

A MODEL-CONSTRAINED TANGENT MANIFOLD LEARNING APPROACH FOR DYNAMICAL SYSTEMS

HAI V. NGUYEN* AND TAN BUI-THANH†

Abstract. Real-time accurate solutions of large-scale complex dynamical systems are in critical need for control, optimization, uncertainty quantification, and decision-making in practical engineering and science applications. This paper contributes in this direction a model-constrained tangent manifold learning (**mcTangent**) approach. At the heart of **mcTangent** is the synergy of several desirable strategies: i) a tangent manifold learning to take advantage of the neural network speed and the time-accurate nature of the method of lines; ii) a model-constrained approach to encode the neural network tangent with the underlying governing equations; iii) sequential learning strategies to promote long-time stability and accuracy; and iv) data randomization approach to implicitly enforce the smoothness of the neural network tangent and its likeness to the truth tangent up second order derivatives in order to further enhance the stability and accuracy of **mcTangent** solutions. Both semi-heuristic and rigorous arguments are provided to analyze and justify the proposed approach. Several numerical results for transport equation, viscous Burger’s equation, and Navier-Stokes equation are presented to study and demonstrate the capability of the proposed **mcTangent** learning approach.

Key words. dynamical systems; model-constrained learning; sequential learning; the method of lines; data randomization; tangent manifold; accuracy and stability; regularization;

1. Introduction. Dynamical systems are pervasive in engineering and science applications. They are typically time-dependent systems of ordinary differential equations (ODEs) or partial differential equations (PDEs). The latter is not different from the former from the method of lines viewpoint in which a PDE reduces to a system of ODEs after a spatial discretization. For practical settings, simulating a dynamical system could be challenging due to a large number of degrees of freedoms, and hence the number of ODEs, interdependent on each other in a highly nonlinear manner. For multi-scale or stiff systems of ODEs, explicit time discretization schemes, though straightforward, are not efficient to due time stepsize limitation to ensure stability. Implicit schemes, on the other hand, are stable but computational expensive as a large linear system of equations needs to be solved at each time step. Though currently infeasible, real-time accurate approximate solutions for the practical complex dynamical system is highly desirable for control, optimization, uncertainty quantification, and decision-making.

Towards achieving real-time solutions for dynamical systems, various pure data-driven deep learning attempts have been made. Autoencoder architecture has been explored to simulate fluid flows [18]. Autoencoder with physics-informed regularization to improve accuracy has been proposed to predict the future sea surface temperature given past series of measurements [9]. The work in [44] proposes a graph network-based model to approximate the forward map and inference model, which is then used to speed up control algorithms. As an effort to combine traditional and machine learning approaches, [30] introduces a deep Koopman model—an auto-encoder architecture of convolution neural network—to predict the dynamics of airflow over a cylinder. Comprehensive overviews of machine learning methods for forecasting dynamical systems can be found in [24] and [2]. The work in [10] presents a review and

*Department of Aerospace Engineering and Engineering Mechanics, the University of Texas at Austin, Texas (hainguyen@utexas.edu, <https://orcid.org/0000-0002-0720-7237>)

†Department of Aerospace Engineering and Engineering Mechanics, The Oden Institute for Computational Engineering and Sciences, the University of Texas at Austin, Texas (tanbui@oden.utexas.edu, <https://users.oden.utexas.edu/~tanbui/>).

aspects of using machine learning techniques to simulate turbulent flows.

Instead of replacing traditional computational approaches with pure data-driven machine learning models, which is debatable and an active research direction, one can use machine learning methods to speed up only computational demanding modules. This could maintain desirable physics constraints as in traditional approaches while gaining computational time. Indeed, a convolution neural network (CNN) can be trained to learn the numerical error between high-resolution and low-resolution simulations. Combining the CNN prediction with low-resolution simulations can then achieve similar high-resolution accuracy while being faster and at that the same time not much compromising the physics. In a different effort, [20] trains neural networks to replace components/terms severely affected by a low-resolution grid. The predictions from neural networks are unrolled over multiple time steps to improve long-time inference performance. A recurrent neural network can also be used to enhance the effectiveness of geometric multigrid methods for simulations of Navier-Stokes equations [28].

Completely replacing traditional methods while respecting governing equations, we argue, is highly desirable for machine learning methods because fast but nonphysical solutions are undesirable. A popular deep learning approach aiming to accomplish this goal is the physics-informed neural network (PINN) [41]. Similar to least squares finite element methods, PINN trains deep learning solution constrained by the PDE residual through a regularization [41, 37, 39, 40, 56, 50]). PINN can learn solutions that attempt to make the PDE residual small. However, PINN approach directly approximates the PDE solution in infinite dimensional spaces. While universal approximation results (see, e.g., [8, 12, 27, 17]) could ensure any desired accuracy with a sufficiently large number of neurons, practical network architectures are moderate in both depth and width, and hence the number of weights and biases, the accuracy of PINN can be limited. Moreover, PINN requires a retrain for new scenarios such as new boundary conditions, or new initial conditions, or new values of the underlying parameters. A physics-informed recurrent neural network has also been studied in [14]. In order to produce physically consistent and better prediction results, energy flow and density-depth constraint laws are integrated into the loss function.

Instead of learning the infinite-dimensional solution as in PINN, learning discretized solutions of dynamical systems is equally popular. The work in [58] uses a neural network to approximate the derivative of the system state in reduced projected subspace. The neural network is then combined with forward Euler and Runge-Kutta time discretization schemes to achieve high accuracy solutions. Alternatively, [32] uses a feed forward neural network directly learn the map from the solution at the current time step to the solution in the next time step. The stability and accuracy of long-time prediction are reinforced by introducing a Jacobian regularization into the loss function. Realizing several drawbacks of the direct learning approach, [53] proposes to learn the tangent manifold with Runge-Kutta schemes. Once trained, the learned tangent manifold can be used with any time discretization schemes and any time step size. In [52], the authors propose to learn a correction neural network that lifts low-resolution solutions to high-resolution accuracy. The training procedure includes the differentiable forward map. Realizing the importance of differential programming, [46] provides a differentiable package that conveniently incorporates numerical solvers with training. Alternatively, [13] develop a differentiable simulations package that wraps a numerical simulator as a gradient kernel for end-to-end back-propagation used in optimization algorithms. Similar to [46], [11] develops a differentiable physic simulations package equipped with the adjoint method for backpropagation that allows

embedding the physical forward model into the training process.

In this paper, aiming at simulating dynamical systems in real-time, we propose a model-constrained tangent manifold deep learning (mcTangent) approach that has several appealing features over existing methods. First, it operates on finite dimensional systems and thus in principle easier to train. However, it is spatial discretization-dependent for systems governed by PDEs. Second, it learns the underlying tangent manifold and thus is semi-discrete in nature. Once trained, it can be deployed with any time discretization schemes with any time step size. Third, it aims to fulfill the governing equations by constraining a fully discrete system in the loss function during training. Fourth, it is equipped with sequential learning strategies and thus promotes stability and accuracy in simulating the underlying dynamical systems far beyond the training time horizon. Fifth, our approach imposes regularizations on the smoothness of the neural network tangent and its derivatives implicitly via data randomization. This provides extra stability and accuracy for mcTangent solutions.

The paper is organized as follows. Section 2 introduces an abstract dynamical system and a model-constrained tangent manifold learning (mcTangent) approach. Both sequential machine learning and sequential model-constrained strategies will be discussed in details in Sections 2.2 and 2.3. Data randomization approach then follows with an in-depth semi-heuristic argument to reveal its implicit regularization nature in Section 2.4. In particular, data randomization induces smoothness regularization for the underlying neural network via the standard machine learning loss. The beauty of the model-constrained loss term is that it not only enforces the likeliness of the neural network and the truth tangent manifolds but also implicitly constrains their likeliness up to second-order derivatives via data randomization. Section 2.5 provides a rigorous estimation for prediction error using mcTangent approach. Several numerical results using the proposed mcTangent approach for transport equation, viscous Burger’s equation, and Navier-Stokes equation are presented in Section 3. Section 4 concludes the paper with future work.

2. Model-constrained tangent manifold deep learning solutions for dynamical systems.

2.1. Motivation. For the concreteness and simplicity of the exposition, let us consider an abstract dynamical system governed by the following time-dependent scalar PDE equation of the form

$$(2.1) \quad \frac{\partial u}{\partial t} = \mathcal{G}(u, \nabla u, \dots) \quad \text{in } \Omega \subset \mathbb{R}^d,$$

where $t \in [0, T]$, $u(\mathbf{x}) \in \mathbb{R}$ for any $\mathbf{x} \in \mathbb{R}^d$, and $d \in \{1, 2, 3\}$. We also assume (2.1) is equipped with appropriate initial conditions and boundary conditions to ensure its well-posedness.

In this paper, we are interested in *parametrized PDEs*. For downstream tasks such as design, control, optimization, inference, and uncertainty quantification, these PDEs need to be solved many times. As such, we wish to approximate solutions of (2.1) in real time for any parameters (e.g. initial conditions or boundary conditions, or some parameter). Training a PINN together with parameters (either by themselves or their neural networks weights and biases as another set of optimization variables) [7, 38, 25, 26] may not be efficient as a new solution (corresponding to new parameters) requires a retrain. We note that attempts using pure data-driven deep learning to learn the parameter-to-solution map have been explored (see, e.g., [21, 54, 34, 49, 35, 21, 48, 15]). On the other hand, standard numerical methods such as finite difference, finite

volume, and finite elements [47, 22, 16] discretize (2.1) both in time and space. One of the most popular approaches is perhaps the method of lines (see, e.g., [45]) in which one performs spatial discretization first to obtain a system of (possibly nonlinear) ordinary differential equations of the form

$$(2.2) \quad \frac{\partial \mathbf{u}}{\partial t} = \mathbf{G}(\mathbf{u}),$$

where \mathbf{u} and \mathbf{G} are vector representations of finite dimensional approximations of u and \mathcal{G} , respectively. Now, either an explicit or implicit (or their combination) can be deployed to discretize the temporal derivative. For the former, the most expensive operation is the evaluation of tangent manifold $\mathbf{G}(\mathbf{u})$. For the latter, evaluating both $\mathbf{G}(\mathbf{u})$ and its (possibly approximate) Jacobian for each time step play a vital role. Implementing the Jacobian, even with the adjoint method [51], is a significant part of the programming effort. Automatic differentiation can mitigate this programming burden at the expense of more memory bandwidth. In summary, computing $\mathbf{G}(\mathbf{u})$ and its Jacobian is a major part, both in implementation and computational time, of existing numerical methods.

To overcome the time burden of estimating the tangent manifold and its Jacobian, we present a model-constrained tangent manifold deep learning approach (**mcTangent**) inspired by the semi-discrete nature of the method of lines. In particular, we first learn the tangent manifold $\mathbf{G}(\mathbf{u})$ using neural network and then use a time discretization to solve for approximations of \mathbf{u} . Our approach thus aims to approximate only the spatial discretization and leaves the temporal discretization for traditional time integrators. At the heart of our approach is the incorporation of the governing equations into the neural network tangent by constraining the learning task to respect a temporal discretization of (2.2). Again, unlike PINN and its siblings which learn the infinite-dimensional solution u , our approach learns the tangent manifold of the finite-dimensional approximation \mathbf{u} . Furthermore, we constrain the physics on the discrete level. Clearly, our approach is discretization-dependent while PINN requires neither spatial discretization nor temporal discretization.

2.2. Model-constrained neural network approach. In this section, we construct a model-constrained neural network $\Psi(\mathbf{u})$ to learn $\mathbf{G}(\mathbf{u})$. This is done in tandem with a time discretization of (2.2). For clarity, we limit our presentation to forward Euler method

$$(2.3) \quad \mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \mathbf{G}(\mathbf{u}^k),$$

as it is straightforward to extend the approach to any time discretization scheme. The task at hand is to train $\Psi(\mathbf{u})$ on a certain spatial mesh \mathcal{T} corresponding to a spatial discretization. To begin, let us denote the numerical solutions of (2.3) at $N_t + 1$ time steps on a finer mesh \mathcal{T}^f as

$$\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^{N_t}\}.$$

which are then down-sampled on \mathcal{T} for training $\Psi(\mathbf{u})$. Doing so has proved to yield more accurate predictions than training directly on the solutions on \mathcal{T} [33, 20, 57]. This is not surprising as the down-sampled training data on \mathcal{T} is more accurate than the solution on \mathcal{T} .

The next idea that we like to incorporate into our approach is sequential training. The key is to feed the predictions back to the neural network model to enable a better long-time predictive capability. Using this idea [55] deploys a mixture of graph

neural network and 3D-U-Net neural network to model fluid flows. Similarly, in [57] sequential learning is used to train a network to obtain the optimal finite difference coefficients from the high-resolution training data. In the context of atmosphere modeling, [5] introduces a stable and highly accurate long-time prediction loss function with sequential training. Following [33, 20, 55, 57], we partition the training data in $N_t - S$ overlapping subsets

$$\mathcal{U} := \{(\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^{S+1}), (\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^{S+2}), \dots, (\mathbf{u}^{N_t-S-1}, \mathbf{u}^{N_t-S}, \dots, \mathbf{u}^{N_t})\}.$$

For convenience in the exposition, we enumerate these $N_t - S$ subsets as

$$\mathcal{U} := \{(\mathbf{u}^{0,0}, \mathbf{u}^{0,1}, \dots, \mathbf{u}^{0,S+1}), (\mathbf{u}^{1,0}, \mathbf{u}^{1,1}, \dots, \mathbf{u}^{1,S+1}), \dots, (\mathbf{u}^{N_t-S-1,0}, \mathbf{u}^{N_t-S-1,1}, \dots, \mathbf{u}^{N_t-S-1,S+1})\},$$

where the second superscript denotes the local index in each subset. To distinguish from the sequential model-constrained learning in Section 2.3, let us call the machine learning approach based on these overlapping subsets as *sequential machine learning*.

We next discuss how we use each subset in our model-constrained approach. Consider the k th subset $(\mathbf{u}^{k,0}, \mathbf{u}^{k,1}, \dots, \mathbf{u}^{k,S+1})$, for $k = 0, \dots, N_t - S - 1$. Starting from $\tilde{\mathbf{u}}^{k,0} = \mathbf{u}^{k,0}$, we can write the sequence of approximate solutions $\{\tilde{\mathbf{u}}^{k,i}\}_{i=1}^{S+1}$ for (2.2) using forward Euler time discretization with the neural network tangent $\Psi(\mathbf{u})$ as

$$(2.4) \quad \tilde{\mathbf{u}}^{k,i+1} = \tilde{\mathbf{u}}^{k,i} + \Delta t \Psi(\tilde{\mathbf{u}}^{k,i}), \quad i = 0, \dots, S.$$

On the other hand, if we feed $\tilde{\mathbf{u}}^{k,i}$ through the forward Euler discretization (2.3) we obtain

$$(2.5) \quad \bar{\mathbf{u}}^{k,i+1} = \tilde{\mathbf{u}}^{k,i} + \Delta t \mathbf{G}(\tilde{\mathbf{u}}^{k,i}), \quad i = 0, \dots, S.$$

As can be seen $\tilde{\mathbf{u}}^{k,i+1} \neq \bar{\mathbf{u}}^{k,i+1}$, though we wish they are the same. *If they were, the approximate solutions using neural network tangent would respect the governing discretized equation exactly.* Obviously, this is not feasible in general. Thus, we resort to requiring $\tilde{\mathbf{u}}^{k,i+1}$ as close as possible to $\bar{\mathbf{u}}^{k,i+1}$. One way to accomplish this is to consider the following loss function for the k th batch:

$$(2.6) \quad \mathcal{J}_k := \frac{1}{S+1} \sum_{i=1}^{S+1} \left(\|\mathbf{u}^{k,i} - \tilde{\mathbf{u}}^{k,i}\|_2^2 + \alpha \|\bar{\mathbf{u}}^{k,i} - \tilde{\mathbf{u}}^{k,i}\|_2^2 \right),$$

where α is a model-constrained penalty (or regularization) parameter. The first term of the loss (2.6)—the ML Loss in Figure 1—ensures the data consistency, while the second term—the MC Loss in Figure 1—is to force approximate solutions of (2.2) using neural network tangent $\Psi(\mathbf{u})$ to best fit the underlying space-time discretization (2.3). The schematic of the **mcTangent** architecture for the k th data subset and $S = 1$ is illustrated in Figure 1.

