FISEVIER

Contents lists available at ScienceDirect

Journal of Computational Physics

www.elsevier.com/locate/jcp



Meta-MgNet: Meta multigrid networks for solving parameterized partial differential equations



Yuyan Chen^a, Bin Dong^{b,*}, Jinchao Xu^c

- ^a School of Mathematical Science, Peking University, China
- ^b Beijing International Center for Mathematical Research & Institute for Artificial Intelligence, Peking University, China
- ^c Department of Mathematics, Pennsylvania State University, United States of America

ARTICLE INFO

Article history:

Received 1 November 2020 Received in revised form 28 December 2021 Accepted 15 January 2022 Available online 20 January 2022

Keywords:
Parameterized PDEs
Multigrid
Deep-learning
MgNet
Meta-learning

ABSTRACT

This paper studies numerical solutions for parameterized partial differential equations (PDEs) with deep learning. Parametrized PDEs arise in many important application areas, including design optimization, uncertainty analysis, optimal control, and inverse problems. The computational cost associated with these applications using traditional numerical schemes can be exorbitant, especially when the parameters fall into a particular range, and the underlying PDE model is required to be solved with high accuracy using a fine spatial-temporal mesh. Recently, solving PDEs with deep learning has become an emerging field in scientific computing. Existing works demonstrate great potentials of the deep learning-based approach in speeding up numerical solutions of various types of PDEs. However, there is still limited research on the deep learning approach for solving parameterized PDEs, If we directly apply existing deep supervised learning models to solving parameterized PDEs, the models need to be constantly fine-tuned or retrained when the parameters of the PDE change. This limits the applicability and utility of these models in practice. To resolve this issue, we propose a meta-learning based method that can efficiently solve parameterized PDEs with a wide range of parameters without retraining. Our key observation is to regard training a solver for the parameterized PDE with a given set of parameters as a learning task. Then, training a solver for the parameterized PDEs with varied parameters can be viewed as a multi-task learning problem, to which meta-learning is one of the most effective approaches. This new perspective can be applied to many existing PDE solvers to make them suitable for solving parameterized PDEs. As an example, we adopt the Multigrid Network (MgNet) [21] as the base solver. To achieve multi-task learning, we introduce a new hypernetwork, called Meta-NN, in MgNet and refer to the entire network as the Meta-MgNet. Meta-NN takes the differential operators and the right-hand-side of the underlying parameterized PDEs as inputs and generates appropriate smoothers for MgNet, which are essential ingredients for multigrid methods and can significantly affect the convergent speed. The proposed Meta-NN is carefully designed so that Meta-MgNet has guaranteed convergence for Poisson's equation. Finally, extensive numerical experiments demonstrate that Meta-MgNet is more efficient in solving parameterized PDEs than the MG methods and MgNet trained by supervised learning.

© 2022 Elsevier Inc. All rights reserved.

E-mail addresses: chenyuyan@pku.edu.cn (Y. Chen), dongbin@math.pku.edu.cn (B. Dong), xu@math.psu.edu (J. Xu).

^{*} Corresponding author.

1. Introduction

Partial differential equations are essential tools in many areas, such as physics, chemistry, biology, and economics. Most PDEs we encounter in practice contain parameters representing the system's physical or geometric properties, e.g., kinetic characteristics, material properties, the shape of the domain, etc. In practice, we often found ourselves in multi-query scenarios where the PDEs need to be solved for numerous different parameters with high accuracy and efficiency. Such scenarios include design optimization, optimal control, uncertainty quantification, inverse problems, etc. Therefore, a uniformly efficient solver for all parameters is urgently needed.

In this paper, we consider the following parameterized steady-state PDEs

$$\begin{cases} \mathcal{A}_{\eta} \overset{u}{\underset{\sim}{\sim}} = f, & \text{in } \Omega, \\ \overset{u}{\underset{\sim}{\sim}} = u_{b}, & \text{on } \partial \Omega, \end{cases}$$
 (1)

where $\Omega \subset \mathbb{R}^d$, $d, n \in \mathbb{N}_+$, \mathfrak{I} , \mathfrak{F} are two function spaces on Ω , and \mathfrak{U}_b is a function space on $\partial \Omega$, $\underline{\mu} = (u^1, u^2, ..., u^n) \in \mathfrak{U}^n$, $\underline{f} = (f^1, f^2, ..., f^n) \in \mathfrak{F}^n$, $\underline{\mu}_b = (u^1_b, u^2_b, ..., u^n_b) \in \mathfrak{U}^n_b$. And $\underline{\mathcal{A}}_{\eta} : \mathfrak{U}^n \to \mathfrak{F}^n$ is a linear differential operator with parameter $\eta = 0$ (η_1,\ldots,η_m) . For convenience, we will omit η when there is not any confusion. In this paper, the specific PDEs we choose for our numerical experiments are 2D/3D steady-state anisotropic diffusion equations and Ossen equations. The steady-state diffusion equations are widely used in fluid mechanics, electronic science, image processing, etc. The Ossen equations play an important role in fluid mechanics. We recall these PDEs as follows:

(1) 2D anisotropic diffusion equations:

$$\begin{cases} -\nabla \cdot (C\nabla u) = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial \Omega, \end{cases}$$

where
$$C = C(\epsilon, \theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$
 is a 2×2 matrix, $\epsilon < 1, \theta \in [0, \pi]$. (2) 3D anisotropic diffusion equations:

$$\begin{cases} -\nabla \cdot (C\nabla u) = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial \Omega, \end{cases}$$

on domain $\Omega = [0, 1]^3$. In this paper, we only concern the case with

$$C = \begin{pmatrix} \epsilon_0 & & \\ & \epsilon_1 & \\ & & \epsilon_2 \end{pmatrix}, \quad \epsilon_0, \epsilon_1, \epsilon_2 > 0.$$

