phase field method for geometric moving interfaces and their numerical approximations*, Handb. Numer. Anal. 21, Elsevier/North-Holland, Amsterdam, 2020, pp. 425–508.
- [16] Y. DU AND T. A. ZAKI, *Evolutional deep neural network*, Phys. Rev. E, 104 (2021), 045303.
- [17] C. DURKAN, A. BEKASOV, I. MURRAY, AND G. PAPAMAKARIOS, *Neural spline flows*, in Proceedings of the 33rd Conference on Neural Information Processing Systems, Adv. Neural Inform. Process. Syst., 32, 2019, pp. 7511–7522.
- [18] W. E, J. HAN, AND A. JENTZEN, *Algorithms for solving high dimensional pdes: From nonlinear Monte Carlo to machine learning*, Nonlinearity, 35 (2022), pp. 278–310.
- [19] W. E, C. MA, AND L. WU, *Machine learning from a continuous viewpoint*, I, Sci. China Math., 63 (2020), pp. 2233–2266.
- [20] W. E, C. MA, L. WU, AND S. WOJTOWYTSCHE, *Towards a mathematical understanding of neural network-based machine learning: What we know and what we don't*, CSIAM Trans. Appl. Math., 1 (2020), pp. 561–615.
- [21] W. E AND B. YU, *The deep ritz method: A deep learning-based numerical algorithm for solving variational problems*, Commun. Math. Stat., 6 (2018), pp. 1–12.
- [22] B. EISENBERG, Y. HYON, AND C. LIU, *Energy variational analysis of ions in water and channels: Field theory for primitive models of complex ionic fluids*, J. Chem. Phys., 133 (2010), 104104.
- [23] J. L. ERICKSEN, *Introduction to the Thermodynamics of Solids*, Appl. Math. Sci. 275, Springer, 1998.

- [24] M.-H. GIGA, A. KIRSSTEIN, AND C. LIU, *Variational modeling and complex fluids*, in Handbook of Mathematical Analysis in Mechanics of Viscous Fluids, Springer, Cham, 2018, pp. 73–113.
- [25] E. D. GIORGI, *Movimenti minimizzanti*, in Proceedings of the Conference on Aspetti e problemi della Matematica oggi, Lecce, 1992.
- [26] X. GLOROT AND Y. BENGIO, *Understanding the difficulty of training deep feedforward neural networks*, in Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [27] J. HAN, A. JENTZEN, AND W. E, *Solving high-dimensional partial differential equations using deep learning*, Proc. Natl. Acad. Sci., 115 (2018), pp. 8505–8510, <https://doi.org/10.1073/pnas.1718942115>.
- [28] B. HANIN, *Which neural net architectures give rise to exploding and vanishing gradients?*, in NIPS'18: Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2018, pp. 580–589.
- [29] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.
- [30] K. HE, X. ZHANG, S. REN, AND J. SUN, *Identity mappings in deep residual networks*, in Proceedings of the European Conference on Computer Vision, Springer, 2016, pp. 630–645.
- [31] G. HINTON, L. DENG, D. YU, G. E. DAHL, A.-R. MOHAMED, N. JAITLEY, A. SENIOR, V. VANHOUCHE, P. NGUYEN, T. N. SAINATH, AND B. KINGSBURY, *Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups*, IEEE Signal Proc. Mag., 29 (2012), pp. 82–97.
- [32] C.-W. HUANG, R. T. CHEN, C. TSIRIGOTIS, AND A. COURVILLE, *Convex Potential Flows: Universal Probability Distributions with Optimal Transport and Convex Optimization*, preprint, <https://arxiv.org/abs/2012.05942>, 2020.
- [33] H. J. HWANG, C. KIM, M. S. PARK, AND H. SON, *The Deep Minimizing Movement Scheme*, preprint, <https://arxiv.org/abs/2109.14851>, 2021.
- [34] K. JIANG AND P. ZHANG, *Numerical methods for quasicrystals*, J. Comput. Phys., 256 (2014), pp. 428–440.
- [35] R. JORDAN, D. KINDERLEHRER, AND F. OTTO, *The variational formulation of the Fokker-Planck equation*, SIAM J. Math. Anal., 29 (1998), pp. 1–17.
- [36] E. F. KELLER AND L. A. SEGEL, *Initiation of slime mold aggregation viewed as an instability*, J. Theor. Biol., 26 (1970), pp. 399–415.
- [37] E. KHARAZMI, Z. ZHANG, AND G. E. KARNIADAKIS, *hp-vpinns: Variational physics-informed neural networks with domain decomposition*, Comput. Methods Appl. Mech. Engrg., 374 (2021), 113547.
- [38] Y. KHOO, J. LU, AND L. YING, *Solving parametric PDE problems with artificial neural networks*, Eur. J. Appl. Math., 32 (2021), pp. 421–435.