Remark 2.1. Note that it is not essential that $\bar{\mathbf{u}}^i$ must be obtained by the forward Euler scheme (2.3). In fact, our approach is flexible in the sense that any one-step explicit scheme, denoted as \mathcal{F} , (including explicit Runge-Kutta) is admissible. In such a case, our neural network can be considered as learning the forward Euler approximation of the ground-truth scheme.

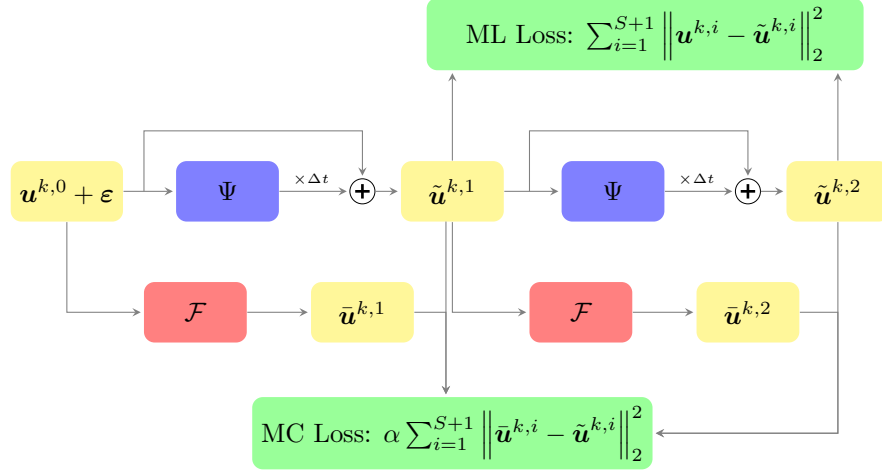


Fig. 1: The schematic of **mcTangent** network architecture for $S = 1$. For the data randomization approach in Section 2.4, the random noise vector, $\boldsymbol{\varepsilon}$, is added to the first input of the neural network.

Taking all the batches into account yields the total loss function

$$(2.7) \quad \mathcal{J} := \frac{1}{(N_t - S)(S + 1)} \sum_{k=0}^{N_t - S - 1} \sum_{i=1}^{S+1} \left(\|\mathbf{u}^{k,i} - \tilde{\mathbf{u}}^{k,i}\|_2^2 + \alpha \|\bar{\mathbf{u}}^{k,i} - \tilde{\mathbf{u}}^{k,i}\|_2^2 \right).$$

To gain insight into our **mcTangent** approach, we consider a linear problem in which $\mathbf{G}(\mathbf{u}) = \mathbf{G}\mathbf{u}$, and a one-layer linear neural network $\Psi(\mathbf{u}^{k,0}) = \mathbf{W}\mathbf{u}^{k,0} + \mathbf{b}$. Under a mild condition, our approach should exactly recover the underlying tangent manifold, i.e. $\Psi(\mathbf{u}) = \mathbf{G}\mathbf{u}$. Indeed, let $S = 0$ so that the loss function (2.7) becomes

$$(2.8) \quad \begin{aligned} \mathcal{J} &= \frac{1 + \alpha}{N_t} \sum_{k=0}^{N_t - 1} \|\mathbf{u}^{k,1} - \tilde{\mathbf{u}}^{k,1}\|_2^2 = \frac{1 + \alpha}{N_t} \|U^1 - \tilde{U}^1\|_F^2 \\ &= \frac{(1 + \alpha) \Delta t^2}{N_t} \|GU^0 - (\mathbf{W}U^0 + \mathbf{b}\mathbf{1}^T)\|_F^2 \end{aligned}$$

where U^{t_i} and \tilde{U}^{t_i} are matrices with true and predictive solutions as columns, respectively, and $\mathbf{1}$ is the unit column vector.

LEMMA 2.2. *The optimal solution $(\mathbf{W}^*, \mathbf{b}^*)$ for the training problem*

$$\min_{\mathbf{W}, \mathbf{b}} \mathcal{J}$$

is given by

$$(2.9) \quad \mathbf{W}^* = G\bar{U}\bar{U}^\dagger, \quad \mathbf{b}^* = G(\mathbf{I} - \bar{U}\bar{U}^\dagger)\bar{\mathbf{u}},$$

where $\bar{\mathbf{u}} := \frac{1}{N_t} U^{t_0} \mathbf{1}$ is the column-average of matrix U^{t_0} , $\bar{U} := U^{t_0} - \bar{\mathbf{u}}\mathbf{1}^T$, and † denotes the pseudo-inverse. Consequently, the optimal network reads

$$\Psi(\mathbf{u}) = G\bar{U}\bar{U}^\dagger \mathbf{u} + G(\mathbf{I} - \bar{U}\bar{U}^\dagger)\bar{\mathbf{u}}.$$

Remark 2.3. Lemma 2.2 tells us that the optimal network exactly recovers the true forward map G if \bar{U} is a full row rank matrix. (In that case, $\bar{U}\bar{U}^\dagger = \mathbf{I}$.) This holds, for example, when the number of independent data samples is equal to the discretized dimension. We would like to point out that the MC loss term is the same as the ML loss term (up to a constant), and thus does not provide any extra information in this simple case. When $S > 0$, at the time of writing this paper, we are not able to find a closed form solution as in Lemma 2.2. We leave it as future work.

2.3. Sequential model-constrained neural network strategy. Section 2.2 presents a sequential machine learning approach for the proposed model-constrained neural network $\Psi(\mathbf{u})$ to learn the tangent manifold while being constrained to provide the best possible approximate solutions for (2.3) for each time step. In order to improve the long-time predictive capability and accuracy, this section constructs, in addition to sequential machine learning, a sequential model-constrained learning strategy for training the neural network $\Psi(\mathbf{u})$. The sequential model-constrained learning is designed to promote the neural network solutions to respect the underlying discretization scheme for multiple time steps concurrently. In particular, starting from $\tilde{\mathbf{u}}^{k,i}$ we can carry out R steps forward in time using the underlying discretization (2.3) as

$$\bar{\mathbf{u}}^{k,i,r+1} = \bar{\mathbf{u}}^{k,i,r} + \Delta t \mathbf{G}(\bar{\mathbf{u}}^{k,i,r}), \quad r = 1, \dots, R,$$

and using the neural network approximation (2.4) as

$$\tilde{\mathbf{u}}^{k,i,r+1} = \tilde{\mathbf{u}}^{k,i,r} + \Delta t \Psi(\tilde{\mathbf{u}}^{k,i,r}), \quad r = 1, \dots, R,$$

where $\bar{\mathbf{u}}^{k,i,1} = \tilde{\mathbf{u}}^{k,i,1} = \tilde{\mathbf{u}}^{k,i}$. Here the third superscript r has been introduced to keep track of R sequential forward steps starting from $\tilde{\mathbf{u}}^{k,i}$ for both exact and neural network tangent manifolds. In order to ensure that these corresponding R sequential predictions closely match each other, we consider the following loss function

$$(2.10) \quad \mathcal{J} := \frac{1}{(N_t - S)(S + 1)} \sum_{k=0}^{N_t - S - 1} \sum_{i=1}^{S+1} \left(\left\| \mathbf{u}^{k,i} - \tilde{\mathbf{u}}^{k,i} \right\|_2^2 + \frac{\alpha}{R} \sum_{r=1}^R \left\| \bar{\mathbf{u}}^{k,i,r} - \tilde{\mathbf{u}}^{k,i,r} \right\|_2^2 \right).$$

Clearly, when $R = 1$ we recover (2.7) from (2.10). In other words, (2.10) is a generalization of (2.7). We would like to point out that it is computationally expensive using large values for both S and R . In the numerical results in Section 3, we study two combinations: $S \geq 1, R = 1$ and $S = 1, R \geq 1$. In order to have a deeper understanding of the role of the loss function (2.10) in training the neural network tangent and its predictive capability, we shall provide an in-depth heuristic argument in Section 2.4 and a rigorous error estimation for `mcTangent` predictions in Section 2.5.

2.4. Data randomization. It has been observed [43] that adding a small amount of noise to training data not only increases the generalization on unseen data but also reduces accumulated errors in predictions. In fact, clean noise-free data does not represent the accumulated error in the predictive state that is fed back to the network to produce subsequent predictions. Moreover, noisy data encourages neural network predictions to be more robust to noise-corrupted inputs and errors. In order to investigate the significance of different noise additions (adding noise to the training inputs, weights of the neural network, and output labels) on the model generalization, [1] demonstrates that the reasonable noise level in the outputs does not influence the trained network. Randomizing training data, on the one hand, prevents the neural

network from overfitting data, and on the other hand, can make the network insensitive to noise in data in the validation phase.

It is well-known that randomization induces a regularization of the gradient of the loss function with respect to the inputs [42]. Consequently, the neural network, if a proper noise level is used, is regularized to be a smooth function of the input data. The smoothness reduces the sensitivity to the variation in the input [29] and can enhance the stability of long-time predictions [36]. The work in [3] showed that adding noise to data is equivalent to introducing a Tikhonov regularization to the loss function (where the regularization parameter is the noise variance) and thus improving the model generalization. However, the analysis is only valid in the context of infinite training data set, as pointed out in [1].

Inspired by the aforementioned work, we randomize the input data for the model-constrained network. We shall heuristically show that randomization induces regularizations not only to promote the smoothness of the network but also to enhance the similarity of the derivatives of the network $\Psi(\mathbf{u})$ and the true tangent manifold $\mathbf{G}(\mathbf{u})$. As shall be seen, the numerical results in Section 3 reveal that randomization improves significantly the long-term stability and accuracy.

In this paper, we randomize the input \mathbf{u} of the network as

$$(2.11) \quad \mathbf{v} = \mathbf{u} + \boldsymbol{\varepsilon},$$

where $\boldsymbol{\varepsilon}$ is a normal random vector $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \delta^2 \mathbf{I})$. Note that the following heuristic arguments also hold for any random vector with independent components, each of which is a random variable with zero mean and variances δ^2 . Let $\mathbb{E}[\cdot]$ denote the expectation with respect to $\boldsymbol{\varepsilon}$. Following [1], for a generic loss function $\mathcal{L}(\mathbf{u})$ we have

$$(2.12) \quad \begin{aligned} \mathbb{E}[\mathcal{L}(\mathbf{v})] &= \mathcal{L}(\mathbf{u}) + \mathbb{E} \left[\frac{\partial \mathcal{L}}{\partial \mathbf{u}} \Big|_{\mathbf{u}} \boldsymbol{\varepsilon} \right] + \frac{1}{2} \mathbb{E} \left[\boldsymbol{\varepsilon}^T \frac{\partial^2 \mathcal{L}}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}} \boldsymbol{\varepsilon} \right] + o(\|\boldsymbol{\varepsilon}\|^2) \\ &\approx \mathcal{L}(\mathbf{u}) + \frac{1}{2} \mathbb{E} \left[\boldsymbol{\varepsilon}^T \frac{\partial^2 \mathcal{L}}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}} \boldsymbol{\varepsilon} \right], \end{aligned}$$

where we have used sufficient small noise $\boldsymbol{\varepsilon}$ (relatively to \mathbf{u}) so that the high-order term $o(\|\boldsymbol{\varepsilon}\|^2)$, using the standard ‘‘small o’’ notation, is negligible and thus can be ignored. We consider $S = 0$ and $R = 1$. (For $S > 0$ and/or $R > 1$, the sequential inputs to the network contain the error which may not satisfy the condition for (2.12) to hold.) In this case the loss function (2.7) becomes

$$(2.13) \quad \mathcal{J} = \frac{1}{N_t} \sum_{k=0}^{N_t-1} \underbrace{\left\| \mathbf{u}^{k,1} - \tilde{\mathbf{u}}^{k,1} \right\|_2^2}_{\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon})} + \alpha \underbrace{\left\| \bar{\mathbf{u}}^{k,1} - \tilde{\mathbf{u}}^{k,1} \right\|_2^2}_{\mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon})}$$

We now study the randomized ML loss term $\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon})$ and the randomized MC loss term $\mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon})$ to gain insights into the role of randomization.

The machine learning loss term reads

$$\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon}) = \left\| \mathbf{u}^{k,1} - (\mathbf{u}^{k,0} + \boldsymbol{\varepsilon} + \Delta t \Psi(\mathbf{u}^{k,0} + \boldsymbol{\varepsilon})) \right\|_2^2$$

which is a function of true input $\mathbf{u}^{k,0}$ plus a random noise vector $\boldsymbol{\varepsilon}$. *It is important to note that we do not randomize the true data $\mathbf{u}^{k,1}$ against which we compare the*

machine prediction $\tilde{\mathbf{u}}^{k,1}$. Replacing \mathcal{L} by \mathcal{L}_{ML} in (2.12) yields

$$(2.14) \quad \mathbb{E} [\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0} + \varepsilon)] \approx \underbrace{\|\mathbf{u}^{k,1} - (\mathbf{u}^{k,0} + \Delta t \Psi(\mathbf{u}^{k,0}))\|_2^2}_{\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0})} + \delta^2 [\mathcal{P}_1(\mathbf{u}^{k,0}) + \Delta t \mathcal{P}_2(\mathbf{u}^{k,0})],$$

where

$$(2.15) \quad \mathcal{P}_1(\mathbf{u}^{k,0}) = \mathbf{Tr} \left[\left(\mathbf{I} + \Delta t \frac{\partial \Psi}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} \right)^T \left(\mathbf{I} + \Delta t \frac{\partial \Psi}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} \right) \right],$$

with $\mathbf{Tr}[\cdot]$ as the trace operator, and

$$(2.16) \quad \mathcal{P}_2(\mathbf{u}^{k,0}) = \mathbf{Tr} \left[\Delta t \frac{\partial^2 \Psi}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}^{k,0}} \odot [(\mathbf{u}^{k,0} + \Delta t \Psi(\mathbf{u}^{k,0})) - \mathbf{u}^{k,1}] \right],$$

where \odot denotes the dot product of the third order tensor $\Delta t \frac{\partial^2 \Psi}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}^{k,0}}$ and the vector $[(\mathbf{u}^{k,0} + \Delta t \Psi(\mathbf{u}^{k,0})) - \mathbf{u}^{k,1}]$.