(3) Ossen equations:

$$\begin{cases} -\mu \Delta \underbrace{u}_{} + (\underbrace{a}_{} \cdot \nabla) \underbrace{u}_{} + \nabla p = \underbrace{f}_{}, & \text{in } \Omega, \\ -\text{div} \underbrace{u}_{} = \underbrace{0}_{}, & \text{in } \Omega, \\ \underbrace{u}_{} = \underbrace{0}_{}, & \text{on } \partial \Omega. \end{cases}$$

where $\mu = \frac{1}{Re}$ and Re is Reynold number, and $\mathbf{g} = (a_{\mathbf{x}}, a_{\mathbf{y}})^{\top}$.

When discretized, equation (1) becomes a linear system of equations

$$\mathbf{A}_{\eta}\mathbf{u} = \mathbf{f}.\tag{2}$$

Linear system (2) is usually of a very large scale in practice, and iterative methods are often used to solve it. The multigrid (MG) method [49,50,19] is one of the most compelling classical methods. The computational complexity of the MG method for the Poisson's equation is only O(n), where n is the size of the matrix; as compared to another popular highperformance numerical method, the spectral method, whose complexity is $O(n\log n)$. However, the MG method still has trouble solving PDEs (1) with η falling into a specified range. Taking the 2D anisotropic diffusion equation as an example [52], the computational cost of the MG method grows rapidly with $\epsilon \to 0$ (see **Fig. 1**).

Therefore, people try to adjust the parameters and components in the MG method (such as smoothers, prolongations, and restrictions) to improve its performance according to η . Although decades of continuing researches are devoted to the convergence theory of the MG method to find the theoretically best parameters, the computational cost of finding such optimal parameters for a given η can be much higher than solving the linear system (2) itself. For example, [51]

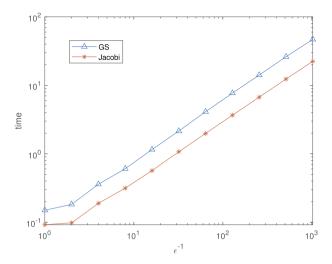


Fig. 1. The wall times of convergence grows rapidly with $\epsilon \to 0$.

derives the theoretically optimal prolongation for a given smoother. However, it requires solving an eigenvalue problem. Crucial parameters such as the damped coefficient of damped Jacobi smoother and the prolongations and restrictions of the algebraic MG method are mostly designed by human experts. However, these manually designed parameters can be rather complicated and have limited improvements for MG methods in practice.

Before the rise of deep learning, machine learning had a similar challenge as MG methods. Classical machine learning heavily relies on feature engineering, where people tried to manually design various types of features that are later fed into a classification or regression model. However, the quality of the features depends on the data set, the downstream task, and also the choice of the classification and regression model. It is extremely difficult to design fully adaptive and good feature extractors purely based on the human experience. This is where deep learning has been proven tremendously effective. After the advent of deep learning, feature extractors can be learned directly from data [28] in an end-to-end fashion. This often leads to feature extractors significantly surpass the previous ones designed by human experts in performance. This enables deep learning to achieve enormous success in many areas of artificial intelligence, such as natural language processing [23] and games [43].

The success of deep learning motivated us to resort to a data-driven approach to determine MG methods' parameters. Among all deep learning approaches, deep supervised learning has become one of the most popular data-driven approaches in scientific computing. Therefore, using deep supervised learning to determine the parameters of the MG method is a natural choice. Since classical MG methods do not work well for certain η , we can learn the parameters in the MG method to make it uniformly efficient for all η .

Deep supervised learning can be easily applied to the MG method if we cast the MG method into a similar formulation as deep neural networks. As first observed in [21], the MG method is an iterative method which can be viewed as a feedforward network. Furthermore, the prolongations, restrictions, and some special smoothers can all be written as convolutions. Therefore, the MG method has a natural connection with the convolutional neural network (CNN). With that, [21] introduced the Multigrid Network (MgNet).

The original MgNet in [21] was proposed for image classification. In this paper, we convert it into a form suitable to solve PDEs and refer to it as PDE-MgNet. PDE-MgNet takes the right-hand-side function f as input and the approximated solution f as output. PDE-MgNet performs very well when it is trained on a data set generated by f f in PDEs (1) and tested on some other f generated from the same distribution f. However, it may generalize poorly (i.e., convergence slows down significantly) beyond the training setting, i.e., for f with f different from f.

This problem is common for deep supervised learning. To resolve it, we need to adopt a more robust learning mechanism. In this paper, we regard learning a solver for PDEs (1) with a given η as one learning task. Then learning solvers for PDEs (1) for all η can be naturally viewed as a multi-task learning problem. In the area of artificial intelligence, meta-learning is an effective approach to solve multi-task learning problems. Thus we propose to use the meta-learning approach to improve PDE-MgNet. Note that PDE-MgNet is just an example we choose in this paper. A similar methodology can be applied to enhance other numerical solvers for parameterized PDEs.