- [39] D. P. KINGMA, T. SALIMANS, R. JOZEFOWICZ, X. CHEN, I. SUTSKEVER, AND M. WELLING, *Improved variational inference with inverse autoregressive flow*, in NIPS'16: Proceedings of the 30th International Conference on Neural Information Processing Systems, Adv. Neural Inf. Process. Syst. 29 2016, pp. 4743–4751.
- [40] J. F. KOLEN AND S. C. KREMER, *Gradient flow in recurrent nets: The difficulty of learning longterm dependencies*, IEEE Press, 2001, pp. 237–243, <https://doi.org/10.1109/9780470544037.ch14>.
- [41] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional neural networks*, Comm. ACM, 60 (2017), pp. 84–90.
- [42] I. E. LAGARIS, A. LIKAS, AND D. I. FOTIADIS, *Artificial neural networks for solving ordinary and partial differential equations*, IEEE Trans. Neural Networks, 9 (1998), pp. 987–1000.
- [43] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to document recognition*, Proc. IEEE, 86 (1998), pp. 2278–2324.
- [44] R. LIFSHTITZ AND H. DIAMANT, *Soft quasicrystals—why are they stable?*, Philos. Mag., 87 (2007), pp. 3021–3030.
- [45] F. H. LIN AND C. LIU, *Static and dynamic theories of liquid crystals*, J. Partial Differential Equations, 14 (2001), pp. 289–330.
- [46] S. LISINI, D. MATTHES, AND G. SAVARÉ, *Cahn–Hilliard and thin film equations with nonlinear mobility as gradient flows in weighted-Wasserstein metrics*, J. Differential Equations, 253 (2012), pp. 814–850.

- [47] G. LITJENS, T. KOOI, B. E. BEJNORDI, A. A. A. SETIO, F. CIOMPI, M. GHAFOORIAN, J. A. VAN DER LAAK, B. VAN GINNEKEN, AND C. I. SÁNCHEZ, *A survey on deep learning in medical image analysis*, Med. Image Anal., 42 (2017), pp. 60–88.
- [48] C. LIU AND H. SUN, *On energetic variational approaches in modeling the nematic liquid crystal flows*, Discrete Contin. Dyn. Syst., 23 (2009), pp. 455–475.
- [49] C. LIU, C. WANG, AND Y. WANG, *A structure-preserving, operator splitting scheme for reaction-diffusion equations with detailed balance*, J. Comput. Phys., 436 (2021), 110253.
- [50] C. LIU AND Y. WANG, *On Lagrangian schemes for porous medium type generalized diffusion equations: A discrete energetic variational approach*, J. Comput. Phys., 417 (2020), 109566.
- [51] C. LIU AND Y. WANG, *A variational Lagrangian scheme for a phase-field model: A discrete energetic variational approach*, SIAM J. Sci. Comput., 42 (2020), pp. B1541–B1569, <https://doi.org/10.1137/20M1326684>.
- [52] L. LU, X. MENG, Z. MAO, AND G. E. KARNIADAKIS, *DeepXDE: A deep learning library for solving differential equations*, SIAM Rev., 63 (2021), pp. 208–228, <https://doi.org/10.1137/19M1274067>.
- [53] Y. LU, J. LU, AND M. WANG, *A priori generalization analysis of the deep Ritz method for solving high dimensional elliptic partial differential equations*, in Proceedings of the Conference on Learning Theory, PMLR, 2021, pp. 3196–3241.
- [54] T. MÜLLER, B. MCWILLIAMS, F. ROUSSELLE, M. GROSS, AND J. NOVÁK, *Neural importance sampling*, ACM Trans. Graphics (TOG), 38 (2019), 145.
- [55] J. NOH, Y. WANG, H.-L. LIANG, V. S. R. JAMPANI, A. MAJUMDAR, AND J. P. LAGERWALL, *Dynamic tuning of the director field in liquid crystal shells using block copolymers*, Phys. Rev. Res., 2 (2020), 033160.
- [56] T. OHTA AND K. KAWASAKI, *Equilibrium morphology of block copolymer melts*, Macromolecules, 19 (1986), pp. 2621–2632.
- [57] L. ONSAGER, *Reciprocal relations in irreversible processes. I.*, Phys. Rev., 37 (1931), pp. 405–426.
- [58] L. ONSAGER, *Reciprocal relations in irreversible processes. II.*, Phys. Rev., 38 (1931), pp. 2265–2279.
- [59] G. PAPAMAKARIOS, E. T. NALISNICK, D. J. REZENDE, S. MOHAMED, AND B. LAKSHMINARAYANAN, *Normalizing flows for probabilistic modeling and inference*, J. Mach. Learn. Res., 22 (2021), 57.
- [60] M. A. PELETIER, *Variational Modelling: Energies, Gradient Flows, and Large Deviations*, preprint, <https://arxiv.org/abs/1402.1990>, 2014.