From (2.14), three observations are in order. First, on average, the randomized ML loss term is approximately the original ML loss term plus two additional terms \mathcal{P}_1 and \mathcal{P}_2 scaled by the variance δ^2 of the noise. Second, the first term \mathcal{P}_1 is positive and thus is a regularization. It enforces the boundedness of the gradient (and hence the smoothness) of the neural network. Third, the second term \mathcal{P}_2 can be either positive or negative. However, when time step Δt is small and/or the ML misfit term $[(\mathbf{u}^{k,0} + \Delta t \Psi(\mathbf{u}^{k,0})) - \mathbf{u}^{k,1}]$ is small (e.g. with sufficient training), the contribution of the second term is expected to be dominated by the first, and thus is negligible. When neither of these two conditions is satisfied, if the training enforces small ‘‘curvature’’ of the neural network (i.e. small $\frac{\partial^2 \Psi}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}^{k,0}}$) then the second term is also negligible. When this happens, training with randomization provides extra smoothness to the network.

Next, from (2.4) and (2.5), the randomized MC loss term can be written as

$$\mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0} + \varepsilon) = \|\bar{\mathbf{u}}^{k,1} - \tilde{\mathbf{u}}^{k,1}\|_2^2 = \Delta t^2 \|G(\mathbf{u}^{k,0} + \varepsilon) - \Psi(\mathbf{u}^{k,0} + \varepsilon)\|_2^2.$$

Applying (2.12) with \mathcal{L}_{MC} in place of \mathcal{L} gives

$$(2.17) \quad \mathbb{E} [\mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0} + \varepsilon)] \approx \underbrace{\Delta t^2 \|G(\mathbf{u}^{k,0}) - \Psi(\mathbf{u}^{k,0})\|_2^2}_{\mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0})} + \delta^2 [\mathcal{Q}_1(\mathbf{u}^{k,0}) + \mathcal{Q}_2(\mathbf{u}^{k,0})],$$

where

$$(2.18) \quad \mathcal{Q}_1(\mathbf{u}^{k,0}) = \Delta t^2 \mathbf{Tr} \left[\left(\frac{\partial G}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} - \frac{\partial \Psi}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} \right)^T \left(\frac{\partial G}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} - \frac{\partial \Psi}{\partial \mathbf{u}} \Big|_{\mathbf{u}^{k,0}} \right) \right],$$

and

$$(2.19) \quad \mathcal{Q}_2(\mathbf{u}^{k,0}) = \Delta t \mathbf{Tr} \left[\left(\frac{\partial^2 G}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}^{k,0}} - \frac{\partial^2 \Psi}{\partial \mathbf{u}^2} \Big|_{\mathbf{u}^{k,0}} \right) \odot \Delta t (\Psi(\mathbf{u}^{k,0}) - G(\mathbf{u}^{k,0})) \right].$$

As can be seen, the randomized MC loss term is approximately a sum of the original ML loss term and two additional terms. The first term \mathcal{Q}_1 is non-negative and behaves

like a regularization to enforce the likeliness of the derivatives with respect to \mathbf{u} of the neural network tangent $\Psi(\mathbf{u})$ and the true tangent $\mathbf{G}(\mathbf{u})$. The second term, though could be either negative or positive, can be negligible with sufficient training so that the MC misfit $\Delta t (\Psi(\mathbf{u}^{k,0}) - \mathbf{G}(\mathbf{u}^{k,0}))$ is relatively small. Another possibility for the insignificance of the second term is when the difference in the ‘‘curvature’’ of the neural network tangent and the true tangent is sufficiently small. In that case, training with randomization promotes the closeness of not only $\Psi(\mathbf{u})$ and $\mathbf{G}(\mathbf{u})$ but their first and second derivatives with respect to \mathbf{u} : *confirming the significant advantages obtained from data randomization.*

Next, combining (2.13), (2.14), and (2.17) yields

$$(2.20) \quad \mathbb{E}[\mathcal{J}] \approx \frac{1}{N_t} \sum_{k=0}^{N_t-1} (\mathcal{L}_{\text{ML}}(\mathbf{u}^{k,0}) + \alpha \mathcal{L}_{\text{MC}}(\mathbf{u}^{k,0})) \\ + \frac{\delta^2}{N_t} \sum_{k=0}^{N_t-1} [\mathcal{P}_1(\mathbf{u}^{k,0}) + \Delta t \mathcal{P}_2(\mathbf{u}^{k,0}) + \alpha (\mathcal{Q}_1(\mathbf{u}^{k,0}) + \mathcal{Q}_2(\mathbf{u}^{k,0}))],$$

in which the first sum is the original loss (without randomization) and the second sum consists of additional terms induced by data randomization. These additional terms play a vital role in stimulating the stability and accuracy of the neural network. Indeed, as discussed above, randomizing the machine learning loss term encourages the smoothness of the neural network tangent by penalizing its first and second derivatives implicitly. Note that explicitly penalizing the first derivative of a neural network as in [32] is possible, but this could be computationally expensive and challenging. Doing so for both the first and second derivatives is not recommended. The above heuristic analysis of data randomization also reveals the *power of the model-constrained term* in training neural network: it promotes the agreement of the neural network tangent and the true tangent up to second order that is otherwise not realizable using the standard data-driven approach with only machine learning loss term.

2.5. Estimation of prediction errors. In this section, we show how data randomization helps improve the stability and accuracy of long-time predictions. We are interested in predicting solutions of the system (2.2) starting from an initial condition \mathbf{u}^0 that is not in the training set. To that end, it is natural to compare the `mcTangent` solutions $\tilde{\mathbf{u}}^i$ in (2.4) with the solutions \mathbf{u}^i obtained from the discretized system (2.3). Let us define the neural prediction error as

$$(2.21) \quad \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^i) = \mathbf{u}^{i+1} - [\tilde{\mathbf{u}}^i + \Delta t \Psi(\tilde{\mathbf{u}}^i)], \quad \varepsilon^{i+1} = \|\mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^i)\|_2.$$

From (2.3), (2.4), and (2.21) we have

$$(2.22) \quad \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^i) = (\mathbf{u}^i + \Delta t \mathbf{G}(\mathbf{u}^i)) - (\tilde{\mathbf{u}}^i + \Delta t \Psi(\tilde{\mathbf{u}}^i)) \\ = \Delta t \mathbf{G}(\tilde{\mathbf{u}}^i + \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^{i-1})) - \Delta t \Psi(\tilde{\mathbf{u}}^i) + \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^{i-1})$$

Applying the Taylor expansion for the first term gives

$$(2.23) \quad \Delta t \mathbf{G}(\tilde{\mathbf{u}}^i + \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^{i-1})) = \Delta t \mathbf{G}(\tilde{\mathbf{u}}^i) + \Delta t \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \mathbf{e}_{\text{ML}}(\tilde{\mathbf{u}}^{i-1}) + o(\varepsilon^i)$$

Substituting back to (2.22), we have

$$\begin{aligned}
(2.24) \quad & e_{\text{ML}}(\tilde{\mathbf{u}}^i) = \Delta t [\mathbf{G}(\tilde{\mathbf{u}}^i) - \Psi(\tilde{\mathbf{u}}^i)] + \Delta t \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} e_{\text{ML}}(\tilde{\mathbf{u}}^{i-1}) + e_{\text{ML}}(\tilde{\mathbf{u}}^{i-1}) \\
& + o(\varepsilon^i) = \Delta t [\mathbf{G}(\tilde{\mathbf{u}}^i) - \Psi(\tilde{\mathbf{u}}^i)] + \Delta t \left[\left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} - \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right] e_{\text{ML}}(\tilde{\mathbf{u}}^{i-1}) \\
& + \left[\mathbf{I} + \Delta t \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right] e_{\text{ML}}(\tilde{\mathbf{u}}^{i-1}) + o(\varepsilon^i)
\end{aligned}$$

Applying triangle inequality and Cauchy–Schwarz inequality for (2.24) and using (2.21) yields

$$\begin{aligned}
(2.25) \quad & \varepsilon^{i+1} \leq \Delta t \|\mathbf{G}(\tilde{\mathbf{u}}^i) - \Psi(\tilde{\mathbf{u}}^i)\|_2 \\
& + \Delta t \left\| \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} - \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2 \varepsilon^i + \left\| 1 + \Delta t \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2 \varepsilon^i + o(\varepsilon^i), \quad i \geq 0.
\end{aligned}$$

We observe in (2.25) that the first term on the right-hand side is the model-constrained loss term being as small as possible at the training data. On the other hand, $\Delta t \left\| \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} - \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2$ and $\left\| 1 + \Delta t \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2$ are regularized to be bounded and/or small by data randomization (see Section 2.4). A heuristic argument reveals that the prediction error is under control at all times. Indeed, suppose $\Delta t \|\mathbf{G}(\tilde{\mathbf{u}}^i) - \Psi(\tilde{\mathbf{u}}^i)\|_2$, $\Delta t \left\| \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} - \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2$, and $\left\| 1 + \Delta t \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2$ are bounded. Since $\varepsilon^0 = 0$, ε^1 is bounded, and by induction ε^i is also bounded for $i \geq 0$. A rigorous version of this argument is given in Theorem 2.4.

THEOREM 2.4. *Assume that the second derivative of $\mathbf{G}(\mathbf{u})$ with respect to \mathbf{u} is uniformly bounded. Let*

$$f^{i+1} := \Delta t \|\mathbf{G}(\tilde{\mathbf{u}}^i) - \Psi(\tilde{\mathbf{u}}^i)\|_2,$$

and

$$g^{i+1} := \Delta t \left\| \left. \frac{\partial \mathbf{G}}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} - \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2 + \left\| 1 + \Delta t \left. \frac{\partial \Psi}{\partial \mathbf{u}} \right|_{\tilde{\mathbf{u}}^i} \right\|_2 + c^i,$$

where $c^i = \mathcal{O}(\varepsilon^i)$. Then, the prediction error ε^n at time t_n satisfies

$$\varepsilon^n \leq \sum_{k=1}^n (\prod_{i=k+1}^n g^i) f^k.$$

Proof. The proof is a simple application of a discrete Gronwall lemma on (2.25).□

Remark 2.5. Note that the boundedness of the second derivative of $\mathbf{G}(\mathbf{u})$ with respect to \mathbf{u} is valid for problem (2.2) with a smooth tangent manifold. The boundedness of f^i and g^i is not too restricted if the prediction scenarios are close to the training data. Indeed, as argued in Section 2.4, data randomization enforces the small values for f^i and g^i at the training points. Now, due to the smoothness of $\Psi(\mathbf{u})$ and $\mathbf{G}(\mathbf{u})$ and their closeness in both values and derivatives (again by randomization), the continuity guarantees the small values for f^i and g^i during the prediction.

Remark 2.6. Theorem 2.4 allows us to bound the error between the neural network prediction with the exact solution of the original PDEs (2.1) provided that an error estimation of the solution of the discretized equation (2.3) is given. Indeed,

suppose the error in the discretized solution \mathbf{u}^n and the exact solution $\mathbf{u}(t_n)$ at time t_n is bounded by $\mathcal{O}(\Delta t + h^p)$, where h is the mesh size and p is the order of accuracy of the underlying spatial discretization. Then by a simple application of triangle inequality we have

$$\tilde{\mathbf{u}}^n - \mathbf{u}(t_n) = \mathcal{O}\left(\Delta t + h^p + \sum_{k=1}^n (\Pi_{i=k+1}^n g^i) f^k\right),$$

which shows that in order to get the optimal accuracy and computational effort we ideally need to balance not only the temporal and spatial discretization errors but also the error in the neural network. Clearly, balancing the former two is not that difficult from a numerical analysis point of view, but balancing also the network error is challenging as it depends on the actual training process and randomization.

3. Numerical results. In this section, we present several numerical results using the proposed model-constrained tangent manifold neural network (**mcTangent**) approach for transport equation, viscous Burger’s equation, and Navier-Stokes equation. As shall be shown, **mcTangent** solutions are—thanks to the model-constrained term and data randomization—stable and capable of producing accurate approximations far beyond the training time horizons.

Five hyperparameters of interest are the number of training samples, noise level δ , sequential machine learning steps S , sequential model-constrained learning steps R , and regularization parameter α . For convenience, we shall conventionally write them in a group. For example the $(d600, 2\%, 1, 1, 0)$ setting means we consider 600 training data samples, 2% noise, $S = 1$, $R = 1$, and $\alpha = 0$. In order to ensure the fairness between simulations and the comparison among approaches, we use fixed random keys for train and test data generation, for adding noise, and for neural network parameter initialization. We implement our approach and perform all computations in JAX [4]. We would like to point out that all computations (training, testing, and predicting) are done with single precision accuracy.

3.1. One-dimensional (1D) wave/transport equation. The 1D wave equation considered in this section is given by

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

with the wave speed $c = 1$, the spatial domain $x \in (0, 1)$, and time horizon $t \in (0, T)$. The equation is equipped with an initial condition $u(x, 0) = u_0(x)$ and periodic boundary condition. We are interested in real-time approximate solutions of the wave equation for any initial condition $u_0(x)$.

Data generation. In this problem, the initial condition samples are drawn from

$$u_0(x) = \sum_{i=1}^5 a_i \sin(2\pi x i) + \sum_{i=1}^5 b_i \cos(2\pi x i),$$

where a_i, b_i are distributed by the standard normal distribution with zero mean and unit variance, i.e., $a_i, b_i \sim \mathcal{N}(0, 1)$. We solve the wave equation with the forward Euler scheme for the temporal derivative and the first order upwind finite difference scheme for spatial derivative. The time horizon is chosen as $T = 5 \times 10^{-2}$. A fine space-time mesh with $n_x = 10000$ points in space and $n_t = 2000$ points in time is deployed to achieve highly accurate solutions. The training data samples are obtained

by extracting the high resolution solutions on a coarser uniform space-time mesh with $n_x = 200$ and $n_t = 100$. In this simple problem, we generate a fixed training data set of 200 initial conditions. Note that we aim to predict long-time solutions, $t \in (0, 3)$, from the short-time training data in the interval $t \in (0, 5 \times 10^{-2})$.

Neural network architecture. Because of the linear nature of the problem and the first order upwind finite difference scheme, a linear neural work is sufficient to approximate the resulting tangent manifold. The linear neural network is defined as

$$\Psi(\mathbf{u}^i) = W\mathbf{u}^i + b,$$

where the weights $W \in \mathcal{R}^{n_x \times n_x}$, and the bias $b \in \mathcal{R}^{n_x}$. To train, we use ADAM [19] optimizer with default parameters and the learning rate of 10^{-3} . We determine the best combination of weights and biases (and hence the final trained network) as the one that provides the lowest accumulated mean square error for 500 time steps for the test sample. Specifically, during the training process, at each epoch, we solve for the predictions from the test initial condition with the current-epoch learned network. The accumulated mean-square error between predictive solutions and ground truth solutions is calculated at the 500th time step to determine the “optimal” network.

Long-time predictions. Shown in Figure 2 is the mean-square error between true (high resolution) solutions and predicted ones obtained by various neural networks, each of which is trained with both randomized and noise-free training data.