First, let us briefly review meta-learning. Meta-learning [45,24], also known as learning to learn, is a science of studying how different machine learning approaches perform on a wide range of tasks and use these experiences to speed up the learning of new tasks. In contrast, supervised learning trains a model M_{η} for each task T_{η} separately. Even though all M_{η} have the same network structures, we need to retrain the model each time η changes. If there are numerous tasks to learn, the computation cost can be unbearable for supervised learning. The meta-learning approach resolves this problem

by gaining experience over multiple learning episodes - often covering the distribution of related tasks - and uses this experience to improve its future learning performance.

There are several strategies in meta-learning. We briefly introduce two strategies that are suitable for our task of interest. Suppose the only difference between M_n is the weights of neural networks, namely

$$M_{\eta} = M(\mathbf{x}; \mathbf{w}_{\eta}),$$

where \mathbf{x} is the input and \mathbf{w}_{η} is the weights of neural networks. The two strategies to quickly find suitable \mathbf{w}_{η} for each task η are given as follows.

(1) Finding a good initialization [13,37]: This strategy makes use of a series of tasks $T_{\eta_1}, T_{\eta_2}, ..., T_{\eta_n}$ to obtain a good initial weights \mathbf{w}_0 of the deep neural network M. Then, we can easily get \mathbf{w}_η for each task η by fine-tuning from \mathbf{w}_0 . For example, suppose the task T_η is sampled from distribution $\mathcal{P}(T)$, the loss of model M on the task T is $L_T(M)$, and we use gradient descent with learning rate α to train the model. After one step of gradient descent starting from \mathbf{w}_0 , we obtain the updated weights as $\mathbf{w}' = \mathbf{w}_0 - \alpha \nabla_{\mathbf{w}} L_{T_\eta}(M_\eta(\cdot; \mathbf{w}_0))$. Thus, the initial weights \mathbf{w}_0 should minimize the following expectation

$$\underset{\mathsf{T}_{\eta} \sim \mathcal{P}(\mathsf{T})}{\mathbb{E}} L_{\mathsf{T}_{\eta}}(\mathsf{M}_{\eta}(\cdot; \mathbf{w}')). \tag{3}$$

This is the main idea of Model-Agnostic Meta-Learning (MAML) proposed by [32], where an algorithm solving (3) is also proposed to estimate \mathbf{w}_0 .

(2) Hypernetwork [17,33,5,54]: This strategy relies on the designs of a network called hypernetwork to infer \mathbf{w}_{η} instead of direct learning of \mathbf{w}_{η} by training. The hypernetwork takes the information on the task η as input and \mathbf{w}_{η} as output. The hypernetwork are trained to make an accurate prediction on \mathbf{w}_{η} for each task η . When the hypernetworks are not powerful enough to make accurate predictions on \mathbf{w}_{η} , we treat the approximation of \mathbf{w}_{η} as a good initial value and resort to fine-tuning on each task η to improve the prediction. In particular, if the output of a hypernetwork is the same for all η , we can think that the hypernetwork gives a good initialization for all tasks T_{η} . In this regard, the previous strategy can be viewed as a special case of the hypernetwork approach.

By viewing solving the parameterized PDE (1) for a given η as a task T_{η} , we adopt the hypernetwork based meta-learning approach to improve upon the PDE-MgNet. With this, we introduce a new model called Meta-MgNet. Compared to PDE-MgNet, the Meta-MgNet uses a carefully designed hypernetwork to infer the model parameters of the MgNet, instead of learning it directly from data. The hypernetwork grants the Meta-MgNet great ability of in-distribution generalization and out-of-distribution transfer. We shall call this hypernetwork Meta-NN. The Meta-NN is used to infer parameters of specific components in the MgNet. In this paper, we select two types of smoother as an example and design the corresponding meta-smoother (i.e., using Meta-NN to infer parameters of the smoother) for the Meta-MgNet. The two types of smoother are the convolution smoother, which is exactly what MgNet uses, and the smoother based on subspace correction. For the convolution smoother, the Meta-NN infers its convolution kernels. For the smoother based on subspace correction, the Meta-NN infers the spanning vectors of the subspace. The parameters of the Meta-NN are first trained on a data set with mixed data from different η . Then, we can fine-tune the Meta-NN while solving a specific T_{η} . However, our numerical experiments show that the weights given by Meta-NN without retraining are good enough. Therefore, we shall skip the retraining step in the experiments.

1.1. Related work

In the area of solving PDEs by machine learning, especially deep learning, the existing algorithms can be divided into the following two categories.

- (1) Using neural networks to approximate the function μ : These algorithms are suitable for a PDE with a fixed differential operator $\mathcal A$ and the right-hand side f. The most notable advantage of this approach is that they can: 1) overcome the curse of dimensionality and solve high-dimension PDEs; 2) resolve complex geometries in the solution due to the meshless representation of neural networks. Successful examples include the nonlinear convection-diffusion PDEs [39,40], Riemann Problem [36], high-dimension PDEs [3,11,20,44], and others [47,53,46,31,42]. Nevertheless, if the parameter f0 or the right-hand-side function f1 is changed, the neural network often needs to be retrained.
- (2) Using neural networks to approximate the solution operator \mathcal{A}^{-1} : The general modeling strategy of methods in this category is to replace a portion of the classical numerical method with neural networks to improve its performance. For example, [41] uses NN to estimate locations of discontinuous, [9] uses NN to introduce an appropriate amount of artificial dissipation in the numerical solver. There are also works using NN to approximate the entire operator \mathcal{A}^{-1} . For example, [25] trains a U-Net as a solver for Poisson's equations. The most related work to the current one is [12], where the authors use a meta-learning approach for parameterized pseudo-differential operators with deep neural networks.