- [61] L. PRONZATO AND W. G. MÜLLER, *Design of computer experiments: Space filling and beyond*, Stat. Comput., 22 (2012), pp. 681–701.
- [62] M. RAISSI, P. PERDIKARIS, AND G. E. KARNIADAKIS, *Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations*, preprint, <https://arxiv.org/abs/1711.10561>, 2017.
- [63] L. RAYLEIGH, *Some general theorems relating to vibrations*, Proc. Lond. Math. Soc., 4 (1873), pp. 357–368.
- [64] D. REZENDE AND S. MOHAMED, *Variational inference with normalizing flows*, in Proceedings of the International Conference on Machine Learning, PMLR, 2015, pp. 1530–1538.
- [65] G. M. ROTSKOFF, A. R. MITCHELL, AND E. VANDEN-EIJNDEN, *Active Importance Sampling for Variational Objectives Dominated by Rare Events: Consequences for Optimization and Generalization*, preprint, <https://arxiv.org/abs/2008.06334>, 2020.
- [66] A. J. SALGADO AND S. M. WISE, *Classical Numerical Analysis: A Comprehensive Course*, Cambridge University Press, 2022.
- [67] Z. SHEN, H. YANG, AND S. ZHANG, *Nonlinear Approximation via Compositions*, preprint, <https://arxiv.org/abs/1902.10170>, 2019.
- [68] J. SIRIGNANO AND K. SPILIOPOULOS, *DGM: A deep learning algorithm for solving partial differential equations*, J. Comput. Phys., 375 (2018), pp. 1339–1364.
- [69] E. G. TABAK AND E. VANDEN-EIJNDEN, *Density estimation by dual ascent of the log-likelihood*, Commun. Math. Sci., 8 (2010), pp. 217–233.
- [70] B. P. VAN MILLIGEN, V. TRIBALDOS, AND J. JIMÉNEZ, *Neural network differential equation and plasma equilibrium solver*, Phys. Rev. Lett., 75 (1995), pp. 3594–3597.
- [71] J. L. VÁZQUEZ, *The Porous Medium Equation: Mathematical Theory*, Oxford University Press, 2007.
- [72] H. WANG, T. QIAN, AND X. XU, *Onsager’s variational principle in active soft matter*, Soft Matter, 17 (2021), pp. 3634–3653, <https://doi.org/10.1039/D0SM02076A>.



- [73] Y. WANG, J. CHEN, C. LIU, AND L. KANG, *Particle-based energetic variational inference*, Stat. Comput., 31 (2021), 34.
- [74] Y. WANG AND C. LIU, *Some recent advances in energetic variational approaches*, Entropy, 24 (2022), pp. 721.
- [75] Y. WANG, C. LIU, P. LIU, AND B. EISENBERG, *Field theory of reaction-diffusion: Law of mass action with an energetic variational approach*, Phys. Rev. E, 102 (2020), 062147.
- [76] Y. WANG, T.-F. ZHANG, AND C. LIU, *A two species micro-macro model of wormlike micellar solutions and its maximum entropy closure approximations: An energetic variational approach*, J. Non-Newton. Fluid Mech., 293 (2021), 104559.
- [77] Q. WEI, Y. JIANG, AND J. Z. Y. CHEN, *Machine-learning solver for modified diffusion equations*, Phys. Rev. E, 98 (2018), 053304.
- [78] E. WEINAN, *Machine Learning and Computational Mathematics*, preprint, <https://arxiv.org/abs/2009.14596>, 2020.
- [79] J. XU, Y. LI, S. WU, AND A. BOUSQUET, *On the stability and accuracy of partially and fully implicit schemes for phase field modeling*, Comput. Methods Appl. Mech. Engrg., 345 (2019), pp. 826–853.
- [80] S. XU, P. SHENG, AND C. LIU, *An energetic variational approach for ion transport*, Commun. Math. Sci., 12 (2014), pp. 779–789.
- [81] S. XU, Z. XU, O. V. KIM, R. I. LITVINOV, J. W. WEISEL, AND M. S. ALBER, *Model predictions of deformation, embolization and permeability of partially obstructive blood clots under variable shear flow*, J. R. Soc. Interface, 14 (2017).
- [82] X. XU, Y. DI, AND M., *Variational method for liquids moving on a substrate*, Phys. Fluids, 28 (2016), 087101.
- [83] P. YUE, J. J. FENG, C. LIU, AND J. SHEN, *A diffuse-interface method for simulating two-phase flows of complex fluids*, J. Fluid Mech., 515 (2004), pp. 293–317.
- [84] Y. ZANG, G. BAO, X. YE, AND H. ZHOU, *Weak adversarial networks for high-dimensional partial differential equations*, J. Comput. Phys., 411 (2020), 109409.