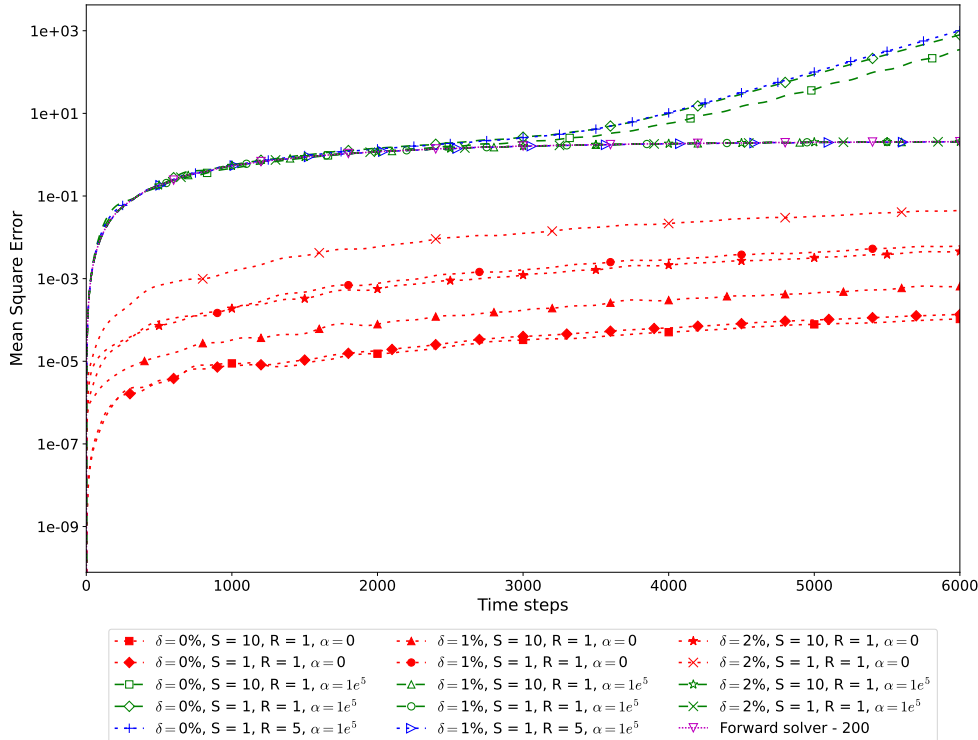


Fig. 2: **Wave/transport equation.** Comparison between different neural network approaches with/without randomization.

For pure data-driven machine learning networks ($\alpha = 0$, and thus no model-constrained term), we observe that noise-free data trained networks outperform those trained with noisy data ones. This is not surprising as for this linear problem, as predicted by Lemma 2.2, one can obtain linear networks that accurately learn the tangent manifold with sufficiently rich data. Therefore, the predictions by the learned linear networks are almost the same as the ground truth solutions. On the contrary, training with noisy data causes the neural network to predict solutions with a small amount of error such that it adapts to (possibly overfits) the amount of noise in the ground-truth solutions. Figure 3 presents the weight matrix and bias vector for two cases $(d200, 0\%, 10, 1, 0)$ and $(d200, 0\%, 1, 1, 0)$ with noise-free data. It can be seen that both networks are almost identical and both have only an upper diagonal with a large magnitude. We also note that the bias vector is relatively small and thus we ignore this bias vector in the subsequent comparisons. We present the test predicted solutions for the setting $(d200, 0\%, 1, 1, 0)$ in Figure 4. As the network fits the tangent manifold for high-resolution data, accurate results are preserved far beyond the training time horizon, while finite difference results on the same coarse grid show a severe diffusion/dissipation effect. Furthermore, settings with large number of sequential steps such as $(d200, 0\%, 10, 1, 0)$, $(d200, 1\%, 10, 1, 0)$ and $(d200, 2\%, 10, 1, 0)$ yield more accurate neural networks than their counterparts with $S = 1$. The reason is that a long sequential training reduces the prediction error.

For model-constrained neural networks, we use $\alpha = 10^5$ as the regularization parameter for all cases. We tested with different values for α and almost the same results are obtained for larger values, while smaller values make neural networks perform similarly to the pure data-driven machine learning networks. It can be seen in Figure 2 that training with randomized data returns neural networks, regardless of S, R values, as good as the coarse finite difference approximation with $n_x = 200$. This is expected as we constrain the training with a coarse finite difference model. The trained weight matrices and bias vectors for these neural networks corresponding to three settings $(d200, 1\%, 10, 1, 10^5)$, $(d200, 1\%, 1, 1, 10^5)$ and $(d200, 1\%, 1, 5, 10^5)$ are shown in Figure 5. Again, the bias vectors does not have significant role in the predictions. Note that, unlike those from purely data-driven in Figure 3 which have arbitrary structure, the model-constrained weight matrices, after ignoring small elements, have the same structure as the first-order upwind scheme matrix. Among these neural networks, the long sequential model-constrained network with $R = 5$ is closest to the first-order upwind scheme matrix. It is not surprising as the neural network is constrained to satisfy the first-order upwind scheme in multiple time steps. On the other hand, the neural network trained with noise-free data shows instability starting from the 2000th time step in long-term predictions. This instability is due to the lack of regularizations as compared to the randomized cases for which regularizations are explicit via the model-constrained term and implicit via randomization (see Section 2.4).

Implicit time integration with learned network. One of the advantages of our proposed tangent manifold learning approach is that once trained the learned tangent manifold can be used with any time discretization method. To demonstrate this, we use the learned neural network tangent for the setting $(d200, 0\%, 1, 1, 0)$ with both backward and forward Euler schemes using a time stepsize $\Delta t' = \frac{50}{3} \Delta t$ which is much larger than the training stepsize. It can be seen in Figure 6 that the forward Euler solutions blow up for both learned and true tangent manifolds, which is obvious as the time stepsize is much larger than the stable time stepsize. Both approaches are stable with implicit integration and the results are comparable (though the learned

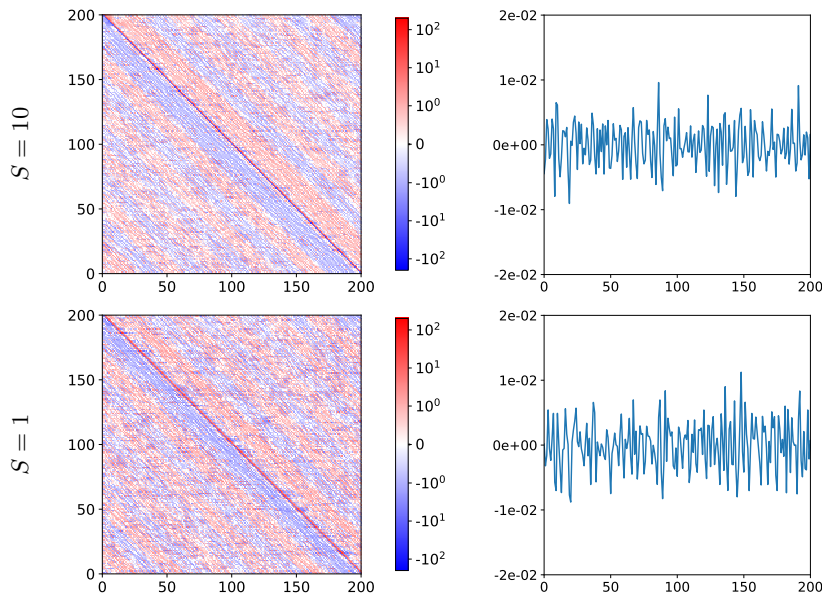


Fig. 3: **Wave/transport equation.** Pure data-driven trained linear neural network parameters: weight matrix heat maps (*left column*) and bias vector magnitudes (*right column*) with $\alpha = 0, \delta = 0\%$.

tangent manifold was trained with a smaller time step size).

Direct learning versus tangent learning. We now compare our tangent manifold learning and direct learning. Here, by direct learning we mean learning the map from \mathbf{u}^i to \mathbf{u}^{i+1} for two consecutive time steps. Clearly, unlike the former, the latter is tailored, and thus limited, to a particular space-time discretization. To be fair, we also use the linear network with zero bias for the direct learning approach. Figure 7 presents the mean-square error of predictions obtained by direct neural networks and tangent manifold networks, both with and without model-constrained terms. As can be seen, both direct and tangent manifold neural networks are comparable in terms of accuracy. However, the learned weight matrices of direct neural networks do not have the pattern of the underlying space-time discretization matrices, and this can be observed from Figure 8. That is, while our tangent manifold approach preserves the structure of spatial discretization, the direct approach, which seems to be natural, does not.

3.2. 2D Burger's equation. We consider the following viscous 2D Burger's equations

$$\begin{aligned} \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} &= \nu \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) \\ \frac{\partial v}{\partial t} + v \frac{\partial v}{\partial x} + u \frac{\partial v}{\partial y} &= \nu \left(\frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right), \end{aligned}$$

where $x, y \in [0, 1]$ and $t \in (0, T]$. The boundary condition is periodic and the initial velocity is given by $v(x, y, 0) = v_0(x, y) = 1$ and $u(x, y, 0) = u_0(x, y)$. We take the

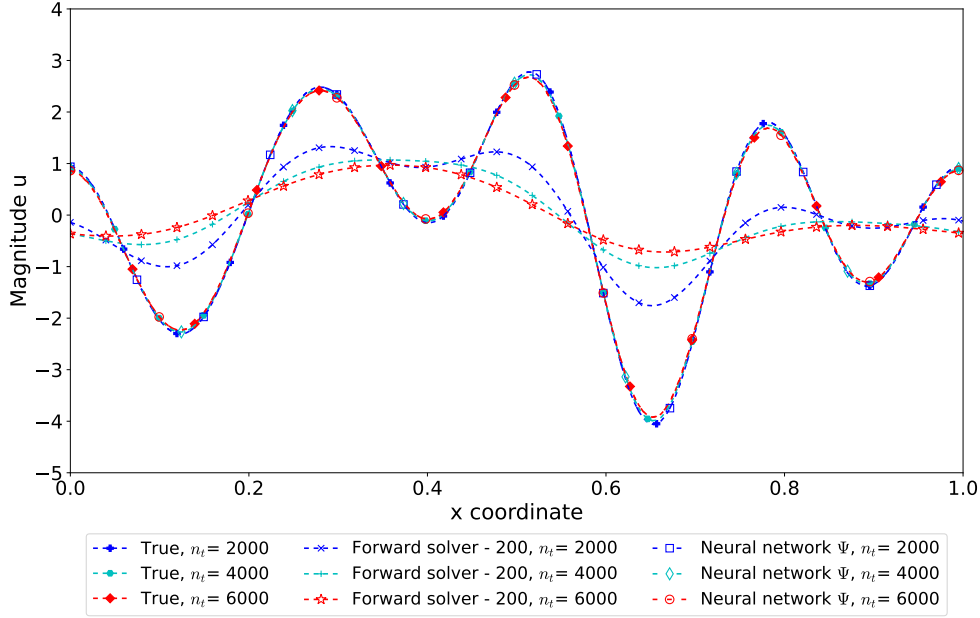


Fig. 4: **Wave/transport equation.** The predicted solutions (Neural network) at time steps $n_t = 2000, 4000, 6000$ by learned neural network corresponding to ($d=200, S = 1, R = 1, \alpha = 0, \delta = 0\%$), finite difference solutions (Forward solver) on coarse grid with $n_x = 200$, and the high resolution solutions (True).

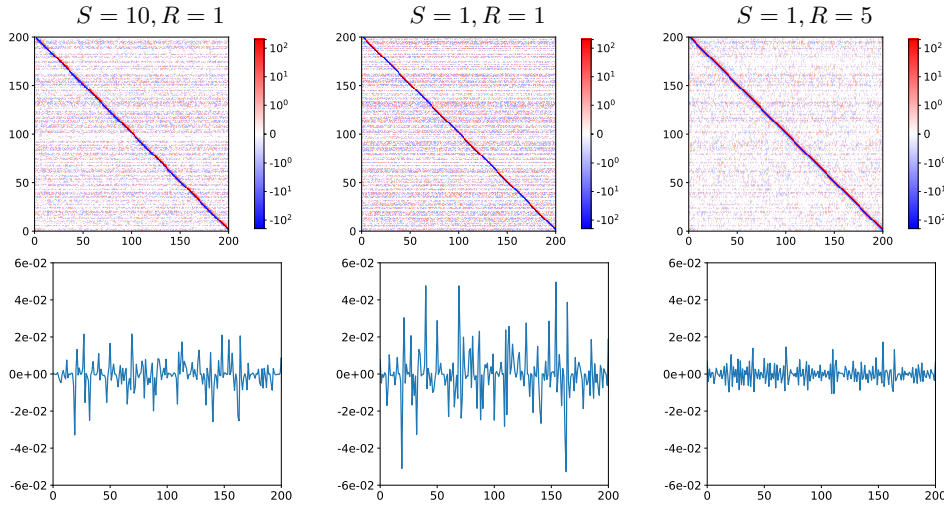


Fig. 5: **Wave/transport equation.** Trained model-constrained linear neural network parameters: weight matrix heat map (*top row*) and bias vector magnitude (*bottom row*) with $\alpha = 1e^5, \delta = 1\%$.

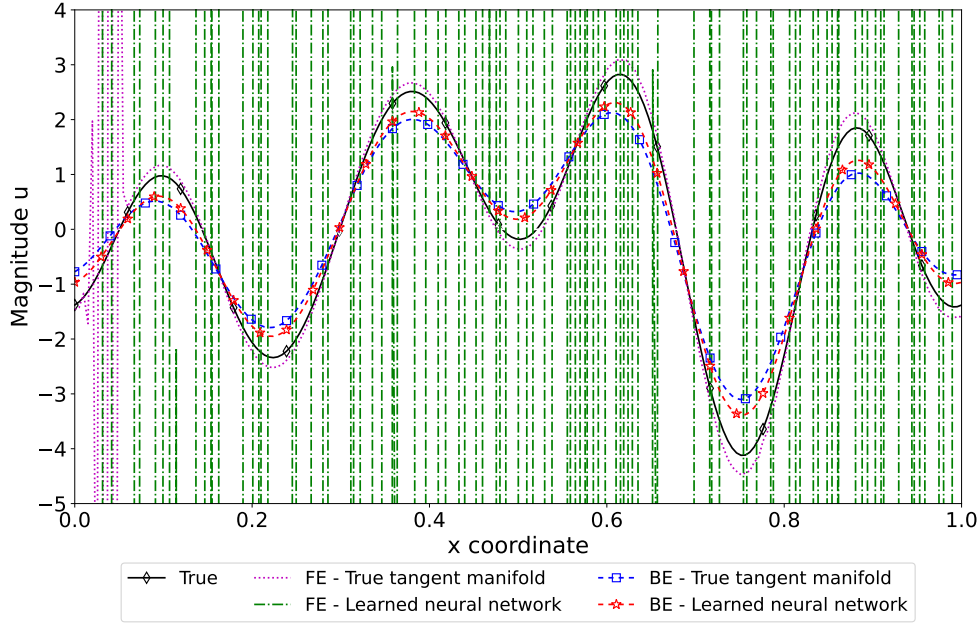


Fig. 6: **Wave/transport equation.** Predicted solutions at $t = 0.1$ obtained by forward Euler scheme (FE) and backward Euler scheme (BE) with the true tangent manifold and learned neural network counterparts with time stepsize $\Delta t' = \frac{50}{3} \Delta t = \frac{25}{3} \times 10^{-3}$.