However, there are two main differences between their work and ours. Firstly, their idea is to find the map $\eta \mapsto \mathcal{A}_{\eta}^{-1}$ directly. Ours is based on the observation that learning solvers for PDE (1) with different η is a multi-task learning problem and then introduce meta-learning to solve it. Secondly, their approach is based on the wavelet method, while ours is based on the MgNet.

There are also several studies on improving the MG method by deep-learning. For instance, [27,15,35] proposed a series of data-driven approaches to design prolongations and restrictions in MG. However, these methods are based on supervised learning while we focus on learning smoothers based on meta-learning.

In addition to the approaches above, the reduced-order modeling (ROM) [1,4] is also a widely used method for solving parameterized PDEs. The objective of ROM is to generate reduced models that are cheaper to solve while still well approximate the original PDEs. For example, [29] proposes a novel framework for projecting dynamical systems onto nonlinear manifolds using minimum-residual formulations at the time-continuous and time-discrete levels; [14] proposes new deep learning-based nonlinear ROM.

The remaining sections are organized as follows. In section 2, we will introduce necessary notations for the rest of the paper and discuss how to represent discrete PDEs by convolutions. In section 3, we review the multigrid method and the PDE-MgNet. In section 4, we introduce the proposed Meta-MgNet and provide a preliminary convergence analysis of the algorithm. In section 5, we present the numerical experiments and comparisons among the classical MG methods, PDE-MgNet, and Meta-MgNet. Concluding remarks are given in section 6.

2. Convolutions and differential operators

The key to solving PDEs is to discretize the differential operators properly. The main goal of this section is to present a theorem that convolutions can express the most meaningful discretizations of differential operators. This theorem plays an essential role for us to rewrite traditional numerical solvers as CNNs. The MgNet introduced by [21] is an example, which is a reformulation from the MG method. Furthermore, we will introduce the definition of convolution operators and then show how to use convolutions to represent discrete forms of differential operators.

2.1. Convolution operators

In this paper, we only consider the convolution of two second-order tensors and the convolution of a fourth-order tensor and a third-order tensor. Consider two second order tensors $K = (K_{j,i})$, with $j, i \in \mathbb{Z}$ and $v = (v_{j,i})$, with $i \in \{0, 1, ..., I\}$, $j \in \{0, 1, ..., I\}$, i.e.

$$\mathsf{K} = \begin{pmatrix} \ddots & \vdots & & \ddots \\ & \mathsf{K}_{-1,-1} & \mathsf{K}_{-1,0} & \mathsf{K}_{-1,1} \\ \cdots & \mathsf{K}_{0,-1} & \mathsf{K}_{0,0} & \mathsf{K}_{0,1} & \cdots \\ & \mathsf{K}_{1,-1} & \mathsf{K}_{1,0} & \mathsf{K}_{1,1} \\ \vdots & \vdots & & \ddots \end{pmatrix} \quad \text{and} \quad \mathsf{v} = \begin{pmatrix} \mathsf{v}_{0,0} & \mathsf{v}_{0,1} & \cdots & \mathsf{v}_{0,I} \\ \mathsf{v}_{1,0} & \mathsf{v}_{1,1} & \cdots & \mathsf{v}_{1,I} \\ \vdots & \vdots & \ddots & \vdots \\ \mathsf{v}_{J,0} & \mathsf{v}_{J,1} & \cdots & \mathsf{v}_{J,I} \end{pmatrix}.$$

The convolution K * v is also a second-order tensor, and the values of its components are

$$(\mathsf{K} \star \mathsf{v})_{j,i} = \sum_{l,t \in \mathbb{Z}} \mathsf{K}_{J,t} \mathsf{v}_{j+J,i+t},$$

where $i \in \{0, 1, ..., I\}$, $j \in \{0, 1, ..., J\}$. If $i + \iota \notin \{0, 1, ..., I\}$ or $j + \jmath \notin \{0, 1, ..., J\}$, we set $v_{j+J, i+\iota} = 0$.

Consider a forth-order tensor K and a third-order tensor v, where $K = (K_{l,k,J,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $v = (v_{k,j,i})$, with $l,k \in \{1,2,...,S\}$, $l,i \in \mathbb{Z}$ and $l,i \in \{1,2,...,S\}$

(1) **Convolution**: The convolution $K \star v$ is a third order tensor and the value of its components are

$$(\mathsf{K} \star \mathsf{v})_{l,j,i} = \sum_{k=1}^{s} \sum_{J,\iota \in \mathbb{Z}} \mathsf{K}_{l,k,J,\iota} \mathsf{v}_{k,j+J,i+\iota}.$$

(2) **Convolution with strides**: Suppose I_s , $J_s \in \mathbb{N}^+$, and we use \star_{J_s,I_s} to express the convolution with stride= (J_s,I_s) and the components of $K \star_{I_s,I_s} v$ are

$$(K \star_{J_s,I_s} \mathsf{v})_{l,j,i} = \sum_{k=1}^{S} \sum_{j,i \in \mathbb{Z}} \mathsf{K}_{l,k,j,i} \mathsf{v}_{k,jJ_s+j,iI_s+i}.$$

If $I_s = J_s$, we can write \star_{I_s,I_s} briefly as \star_{I_s} .