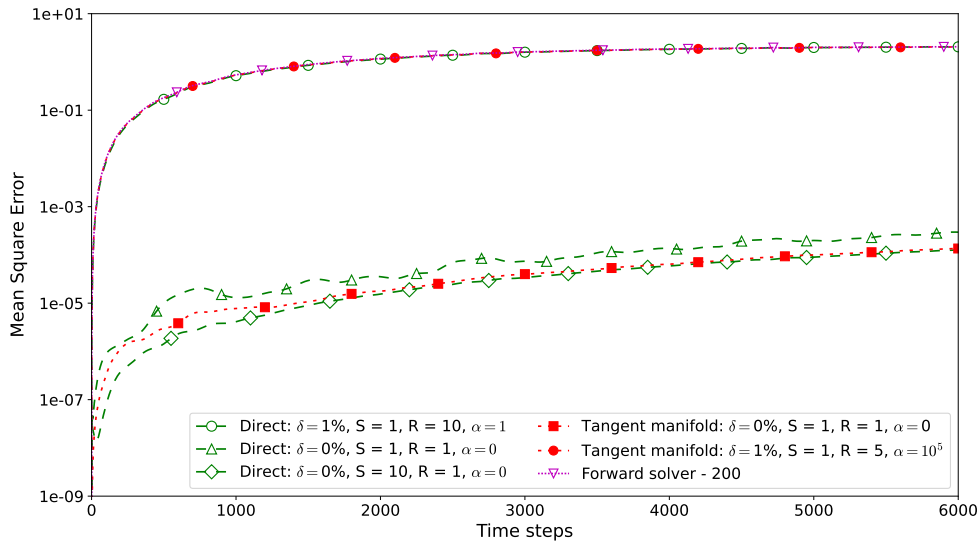


Fig. 7: **Wave/transport equation.** Comparison between direct neural networks (Direct) and tangent manifold neural networks (Tangent manifold).

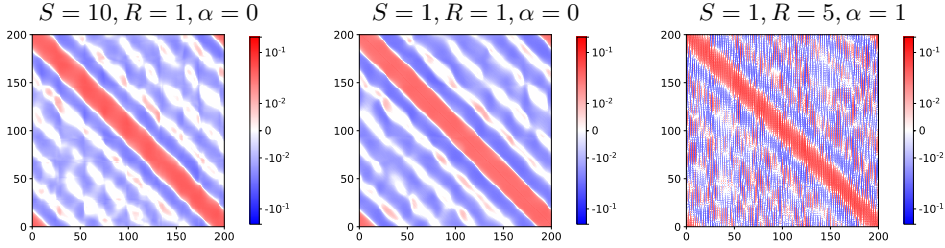


Fig. 8: **Wave/transport equation.** Heat map of weight matrix of direct linear neural networks.

viscosity coefficient to be $\nu = 10^{-2}$. We aim to predict x -velocity u in the time interval $t \in (0, 1.5)$ given an initial velocity $u_0(x, y)$ at $t = 0$.

Data generation. We draw periodic samples of \mathbf{u} using the truncated Karhunen-Loève expansion

$$u_0(x, y) = \exp \left(\sum_{i=1}^{15} \sqrt{\lambda_i} \omega_i(x, y) z_i \right),$$

where $\mathbf{z} = \{z_i\}_{i=1}^{15} \sim \mathcal{N}(0, \mathbf{I})$, and (λ, ω) are eigenpairs obtained by the eigendecomposition of the covariance operator $7^{\frac{3}{2}} (-\Delta + 49\mathbf{I})^{-2.5}$, where Δ is the Laplacian operator, with periodic boundary conditions. Training data corresponding to each initial velocity is generated from a 128×128 high-resolution spatial mesh and 1000 time steps for the time horizon $T = 0.1$ using finite difference method. These high resolution solutions are down-sampled on a coarser mesh of 100 time steps ($\Delta t = 10^{-3}$) and 32×32 spatial mesh. These down-sampled solutions are treated as true solutions for the training process. Meanwhile, we draw 10 test initial velocity samples independently, and the corresponding test data set of 10 samples is created in the same manner. However, the time horizon $T = 1.5$ for test samples is chosen—much larger than the trained time horizon—with time stepsize $\Delta t = 10^{-3}$. This helps us test the accuracy and stability of neural network solutions beyond the training regime.

Neural network architecture. We use a shallow network of one layer with 5000 neurons for all cases to approximate the tangent manifold of Burger’s equations. Note that we have compared the one-layer network with two- and three-layer networks with different numbers of neurons ranging from 100 to 5000. These deeper networks perform poorly with small data sets and are improved with large data sets in which the shallow one has comparable performance. Note that one-layer neural network approximation capability are rigorously justified by past universal approximation theories (see, e.g. [8, 12, 27, 17]) and our current work [6]. Thus we shall use a one-layer neural network for all numerical results. In addition, ReLU [31] is used as the activation function. ADAM optimizer is used with the learning rate of 10^{-4} and the training batch size is 40 samples. For this example, reasonably optimal weights/biases are the ones giving the lowest accumulated mean square error after 1500 time steps for 10 test data. We take $\alpha = 10^5$ for the regularization parameter as this gives the best results from our numerical experiments (not shown here).

Comparison of different learned neural networks. Figure 9 presents the comparison of mean square error obtained by different learned neural networks with the data set of 200 samples. It can be seen that, in general, the model-constrained

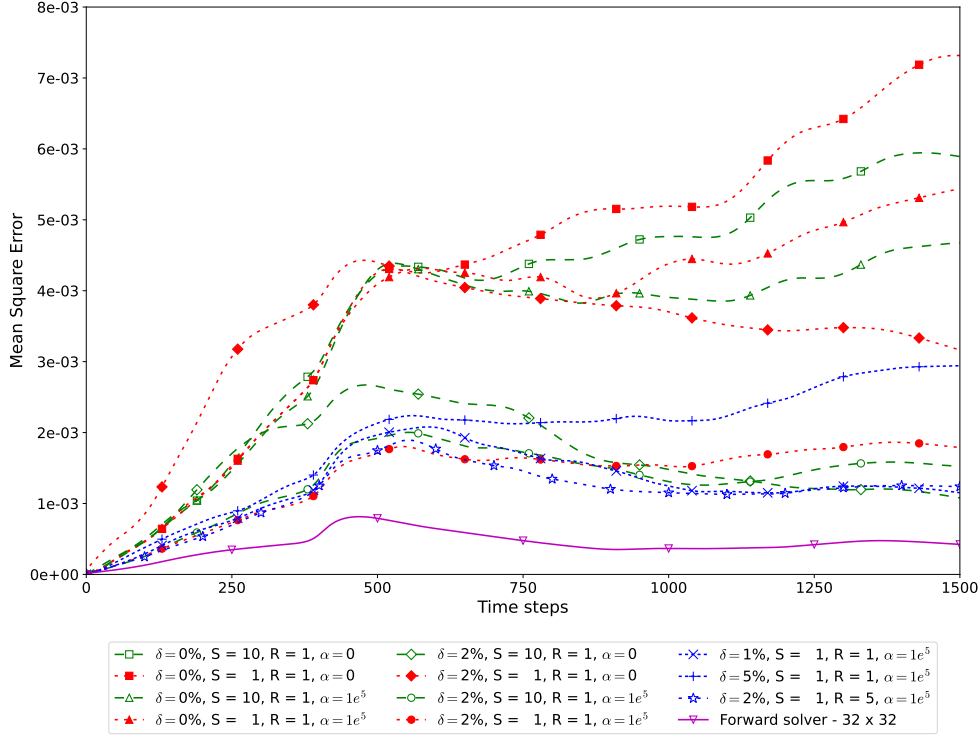


Fig. 9: **Burger’s equations.** Comparison of mean square error among different neural networks trained with 200 data with/without noise. Recall that $\alpha = 0$ corresponds to the pure data-driven neural network training without model-constrained terms. Forward solver denotes the numerical solution on 32×32 spatial mesh.

neural networks are far better than their pure data-driven counterparts (i.e. with $\alpha = 0$). Additionally, long sequential machine learning trainings with $S = 10$ provide slightly better accuracy than $S = 1$, except for the noisy data with pure data-driven network in which the improvement is significant.

For model-constrained neural networks, long sequential training results with $S = 10$ in two settings ($(d200, 0\%, 10, 1, 10^5)$ and $(d200, 2\%, 10, 1, 10^5)$) show a marginal improvement compared to short sequential training with $S = 1$ in two settings ($(d200, 0\%, 1, 1, 10^5)$ and $(d200, 2\%, 1, 1, 10^5)$). Therefore, $S = 1$ is sufficient and we use it for the rest of numerical results with model-constrained neural networks. Figure 9 shows that using 5% noise causes the neural network corresponding to $(d200, 5\%, 1, 1, 10^5)$ to perform poorly, while 1% noise gives almost the same accuracy as 2% noise. It is noticeable that the long sequential model-constrained training with $R = 5$ ($(d200, 2\%, 1, 5, 10^5)$) yields higher accuracy than the others. However, large R is more computationally expensive since many passes through the back-propagation computational graph are needed.

Long-time predictions with small and large training data sets. As discussed above, since long sequential machine learning training does not provide significant improvement, we consider $S = 1$ for numerical results using large data sets in Figure 10. As can be seen, compared to 200 data samples, training with 600 data

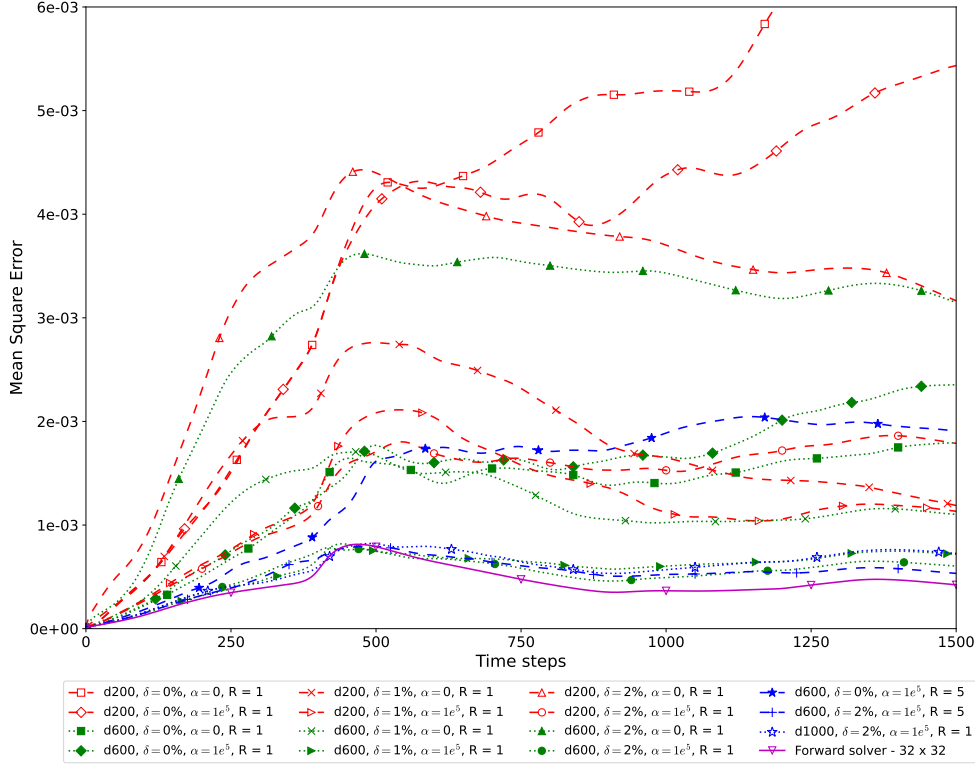


Fig. 10: **Burger’s equations**. The mean square error versus the number of time steps for various learned neural networks using 200, 600, and 1000 data samples with $S = 1$.

samples provides more accurate predictions. Moreover, model-constrained neural networks with randomized data are the most accurate among others (model-constrained with noise-free data and pure machine learning with/without randomized data). We can also observe that using more than 600 data samples does not provide significant improvements but is more expensive. Unlike the case with 200 data samples, long and short sequential model-constrained trainings with $R = 5$ and $R = 1$, respectively, provide similar results for 600 data samples. This is expected as richer data reduces the significance of the model-constrained term.

As shown in Figure 11, predicted solutions obtained by the model-constrained approach (the fifth row) with data randomization are in good agreement with the ground-truth counterparts. On the contrary, the pure data-driven approach with data randomization (the third row) shows poor long-time predictions. We also observe that both pure data-driven learning solutions and model-constrained solutions (the second and fourth rows, respectively) without randomization are unstable for long-time predictions. It is not surprising since both do not have sufficient regularizations compared to the randomized cases in which extra regularizations are implicitly performed (see Section 2.4). Moreover, regularizations induced by data randomization shown in Section 2.4 stabilize the network predictions and this can be clearly seen by comparing the third and the second rows for the pure data-driven learning

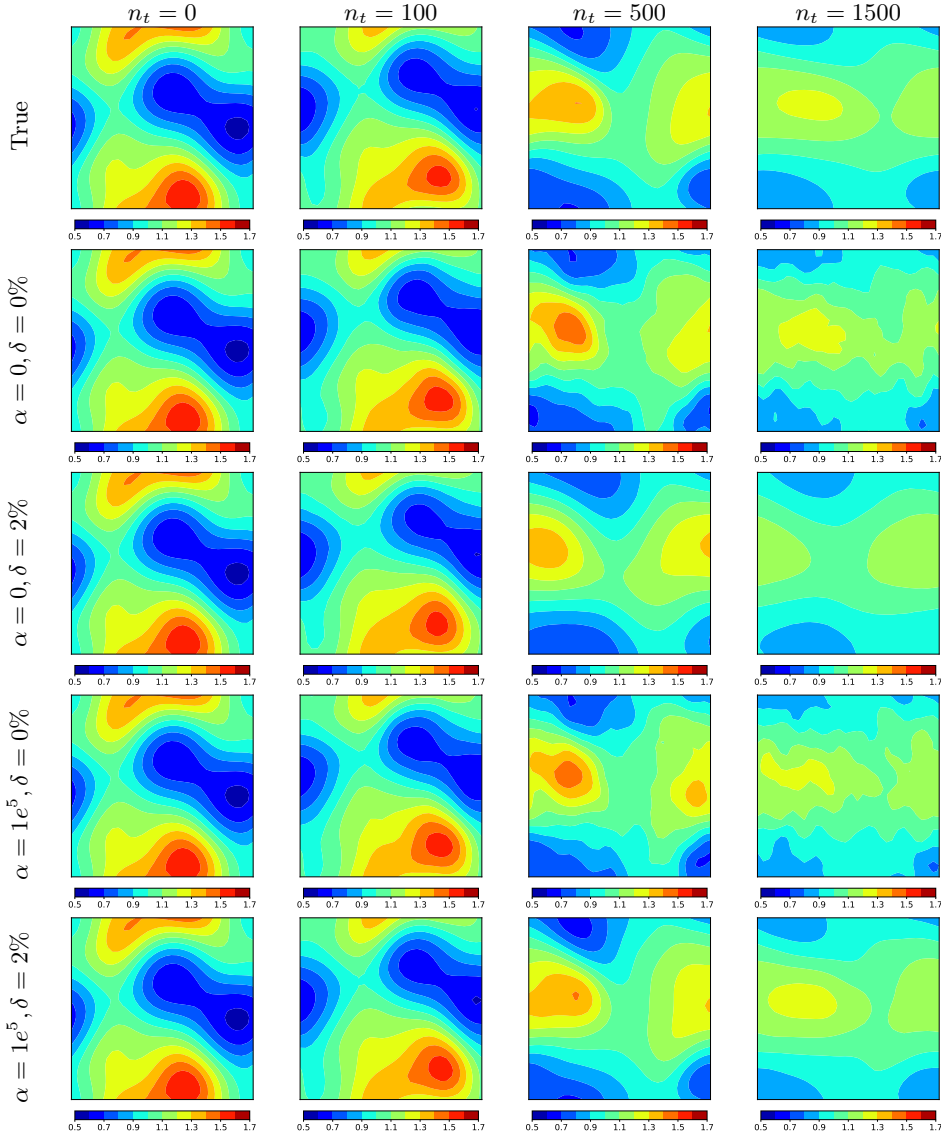


Fig. 11: **Burger's equations.** Predicted solutions at different time steps (n_t) obtained by various learned neural network tangents with 600 training data samples and $S = 1$ and $\Delta t = 10^{-3}$. *Top row:* True means high-resolution solution; *Second row:* pure data-driven network without data randomization; *Third row:* pure data-driven network with noisy data; *Fourth row:* model-constrained network without data randomization; *Fifth row:* model-constrained network with noisy data.

approach, and by comparing the fourth and the fifth rows for the model-constrained learning approach.