(3) **Deconvolution**: Suppose I_s , $J_s \in \mathbb{N}^*$, and we use \star^{J_s,I_s} to express deconvolution with strides= (J_s,I_s) and the components of $K \star^{J_s,I_s} V$ are

$$(\mathsf{K} \star^{J_{s},I_{s}} \mathsf{v})_{l,jJ_{s}+j',iI_{s}+i'} = \sum_{k=1}^{S} \sum_{J,\iota \in \mathbb{Z}} \mathsf{K}_{l,k,JJ_{s}+j',\iota I_{s}+i'} \mathsf{v}_{\kappa,j+J,i+\iota}.$$

If $I_s = J_s$, we can write \star^{J_s,I_s} briefly as \star^{J_s} . We can also regard deconvolution as a convolution after upsampling. For example, we have

$$K \star (v \uparrow) = K \star^2 v$$
.

2.2. Representing discretized PDEs in convolutions

We discuss how to transform the \mathcal{A}_{η} to the convolution form. We first write components of (1) in the following matrix form

$$\begin{pmatrix} \mathcal{A}_{1,1} & \mathcal{A}_{1,2} & \cdots & \mathcal{A}_{1,n} \\ \mathcal{A}_{2,1} & \mathcal{A}_{2,2} & \cdots & \mathcal{A}_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{A}_{n,1} & \mathcal{A}_{n,2} & \cdots & \mathcal{A}_{n,n} \end{pmatrix} \begin{pmatrix} u^1 \\ u^2 \\ \vdots \\ u^n \end{pmatrix} = \begin{pmatrix} f^1 \\ f^2 \\ \vdots \\ f^n \end{pmatrix}, \tag{4}$$

where each linear differential operator $\mathcal{A}_{i,j}(i,j=1,2,...,n)$ is a component of \mathcal{A}_{η} . Thanks to the linear superposition property of each component in (4), we only need to consider each component separately, i.e., for given linear differential operator \mathcal{K} and a functional f in \mathfrak{V}' , find $v \in \mathfrak{V}$ to satisfy

$$\mathcal{K}\mathbf{v} = f. \tag{5}$$

Our goal is to find a kernel K, and tensors v and f to represent the discretization of (5) as

$$K \star v = f. \tag{6}$$

Suppose that the Galerkin method, e.g. the finite difference method (FDM) or finite element method (FEM), is used to discretize the PDEs (5). Now, we propose a general way to convert (5) into the convolution form (6). According to the Galerkin method, we first convert PDEs (5) to its weak form:

find
$$v \in \mathfrak{V}$$
, such that $\forall w \in \mathfrak{V}$, $K(v, w) = f(w)$. (7)

Then, we choose a finite dimensional subspace $\mathfrak{V}_h \subset \mathfrak{V}$ to discretize (7) and solve the projected problem:

find
$$v_h \in \mathfrak{V}_h$$
, such that $\forall w_h \in \mathfrak{V}_h$, $K(v_h, w_h) = f(w_h)$.

Let Φ be a set of basis of \mathfrak{V}_h and assume that it satisfies the following assumptions.

Assumption 2.1. Suppose Φ can be divided into S groups $\Phi_1,...,\Phi_S$, where $\Phi_k = \{\phi_{k,j,i}|i=0,1,2...,I_k,j=0,1,2,...,J_k\}$ and each Φ_k can be generated by translations along the grid-lines from a compact support function φ_k , i.e. $\exists \varphi_k \ \forall i \in \{0,1,2,...,I_k\}$, $j \in \{0,1,2,...,J_k\}$ such that

$$\phi_{k,j,i}(x, y) = \varphi_k(x - ih, y - jh).$$

Then, we have the following theorem stating that the PDE (5) discretized by FDM or FEM can be expressed in convolution form.

Theorem 1. If a discretized scheme in FDM or FEM satisfies Assumption 2.1, the discretized PDE Kv = f can be written in the form $K \star v = f$ with K and K given as follows

$$K = (K_{l,k,j,i}) = K(\varphi_k(x - ih, y - jh), \varphi_l(x, y)) \quad and \quad f = (f_{l,i,j}) = f(\varphi_{l,i,i}).$$
 (8)

Proof. With Assumption 2.1, we can write v_h as

$$v_h = \sum_{k=1}^{S} \sum_{j=0}^{J_k} \sum_{i=0}^{I_k} v_{k,j,i} \phi_{k,j,i}.$$

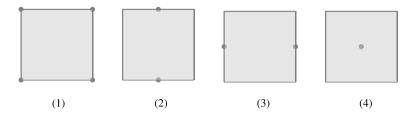


Fig. 2. Different case of point-wise discretion.

According to Galerkin method, we have

$$\begin{split} K(v_h,\phi_{l,j,i}) &= K(\sum_{k=1}^S \sum_{j=0}^{J_k} \sum_{\iota=0}^{I_k} \mathsf{v}_{k,j,\iota} \phi_{k,j,\iota}, \phi_{l,j,i}) = \sum_{k=1}^S \sum_{j=0}^{J_k} \sum_{\iota=0}^{I_k} \mathsf{v}_{k,j,\iota} K(\phi_{k,j,\iota},\phi_{l,j,i}) \\ &= \sum_{k=1}^S \sum_{j=0}^{J_k} \sum_{\iota=0}^{I_k} \mathsf{v}_{k,j,\iota} K(\varphi_k(x-\iota h,y-jh), \varphi_l(x-ih,y-jh)) \\ &= \sum_{k=1}^S \sum_{j=-j}^{J_k-j} \sum_{\iota=-i}^{I_k-i} \mathsf{v}_{k,j+j,\iota+i} K(\varphi_k(x-(\iota+i)h,y-(j+j)h), \varphi_l(x-ih,y-jh)) \\ &= \sum_{k=1}^S \sum_{j=-j}^{J_k-j} \sum_{\iota=-i}^{I_k-i} \mathsf{v}_{k,j+j,\iota+i} K(\varphi_k(x-\iota h,y-jh), \varphi_l(x,y)). \end{split}$$

Let $K = (K_{l,k,j,i})$, where

$$K_{l,k,j,j} = K(\varphi_k(x - ih, y - jh), \varphi_l(x, y)).$$
 (9)

While $i \notin \{0, 1, ..., I_k\}$ or $j \notin \{0, 1, ..., J_k\}$, we set $v_{k,j,i} = 0$, $K_{l,k,j,i} = 0$. Then, we have

$$K(v_h, \phi_{l,j,i}) = \sum_{k=1}^{S} \sum_{l=-i}^{J_k-j} \sum_{l=-i}^{I_k-i} \mathsf{v}_{k,j+j,\iota+i} \mathsf{K}_{l,k,j,\iota} = \sum_{k=1}^{S} \sum_{l \in \mathbb{Z}} \sum_{j \in \mathbb{Z}} \mathsf{v}_{k,j+j,\iota+i} \mathsf{K}_{l,k,j,\iota} = (\mathsf{K} \star \mathsf{v})_{l,j,i}.$$

Let $f = (f_{l,i,i})$ with $f_{l,i,i} = f(\phi_{l,i,i})$. We obtain

$$K(v_h, \phi_{l,j,i}) = f(\phi_{l,j,i}), \quad \forall l = 1, 2, ..., S,$$

which is the same as (6). \square

Although we have demonstrated a generic method to convert PDEs (5) into the convolution form, it is sometimes inconvenient to use. We can calculate K in an easier way for most discretization schemes, as will be described in the remaining part of this subsection. For both FDM and FEM discretization, we assume that the mesh \mathcal{T} is $N \times N$ uniform triangular or rectangular mesh, and let $h = \frac{1}{N}$

2.2.1. Finite difference methods (FDM)

FDM is one of the most popular discretizations used for solving PDEs by approximating them with difference equations. The basis functions of FDM are Legendre polynomials. Thus, it is not convenient to use (8) to calculate K and f. Fortunately, we can use Taylor's expansion to compute entries of K and f. Furthermore, the functions v and f can be easily discretized by using their values restricted on the mesh \mathcal{T} . As examples, we present four commonly seen cases as follows (see Fig. 2).

- (1) Vertex of an element: set $v_{i,j} = v(ih, jh)$, i, j = 0, 1, ..., N;
- (2) Midpoint of a horizontal edge: set $v_{j,i} = v((i-0.5)h, jh), i = 1, 2, ..., N, j = 0, 1, ..., N;$ (3) Midpoint of a vertical edge: set $v_{j,i} = v(ih, (j-0.5)h), i = 0, 1, ..., N, j = 1, 2, ..., N;$ (4) Center of an element, set $v_{j,i} = v((i-0.5)h, (j-0.5)h), i, j = 1, 2, ..., N.$

After choosing a discretization for v and f, we use Taylor expansion to discretize \mathcal{K} . Suppose $p \in \mathbb{N}$, $\alpha = (\alpha_1, \alpha_2)$, $|\alpha| = (\alpha_1,$ $\alpha_1 + \alpha_2, \partial^{\alpha} = \frac{\partial^{\alpha_1 + \alpha_2}}{(\partial x)^{\alpha_1} (\partial y)^{\alpha_2}}$ and $\mathcal{K} = \sum_{|\alpha| \leqslant p} a_{\alpha} \partial^{\alpha}$. Then, we can obtain the kernel K from a difference scheme. If $\forall \alpha$ such that $\alpha \le p$, we have a finite difference approximation of ∂^{α} as $\partial^{\alpha}_{h} v_{j,i} \approx \sum_{j,i} q^{\alpha}_{j,i} v_{j+j,i+i}$ and the expression of each

Table 1
Some common finite difference schemes and its corresponding kernels K.

K	Difference scheme	Kernel K
∂_X	$ \frac{1}{\hbar} (v_{i,j} - v_{i-1,j}) \\ \frac{1}{\hbar} (v_{i+1,j} - v_{i,j}) $	$ \frac{\frac{1}{h}\begin{pmatrix} -1 & 1 & 0 \end{pmatrix}}{\frac{1}{h}\begin{pmatrix} 0 & -1 & 1 \end{pmatrix}} $
∂_y	$\frac{1}{\hbar}(v_{i,j}-v_{i,j-1})$	$\frac{1}{h} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}$
	$\frac{1}{\hbar}(v_{i,j+1}-v_{i,j})$	$\frac{1}{h} \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}$
∂_{xx}	$\frac{1}{h^2}(v_{i-1,j} + v_{i+1,j} - 2v_{i,j})$	$\frac{1}{h^2}\begin{pmatrix} 1 & -2 & 1 \end{pmatrix}$
∂_{xy}	$\frac{1}{4h^2}(v_{i-1,j-1}+v_{i+1,j+1}-v_{i+1,j-1}-v_{i-1,j+1})$	$\frac{1}{4h^2} \left(\begin{array}{rrr} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{array} \right)$
∂_{yy}	$\frac{1}{\hbar^2}(v_{i,j-1} + v_{i,j+1} - 2v_{i,j})$	$\frac{1}{h^2} \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}$
Δ	$\frac{1}{h^2}(v_{i-1,j} + v_{i+1,j} + v_{i,j-1} + v_{i,j+1} - 4v_{i,j})$	$\frac{1}{h^2} \begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}$

Fig. 3. Some usual support of base function in FEM on rectangular mesh.