Figure 12 plots the contours of the learned and the true tangent manifolds. Clearly, the learned model-constrained tangent manifold with data randomization

provides the best agreement with the true tangent manifold. This is not surprising as both the governing equations (explicit via model-constrained term) and sufficient regularizations (implicit via data randomization) are incorporated.

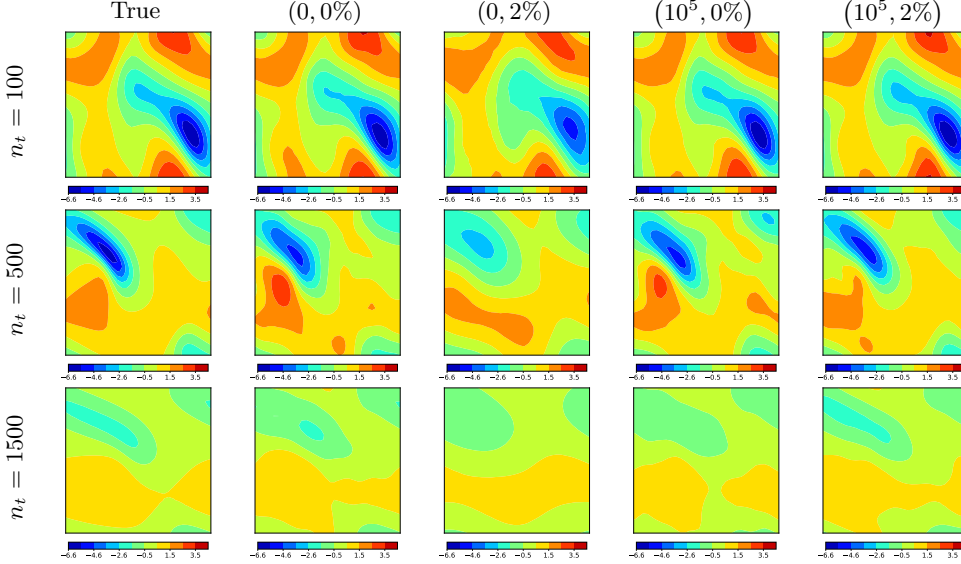


Fig. 12: **Burger's equations.** Contours of True and various learned tangent manifolds. Contours are plotted at different time steps n_t for four different combinations of regularization parameter α and noise level δ . For all cases, we use 600 data samples, $S = 1$, and $\Delta t = 10^{-3}$. *First column:* True tangent manifold (0,0%); *Second column:* pure data-driven tangent manifold without data randomization (0,2%); *Third column:* pure data-driven tangent manifold with data randomization (10^5 ,0%); *Fourth column:* model-constrained tangent manifold without data randomization; *Fifth column:* model-constrained tangent manifold with data randomization (10^5 ,2%).

Predictive flexibility in time for mcTangent approach As discussed above, one appealing feature of tangent manifold learning is that once trained it can be used to solve for approximate solutions with smaller or larger time stepsizes, despite the fact that it is trained based on a particular spatial discretization. On the contrary, direct learning is attached to a space-time discretization. Figure 13 shows the model-constrained tangent manifold learning solutions and contours of the corresponding learned tangent manifold at various times for the setting $(d600, 2\%, 1, 1, 10^5)$. Here we use half of the training time stepsize $\Delta t' = \frac{1}{2}\Delta t = 5 \times 10^{-4}$. It can be seen that these predictions are indistinguishable from ones (the fifth row in Figure 11 for prediction solutions and the fifth column in Figure 12 for predicted tangent manifolds) obtained by using the training time stepsize $\Delta t = 10^{-3}$ with the same learned network.

Implicit time integration with learned network. Another appealing feature of tangent manifold learning is that once trained it can be deployed with any time discretization schemes. We use the learned network from the setting $(d600, 2\%, 1, 1, 10^5)$ together with the backward Euler method with a larger time stepsize $\Delta t' = 12.5\Delta t = 1.25 \times 10^{-2}$, where $\Delta t = 10^{-3}$ is the training stepsize. Shown in Figure 14 are predicted solutions at $t = \{0, 0.1, 0.5, 1.5\}$ corresponding to 0, 100, 500, 1500th time steps.

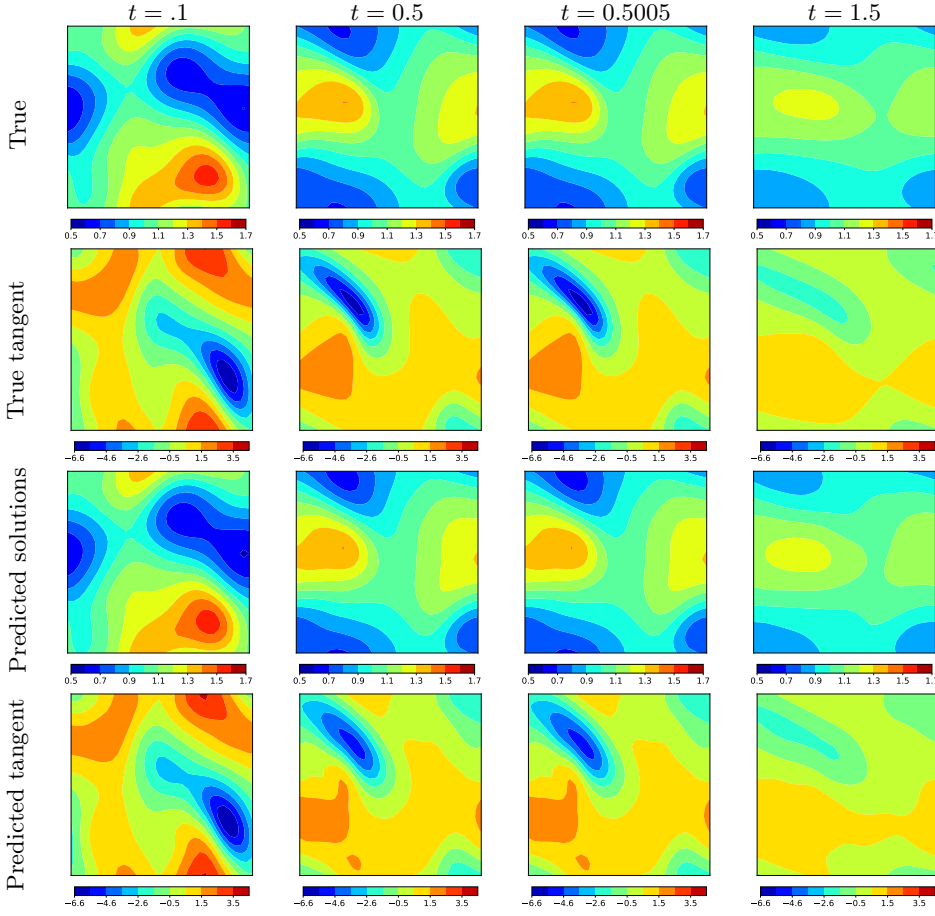


Fig. 13: **Burger’s equations**. Predicted solutions and tangent manifold using by `mcTangent` neural networks with $(d600, 2\%, 1, 1, 10^5)$, and time step $\Delta t'' = \frac{1}{2}\Delta t = 5 \times 10^{-4}$. *First row*: True means high-resolution solutions; *Second row*: contours of True tangent manifold (True tangent); *Third row*: Predicted `mcTangent` solutions; *Fourth row*: contours of `mcTangent` tangent manifold (Predicted tangent).

We observe that solutions using the forward Euler scheme, regardless of using the true tangent manifold or learned one (second and third rows, respectively), are unstable as the time stepsize $\Delta t'$ is too big for stability. On the contrary, using the backward Euler scheme, `mcTangent` solutions are comparable to the true counterparts (fourth and fifth rows, respectively). Clearly, due to large time stepsize, both are more diffusive compared to the true solutions with small time stepsize Δt in the first row.

Direct learning versus tangent learning. Recall that by direct learning we mean learning the map from \mathbf{u}^i to \mathbf{u}^{i+1} for two consecutive time steps. We investigate the difficulty and complexity of direct learning. Specifically, we use a data set with 600 samples with/without data randomization to learn the neural network with one layer of 5000 neurons that maps velocities from one step to the next. As shown in Figure 15, the direct learning approach (with the best combination of hyperparameters) for

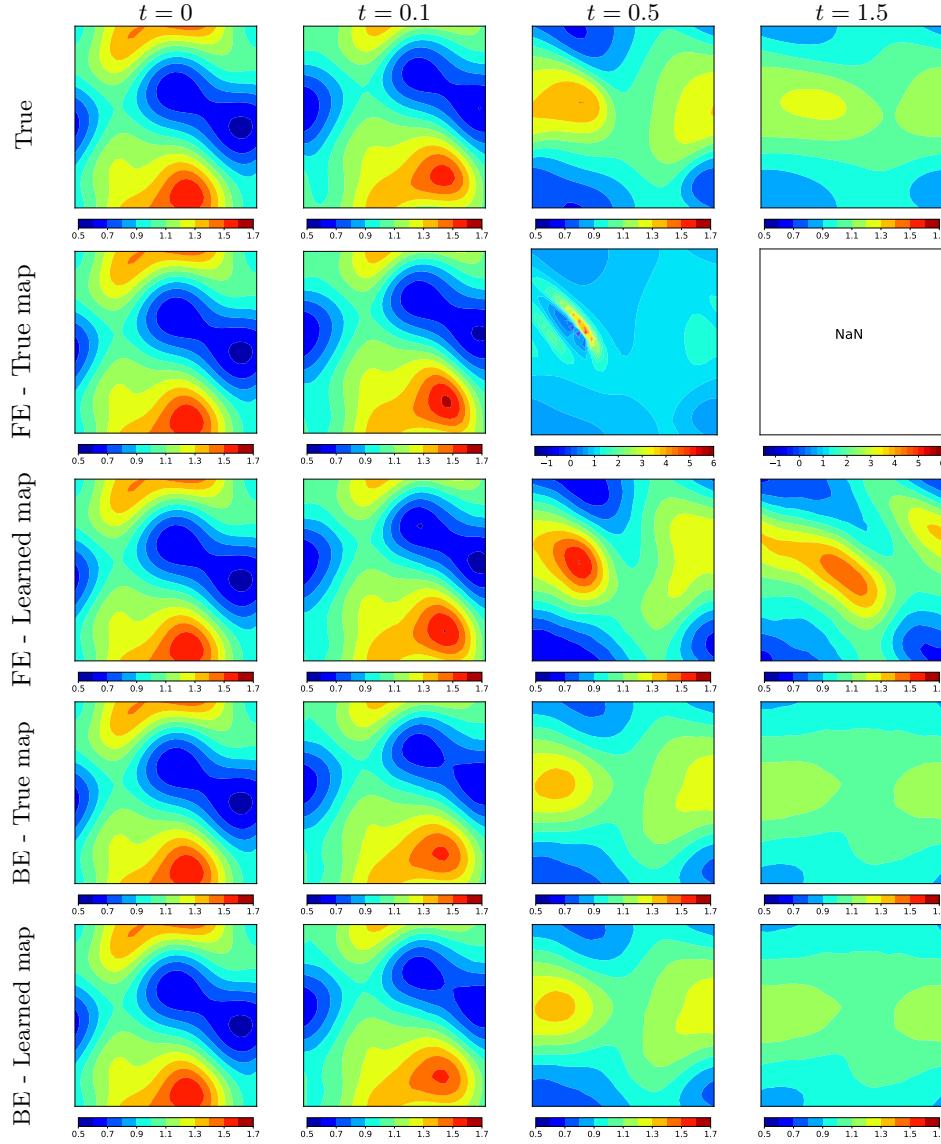


Fig. 14: **Burger's equations.** Predicted solutions at different times obtained by Forward Euler (FE) scheme and Backward Euler (BE) scheme using large stepsize $\Delta t' = 12.5\Delta t = 1.25 \times 10^{-2}$ with the true tangent manifold (True map) and the learned one (Learned map) for the setting $(d600, 2\%, 1, 1, 10^5)$.

the setting $(d600, 2\%, 1, 3, 2)$ is less accurate for both short-time and long-time predictions compared to the tangent learning counterpart with even a smaller data set of 200 samples with the setting $(d200, 2\%, 1, 1, 10^5)$. Interestingly, unlike the tangent learning approach, the direct learning approach, both pure data-driven and model-constrained approaches, trained with randomized data is less accurate compared to noise-free data in short sequential training $S = 1$. On the other hand, data random-

ization does not have visible benefits on long sequential training $S = 10$. Specifically, both $(d600, 0\%, 10, 1, 0)$ and $(d600, 2\%, 10, 1, 0)$ settings behave similarly.

Also seen in Figure 15, among pure data-driven networks ($\alpha = 0$) with direct learning, long sequential machine learning training with $S = 10$ is the most accurate. Model-constrained network with direct learning for the setting $(d600, 2\%, 1, 1, 10)$ is much more accurate compared to the pure data-driven network with direct learning for the same setting. Moreover, sequential model-constrained networks for $R = 2, 3$ corresponding to two settings $(d600, 2\%, 1, 2, 2)$ and $(d600, 2\%, 1, 3, 2)$ are comparable to much longer sequential machine learning network with $S = 10$ for the setting $(d600, 2\%, 10, 1, 0)$. In the presented results, it is important to point out that for direct learning, care must be taken in choosing a good regularization parameter α . For example, $\alpha = 2$ is good for $R = 2, 3$, but $\alpha = 10$ is good for $R = 1$. On the contrary, tangent learning is more robust. In particular, a single $\alpha = 10^5$ works well for all settings. Solutions predicted by direct and tangent learnings (both with model-constrained terms) for $(d600, 2\%, 1, 3, 2)$ and $(d200, 2\%, 1, 1, 10^5)$, respectively, are shown in Figure 16. As can be observed, tangent learning solutions with even smaller data set $(d200, 2\%, 1, 1, 10^5)$ are much more accurate than the direct learning with $(d600, 2\%, 1, 3, 2)$. This is due to the fact that direct learning tries to learn a mixed space-time discretization, which is more difficult than learning only the spatial discretization in the tangent learning.