(3)

(4)

(2)

component of K is $K_{j,i} = \sum_{|\alpha| \leq p} \sum_{j,i} a_{\alpha} q_{j,i}^{\alpha}$. Common finite difference approximations of partial derivatives and their corresponding convolution kernels are listed in **Table 1**.

2.2.2. Finite element methods (FEM)

(1)

FEM is more complicated than FDM because a variable v usually contains several types of basis functions. Each type of basis functions is a channel of the tensor v. For most FEM methods on the rectangle mesh, we can divide each basis functions into the following 4 cases based on the support of the functions, i.e. 2×2 , 2×1 , 1×2 , and 1×1 elements (see **Fig. 3**). These four cases can be reduced to the case for FDM.

On the triangular mesh, each basis function can be divided into 6 cases, which are shown in **Fig. 4** according to the shape of the support of the function. For case $(1)\sim(4)$, they can also be reduced to the case of FDM, while for case (5) or (6), we may have to apply **Theorem 1** to calculate K and f.

3. The multigrid method and PDE-MgNet

In this section, we briefly describe the geometric MG method and PDE-MgNet. MG method is one of the most high-efficiency methods for solving PDEs. We consider the discrete form of the parameterized PDEs (2).

3.1. Multigrid method

Iterative method is one of the basic numerical methods for solving the linear system (2). Given an initial guess \mathbf{u}_0 and an update scheme represented by \mathcal{H} , we can write the iterative method generically as

$$\mathbf{u}_{t+1} = \mathbf{u}_t + \mathcal{H}(\mathbf{u}_t), \ t = 0, 1 \dots, T.$$
 (10)

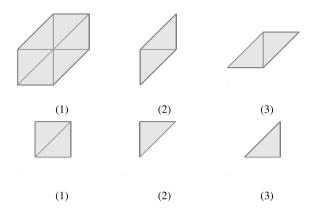


Fig. 4. Some usual support of base function in FEM on triangular mesh.

Note that we can regard (10) as a dynamical system or a feed-forward network. Such perspective is the key to connect deep neural networks (e.g. ResNet [22]) with dynamic systems, PDEs and optimal control [16,8,10,18,34,7].

In addition to the iterative scheme (10), the residual correction scheme

$$\mathbf{u}_{t+1} = \mathbf{u}_t + \mathcal{H}(\mathbf{f} - \mathbf{A}\mathbf{u}_t), \ t = 0, 1 \dots, T, \tag{11}$$

is one of most important types of iterative method. The MG method is one of such iterative methods, which is written as:

$$\mathbf{u}_{t+1} = \mathbf{u}_t + Mg(\mathbf{f} - \mathbf{A}\mathbf{u}_t), \ t = 0, 1..., T,$$
 (12)

where the Mg operator in (12) is given by **Algorithm 1**.

The MG operator can be divided into two steps, i.e. smoothing and coarse grid correction. The smoothing step is to use a smoother to eliminate high-frequency errors. A smoother \mathcal{B} (or its matrix form \mathbf{B}) is often an iterative scheme by itself taking the form

$$\mathbf{u}_{+} = \mathbf{u}_{0} + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}_{0}).$$

Popular choices of the smoother in MG include the Jacobi and Gauss-Seidel (GS) smoother, which are listed below

$$\mathbf{B} = \begin{cases} \operatorname{diag}(\mathbf{A})^{-1} & \operatorname{Jacobi} \\ \operatorname{tril}(\mathbf{A})^{-1} & \operatorname{GS} \end{cases}$$
 (13)

The Jacobi or GS smoother can efficiently eliminate high-frequency approximation errors. However, they are ineffective for low-frequency errors, which is where the coarse grid correction is needed. The MG methods utilize the solution on a coarse grid to approximate the low-frequency error. As a simple example, we illustrate the steps of coarse grid correction in a two-level MG in **Fig. 5**.

For the multi-level MG method, if we have a sequence of grids $\mathcal{T}^1(=\mathcal{T}), \mathcal{T}^2, ..., \mathcal{T}^J$ and assume that the prolongation and restriction between \mathcal{T}^ℓ and $\mathcal{T}^{\ell+1}$ are $\mathbf{P}^\ell_{\ell+1}$ and $\mathbf{R}^{\ell+1}_\ell$. Let $\mathbf{A}^{\ell+1} = \mathbf{R}^{\ell+1}_\ell \mathbf{A}^\ell \mathbf{P}^\ell_{\ell+1}$. The **Algorithm 1** presents the algorithm of multi-level MG.

3.2. PDE-MgNet

The original MgNet proposed by [21] is a new CNN model for image classification, which is inspired by the connection between CNNs and the MG method. As observed by [21] that the smoothers, prolongations, and restrictions can be represented by convolutions. Thus, we can write the corresponding kernels as B, P, and R. Therefore, the MG method can be naturally reformulated as a CNN model, which was called the MgNet by [21].