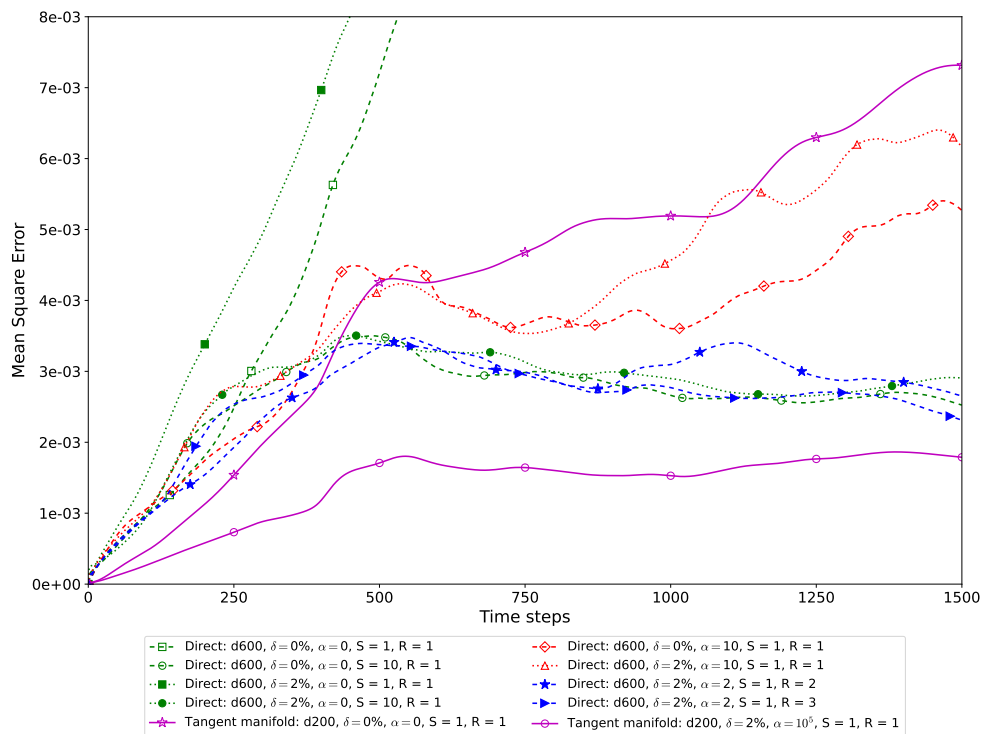


Fig. 15: **Burger's equations.** The mean square error versus time steps obtained by the direct learning approach (Direct) using 600 training samples and the tangent learning approach (Tangent manifold) using 200 training samples.

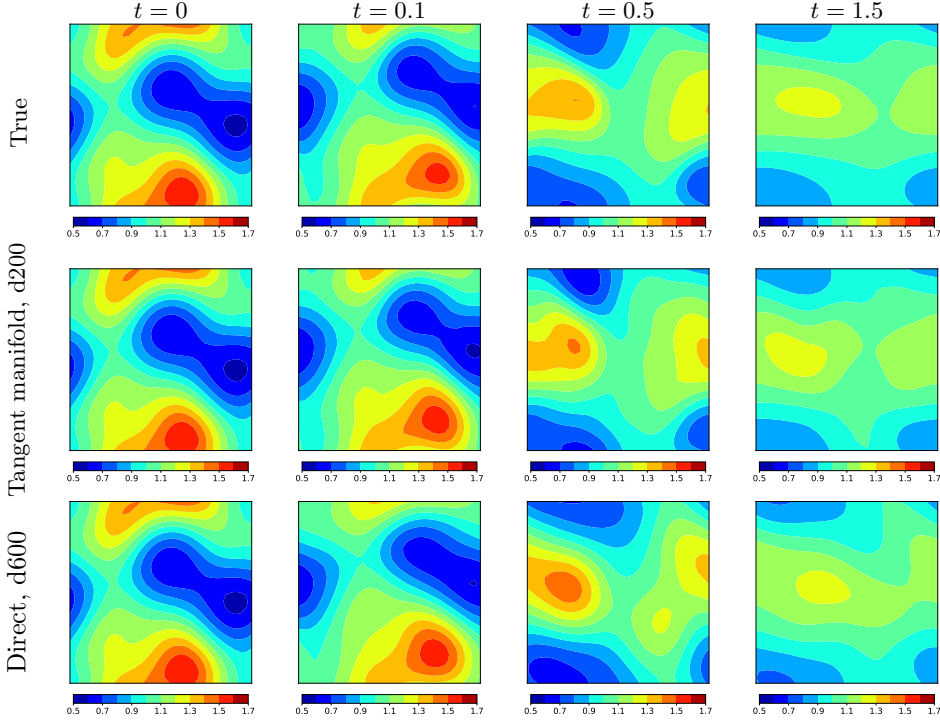


Fig. 16: **Burger’s equations.** Solutions at different times, *First row*: True means high-resolution solutions; *Second row*: learned tangent manifold neural network solutions with $(d200, 2\%, 1, 1, 10^5)$; *Third row*: learned direct neural network solutions with $(d600, 2\%, 1, 3, 2)$.

3.3. Navier-Stokes equation. The vorticity form of the 2D Navier-Stokes equation for viscous and incompressible fluid [23] can be written as

$$\begin{aligned} \partial_t u(x, t) + v(x, t) \cdot \nabla u(x, t) &= \nu \Delta u(x, t) + f(x), & x \in (0, 1)^2, t \in (0, T] \\ \nabla \cdot v(x, t) &= 0, & x \in (0, 1)^2, t \in (0, T] \\ u(x, 0) &= u_0(x), & x \in (0, 1)^2 \end{aligned}$$

where $v(x, t)$ is the velocity field, $u = \nabla \times v$ is the vorticity, u_0 is the initial vorticity, $f(x) = 0.1(\sin(2\pi(x_1 + x_2)) + \cos(2\pi(x_1 + x_2)))$ is the forcing function and $\nu = 10^{-3}$ is the viscosity coefficient. Our goal is to solve for the vorticity $u(x, y, t)$ given the initial condition u_0 at $t = 0$ by a trained tangent network Ψ .

Data generation. Data pair (\mathbf{u}, \mathbf{y}) is generated by a similar procedure outlined for Burger’s equation problem in Section 3.2. In particular, we draw samples of \mathbf{u}_0 using the truncated Karhunen-Loève expansion

$$\mathbf{u}_0 = \sum_{i=1}^{15} \sqrt{\lambda_i} \omega_i(x) z_i,$$

where $z_i \sim \mathcal{N}(0, 1)$, $i = 1, \dots, 15$, and (λ, ω) is eigenpairs obtained by the eigen-decomposition of the covariance operator $7^{\frac{3}{2}}(-\Delta + 49\mathbf{I})^{-2.5}$ with periodic boundary

conditions. Next, given initial vorticity \mathbf{u}_0 , we solve the Navier-Stokes equation by the stream-function formulation with a pseudospectral method [23]. High resolution solutions are obtained on a uniform 128×128 spatial mesh and uniform 1000 time steps in $(0, 2)$. The high-resolution solutions are then down-sampled on a coarser mesh 32×32 in space and 200 uniform time steps, and they are used as the training data. To verify the accuracy of the learned neural network, we draw 10 test samples independently. It turns out that the Navier-stokes equation is much more challenging than Burger’s equation, thus we use 200 time steps for each training data as opposed to 100 for the Burger equation. Similar to the above, to challenge the learned network we use 1500 time steps for testing, and thus the testing time horizon is far beyond the training time horizon.

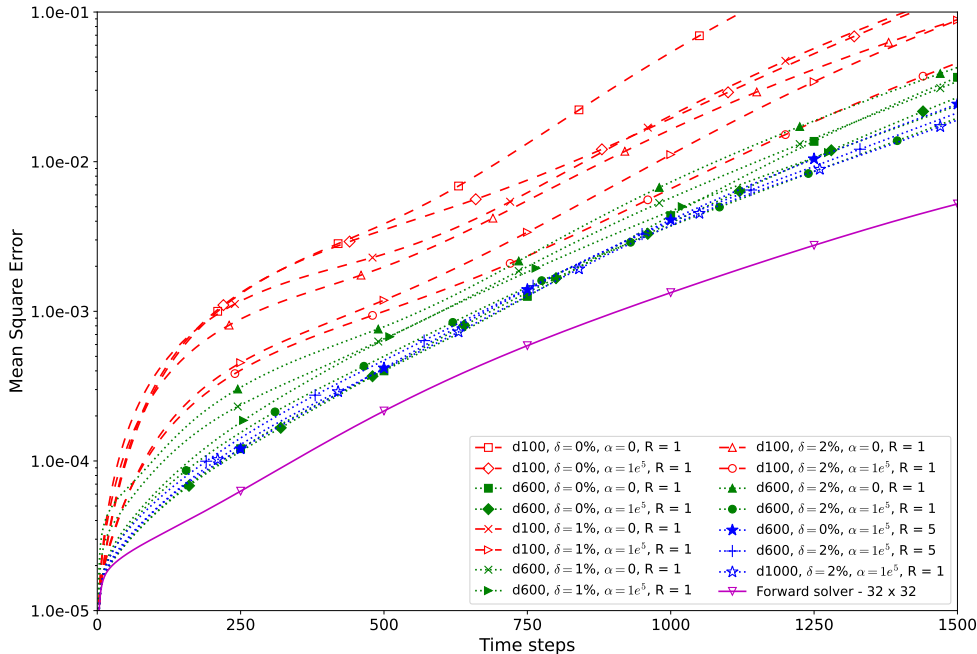


Fig. 17: **Navier-Stokes equation.** The mean-square error versus time steps obtained by the various learned neural network with $S = 1$.

Neural network architecture. With the same observation for the Burger equation in Section 3.2, we use a shallow network of one layer with 5000 neurons using ReLU activation function. ADAM optimizer with default parameters is used with the learning rate of 2×10^{-4} , while the training batch size is 2 samples. The chosen “optimal” network is the one having the lowest accumulated mean square error after 1500 time steps for 10 testing samples. Following the wave and Burger examples, we pick a relatively large value for the model-constrained regularization parameter $\alpha = 10^5$.

Long-time predictions. Figure 17 shows the mean-square error of predictions and ground truth solutions as a function of time steps. It can be seen that training with a large data set with 600 samples provides much more accurate solutions than with small data set with 100 samples. On the one hand, among learned neu-

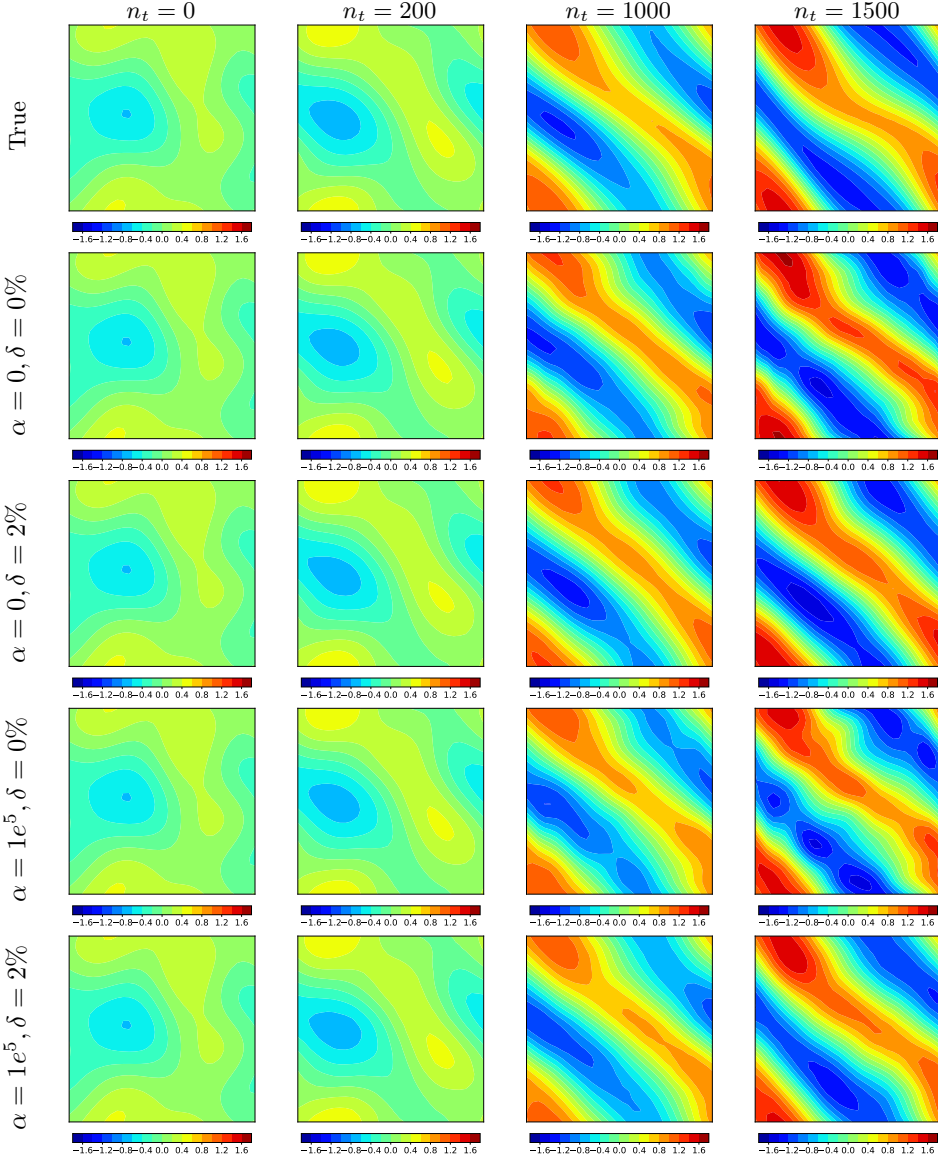


Fig. 18: **Navier-Stokes equation.** Predicted solutions at different time steps obtained by various trained networks with 600 data samples. *First row:* ground truth; *Second row:* pure data-driven network solutions with noise-free data; *Third row:* pure data-driven network solutions with randomized data; *Fourth row:* model-constrained network solutions with noise-free data; *Fifth row:* model-constrained network solutions with randomized data.

ral networks trained with 100 data samples, the model-constrained network with data randomization for $(d100, 2\%, 1, 1, 10^5)$ setting is far closer to the true solution than the other networks. This implies that the model-constrained approach has a significant

contribution to producing accurate predictions in the context of small data. In the case of richer data set with 600 samples, networks with two settings ($d600, 0\%, 1, 1, 0$) and ($d600, 0\%, 1, 1, 10^5$) trained with noise-free data show a good performance in the short time predictions, while the long-time predictions deteriorate. Noticeably, between these two networks, the model-constrained one has more accurate predictions starting from the 500th time step. In the meantime, with the same data set with 600 samples, pure data-driven neural networks trained with higher noise level data give a higher error, for example, ($d600, 2\%, 1, 1, 0$) neural network predictions are less accurate than those obtained from ($d600, 1\%, 1, 1, 0$). In contrast, model-constrained network with 2% noise level ($d600, 2\%, 1, 1, 10^5$) is superior to 1% noise level ($d600, 1\%, 1, 1, 10^5$). Another point is that as we increase the sequential model-constrained value to $R = 5$, we obtain good predictions in both short-time and long-time intervals. Two model-constrained networks with ($d600, 0\%, 1, 5, 10^5$) and ($d600, 2\%, 1, 5, 10^5$) are comparable to the network with much larger data set ($d1000, 0\%, 1, 1, 10^5$) without randomization. However, the noisy data network ($d600, 2\%, 1, 5, 10^5$) outperforms the noise-free

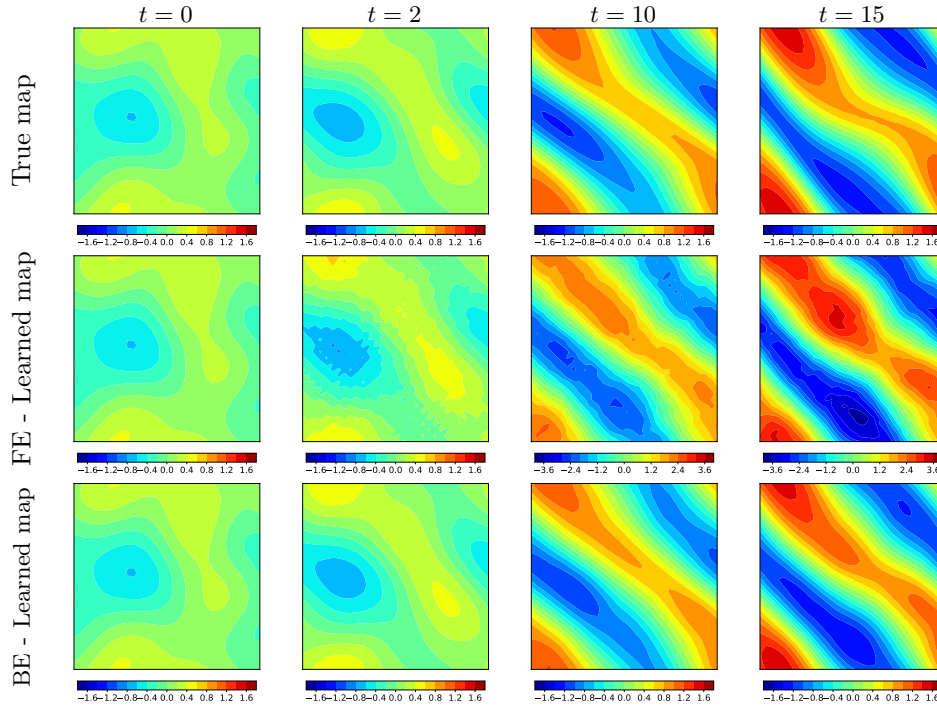


Fig. 19: **Navier-Stokes equation.** Comparison of various neural network solutions. *Top row:* Spectral method (Crank–Nicolson time integration scheme) with the true tangent manifold; *Second row:* Forward Euler (FE) scheme (a different scale bar in the third and fourth columns); *Third row:* Backward Euler (BE) scheme with the learned neural network for 20 times larger time stepsize $\Delta t' = 20\Delta t = 0.2$.

one ($d600, 0\%, 1, 5, 10^5$) in the long-time predictions. In summary, model-constrained network with data randomization outperforms all other networks. Given a test initial vorticity, the plots of predicted solutions obtained by different learned networks are

shown in Figure 18. As can be seen, the model-constrained network with the setting $(d600, 2\%, 1, 1, 10^5)$ provides the most accurate solutions as opposed to others trained from the same data set.

Implicit time integration with learned network. We used the learned network for backward Euler scheme with 20 times larger time stepsize, $\Delta t' = 20\Delta t = 0.2$, compared to training stepsize $\Delta t = 0.01$. As shown in Figure 19, forward Euler scheme with the learned network shows severe instability, while the backward Euler scheme with the learned network solutions are in good agreement with the spectral solution with Crank–Nicolson scheme with much smaller time stepsize.

4. Conclusions. We have presented a model-constrained tangent manifold learning (**mcTangent**) approach to simulate dynamical systems in real-time. At the heart of **mcTangent** is a careful craft synergizing several desirable strategies: i) a tangent manifold learning to take advantage of the neural network speed and time-accurate nature of the method of lines; ii) a model-constrained approach to encode the neural network tangent with the underlying physics; iii) sequential learning strategies to promote long-time stability and accuracy; and iv) data randomization approach to implicitly regularize the smoothness of the neural network tangent and its likeliness to the truth tangent up second order derivatives in order to further enhance the stability and accuracy of **mcTangent** solutions. Both semi-heuristic and rigorous arguments are provided to analyze and justify the proposed approach. Several numerical results for transport equation, viscous Burger’s equation, and Navier-Stokes equation are presented to study and demonstrate the capability of the proposed **mcTangent** learning approach. Further theoretical analysis of **mcTangent** with both sequential learning strategies is ongoing to provide a deeper understanding of the approach. Strategies to improve the accuracy and to strongly encode the underlying governing equations are also part of future work.

Disclosure statement. No potential conflict of interest was reported by the author(s).

REFERENCES

- [1] Guozhong An. The effects of adding noise during backpropagation training on a generalization performance. *Neural computation*, 8(3):643–674, 1996.
- [2] Konstantinos Benidis, Syama Sundar Rangapuram, Valentin Flunkert, Yuyang Wang, Danielle Maddix, Caner Turkmen, Jan Gasthaus, Michael Bohlke-Schneider, David Salinas, Lorenzo Stella, et al. Deep learning for time series forecasting: Tutorial and literature survey. *ACM Computing Surveys (CSUR)*, 2018.
- [3] Chris M Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1):108–116, 1995.
- [4] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [5] Noah D Brenowitz and Christopher S Bretherton. Prognostic validation of a neural network unified physics parameterization. *Geophysical Research Letters*, 45(12):6289–6298, 2018.
- [6] Tan Bui-Thanh. A unified and constructive framework for the universality of neural networks, 2021.
- [7] Yuyao Chen, Lu Lu, George Em Karniadakis, and Luca Dal Negro. Physics-informed neural networks for inverse problems in nano-optics and metamaterials. *Opt. Express*, 28(8):11618–11633, Apr 2020.
- [8] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, Dec 1989.
- [9] Emmanuel De Bézenac, Arthur Pajot, and Patrick Gallinari. Deep learning for physical processes: Incorporating prior scientific knowledge. *Journal of Statistical Mechanics: Theory*

- and *Experiment*, 2019(12):124009, 2019.
- [10] Karthik Duraisamy. Perspectives on machine learning-augmented reynolds-averaged and large eddy simulation models of turbulence. *Physical Review Fluids*, 6(5):050504, 2021.
 - [11] Philipp Holl, Vladlen Koltun, Kiwon Um, and Nils Thuerey. phiflow: A differentiable pde solving framework for deep learning via physical simulations. In *NeurIPS Workshop*, volume 2, 2020.
 - [12] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
 - [13] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-Kelley, and Frédo Durand. DiffTaichi: Differentiable programming for physical simulation. *ICLR*, 2020.
 - [14] Xiaowei Jia, Jared Willard, Anuj Karpatne, Jordan Read, Jacob Zwart, Michael Steinbach, and Vipin Kumar. Physics guided rnns for modeling dynamical systems: A case study in simulating lake temperature profiles. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages 558–566. SIAM, 2019.
 - [15] Jiaqi Jiang, Mingkun Chen, and Jonathan A. Fan. Deep neural networks for the evaluation and design of photonic devices. *Nature Reviews Materials*, Dec 2020.
 - [16] Claes Johnson. *Numerical solution of partial differential equations by the finite element method*. Courier Corporation, 2012.
 - [17] Jesse Johnson. Deep, skinny neural networks are not universal approximators. In *International Conference on Learning Representations*, 2019.
 - [18] Byungsoo Kim, Vinicius C Azevedo, Nils Thuerey, Theodore Kim, Markus Gross, and Barbara Solenthaler. Deep fluids: A generative network for parameterized fluid simulations. In *Computer graphics forum*, volume 38, pages 59–70. Wiley Online Library, 2019.
 - [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - [20] Dmitrii Kochkov, Jamie A Smith, Ayya Alieva, Qing Wang, Michael P Brenner, and Stephan Hoyer. Machine learning-accelerated computational fluid dynamics. *Proceedings of the National Academy of Sciences*, 118(21):e2101784118, 2021.
 - [21] Keisuke Kojima, Bingnan Wang, Ulugbek Kamilov, Toshiaki Koike-Akino, and Kieran Parsons. Acceleration of fdt-d-based inverse design using a neural network approach. In *Advanced Photonics 2017 (IPR, NOMA, Sensors, Networks, SPPCom, PS)*, page ITu1A.4. Optical Society of America, 2017.
 - [22] Randall J LeVeque. *Finite volume methods for hyperbolic problems*, volume 31. Cambridge university press, 2002.
 - [23] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations, 2020.
 - [24] Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194):20200209, 2021.
 - [25] Lu Lu, Xuhui Meng, Zhiping Mao, and George Em Karniadakis. Deepxde: A deep learning library for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
 - [26] Lu Lu, Raphael Pestourie, Wenjie Yao, Zhicheng Wang, Francesc Verdugo, and Steven G. Johnson. Physics-informed neural networks with hard constraints for inverse design, 2021.
 - [27] Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The expressive power of neural networks: A view from the width. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
 - [28] Nils Margenber, Dirk Hartmann, Christian Lessig, and Thomas Richter. A neural network multigrid solver for the navier-stokes equations. *Journal of Computational Physics*, 460:110983, 2022.
 - [29] Kiyotoshi Matsuoka. Noise injection into inputs in back-propagation learning. *IEEE Transactions on Systems, Man, and Cybernetics*, 22(3):436–440, 1992.
 - [30] Jeremy Morton, Antony Jameson, Mykel J Kochenderfer, and Freddie Witherden. Deep dynamical modeling and control of unsteady fluid flows. *Advances in Neural Information Processing Systems*, 31, 2018.
 - [31] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
 - [32] Shaowu Pan and Karthik Duraisamy. Long-time predictive modeling of nonlinear dynamical systems using neural networks. *Complexity*, 2018, 2018.
 - [33] Jaideep Pathak, Mustafa Mustafa, Karthik Kashinath, Emmanuel Motheau, Thorsten Kurth,

- and Marcus Day. Using machine learning to augment coarse-grid computational fluid dynamics simulations. *arXiv preprint arXiv:2010.00072*, 2020.
- [34] Raphaël Pestourie, Youssef Mroueh, Thanh V. Nguyen, Payel Das, and Steven G. Johnson. Active learning of deep surrogates for pdes: application to metasurface design. *npj Computational Materials*, 6(1):164, Oct 2020.
- [35] John Peurifoy, Yichen Shen, Li Jing, Yi Yang, Fidel Cano-Renteria, Brendan G. DeLacy, John D. Joannopoulos, Max Tegmark, and Marin Soljačić. Nanophotonic particle simulation and inverse design using artificial neural networks. *Science Advances*, 4(6), 2018.
- [36] Tomaso Poggio and Federico Girosi. Networks for approximation and learning. *Proceedings of the IEEE*, 78(9):1481–1497, 1990.
- [37] M. Raissi, P. Perdikaris, and G.E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686 – 707, 2019.
- [38] M. Raissi, P. Perdikaris, and G.E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [39] Maziar Raissi and George Em Karniadakis. Hidden physics models: Machine learning of nonlinear partial differential equations. *Journal of Computational Physics*, 357:125 – 141, 2018.
- [40] Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Machine learning of linear differential equations using gaussian processes. *Journal of Computational Physics*, 348:683 – 693, 2017.
- [41] Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*, 2017.
- [42] Russell Reed, Seho Oh, RJ Marks, et al. Regularization using jittered training data. In *International Joint Conference on Neural Networks*, volume 3, pages 147–152, 1992.
- [43] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and Peter Battaglia. Learning to simulate complex physics with graph networks. In *International Conference on Machine Learning*, pages 8459–8468. PMLR, 2020.
- [44] Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, and Peter Battaglia. Graph networks as learnable physics engines for inference and control. In *International Conference on Machine Learning*, pages 4470–4479. PMLR, 2018.
- [45] William E Schiesser. *The numerical method of lines: integration of partial differential equations*. Elsevier, 2012.
- [46] Samuel S. Schoenholz and Ekin D. Cubuk. Jax m.d. a framework for differentiable physics. In *Advances in Neural Information Processing Systems*, volume 33. Curran Associates, Inc., 2020.
- [47] Gordon D Smith, Gordon D Smith, and Gordon Dennis Smith Smith. *Numerical solution of partial differential equations: finite difference methods*. Oxford university press, 1985.
- [48] Sunae So, Trevon Badloe, Jaebum Noh, Jorge Bravo-Abad, and Junsuk Rho. Deep learning enabled inverse design in nanophotonics. *Nanophotonics*, 9(5):474, February 2020.
- [49] Mohammad H. Tahersima, Keisuke Kojima, Toshiaki Koike-Akino, Devsh Jha, Bingnan Wang, Chungwei Lin, and Kieran Parsons. Deep neural network inverse design of integrated photonic power splitters. *Scientific Reports*, 9(1):1368, Feb 2019.
- [50] Rohit K. Tripathy and Ilias Billionis. Deep uq: Learning deep neural network surrogate models for high dimensional uncertainty quantification. *Journal of Computational Physics*, 375:565 – 588, 2018.
- [51] Jeroen Tromp, Carl Tape, and Qinya Liu. Seismic tomography, adjoint methods, time reversal and banana-doughnut kernels. *Geophysical Journal International*, 160(1):195–216, 2005.
- [52] Kiwon Um, Robert Brand, Yun Raymond Fei, Philipp Holl, and Nils Thuerey. Solver-in-the-loop: Learning from differentiable physics to interact with iterative pde-solvers. *Advances in Neural Information Processing Systems*, 33:6111–6122, 2020.
- [53] Yi-Jen Wang and Chin-Teng Lin. Runge-kutta neural network for identification of dynamical systems in high accuracy. *IEEE transactions on neural networks*, 9 2:294–307, 1998.
- [54] Daniel A. White, William J. Arrighi, Jun Kudo, and Seth E. Watts. Multiscale topology optimization using neural network surrogate models. *Computer Methods in Applied Mechanics and Engineering*, 346:1118–1135, 2019.
- [55] Tailin Wu, Qinchen Wang, Yiman Zhang, Rex Ying, Kaidi Cao, Rok Sosič, Ridwan Jalali, Hassan Hamam, Marko Maucec, and Jure Leskovec. Learning large-scale subsurface simulations with a hybrid graph network simulator. *arXiv preprint arXiv:2206.07680*, 2022.

- [56] Yibo Yang and Paris Perdikaris. Adversarial uncertainty quantification in physics-informed neural networks. *Journal of Computational Physics*, 2019.
- [57] Jiawei Zhuang, Dmitrii Kochkov, Yohai Bar-Sinai, Michael P Brenner, and Stephan Hoyer. Learned discretizations for passive scalar advection in a two-dimensional turbulent flow. *Physical Review Fluids*, 6(6):064605, 2021.
- [58] Qinyu Zhuang, Juan Manuel Lorenzi, Hans-Joachim Bungartz, and Dirk Hartmann. Model order reduction based on runge–kutta neural networks. *Data-Centric Engineering*, 2, 2021.