In this paper, we focus on solving PDEs. Thus, we need to modify the original MgNet in [21] to be suitable for solving PDEs. We shall call the modified MgNet the PDE-MgNet. We formulate every component of PDE-MgNet in the convolution form. For that, we consider the convolution form of PDEs (1):

$$A_n \star u = f. \tag{14}$$

PDE-MgNet replaces the smoother \mathcal{B}^{ℓ} , the prolongation \mathcal{P} , the restriction \mathcal{R} and \mathbf{A}^{ℓ} in the MG methods described in **Algorithm 1** with trainable convolution operators. **Figure 6** shows the architecture of a two-level \-Cycle PDE-MgNet and the multi-level case of PDE-MgNet is presented in **Algorithm 2**.

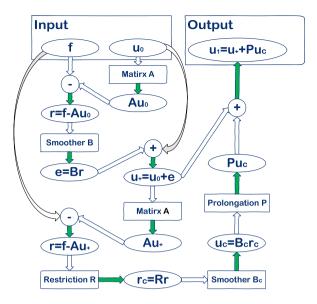


Fig. 5. Two level \-Cycle MG.

Algorithm 1 $\mathbf{u} = \mathrm{Mg}(\mathbf{f}; J, \nu_1, \cdots, \nu_I).$

```
Hyper-parameters: number of grids J, times of smooth in each grid: v_1, \dots, v_J
Input: right hand side f
Output: approximate solution {\bf u}
Initialization:
               \mathbf{f}^1 \leftarrow \mathbf{f}, \quad \mathbf{u}^{1,0} \leftarrow \mathbf{0}, \quad \mathbf{r}^{1,0} \leftarrow \mathbf{f}.
Smoothing and restriction from fine to coarse level (nested)
for \ell = 1 : J do
      Smoothing
      if \ell = J then
               \boldsymbol{u}^{\ell,1} \leftarrow (\boldsymbol{A}^{\ell})^{-1} \boldsymbol{r}^{\ell,0}
      else
            for i = 1 : \nu_{\ell} do
               \mathbf{u}^{\ell,i} \leftarrow \mathbf{u}^{\ell,i-1} + \mathbf{B}^{\ell} \mathbf{r}^{\ell,i-1}
               \mathbf{r}^{\ell,i} \leftarrow \mathbf{f}^{\ell} - \mathbf{A}^{\ell} \mathbf{u}^{\ell,i}.
            end for
      end if
      Form restricted residual
               \mathbf{f}^{\ell+1} \leftarrow \mathbf{R}_{\ell}^{\ell+1} \mathbf{r}^{\ell, \nu_{\ell}}, \quad \mathbf{u}^{\ell+1, 0} \leftarrow \mathbf{0}, \quad \mathbf{r}^{\ell+1, 0} \leftarrow \mathbf{f}^{\ell+1}.
end for
Prolongation and restriction from coarse to fine level
for \ell = J - 1:1 do
      Coarse grid correction
               \boldsymbol{u}^{\ell,\nu_\ell} \leftarrow \boldsymbol{u}^{\ell,\nu_\ell} + P^\ell_{\ell+1} \boldsymbol{u}^{\ell+1,\nu_\ell}.
end for
 return \ u = u^{1,\nu_1}.
```

With **Algorithm 2**, we obtain the iterative PDE-MgNet for solving (14):

$$u_{t+1} = u_t + PDE-MgNet(f - A \star u_t), \ t = 0, 1..., T,$$
 (16)

with $u_0 = 0$. For convenience, we shall refer to the iterative PDE-MgNet (16) simply as the PDE-MgNet. Note that the PDE-MgNet (16) is precisely the MG method (12) with \mathcal{A} , \mathcal{B} , \mathcal{P} and \mathcal{R} replaced by convolutions, and it reduces to the original MgNet proposed by [21] when the prolongation step is replaced by a classifier.

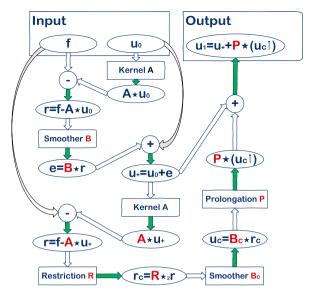


Fig. 6. Two level \-Cycle PDE-MgNet.

Algorithm 2 $u = PDE-MgNet(f; J, v_1, \dots, v_J).$

```
Hyper-parameters: number of grids J, times of smooth in each grid: v_1, \dots, v_I
Input: right-hand side f
Output: approximate solution u
Initialization
             f^1 \leftarrow f, u^{1,0} \leftarrow 0, r^{1,0} \leftarrow f.
Smoothing and restriction from fine to coarse level
for \ell = 1 : J do
     Smoothing:
     if \ell = J then
       Convert r^{\ell,0} into vector form \mathbf{r}^{\ell,0} and A^{\ell} into matrix form \mathbf{A}^{\ell}.
             \boldsymbol{u}^{\ell,1} \leftarrow (\boldsymbol{A}^{\ell})^{-1} \boldsymbol{r}^{\ell,0}.
       Convert \mathbf{u}^{\ell,1} into tensor form \mathbf{u}^{\ell,1}.
           for i = 1 : \nu_{\ell} do
                                                                                              \mathbf{u}^{\ell,i} \leftarrow \mathbf{u}^{\ell,i-1} + \mathbf{B}^{\ell,i-1} \star \mathbf{r}^{\ell,i-1}.
                                                                                                                                                                                                                                                               (15)
             \mathbf{r}^{\ell,i} \leftarrow \mathbf{f}^{\ell} - \mathbf{A}^{\ell} \star \mathbf{u}^{\ell,i}.
           end for
     end if
     Form restricted residual
             f^{\ell+1} \leftarrow R_\ell^{\ell+1} \star_2 r^{\ell,\nu_\ell}, \quad u^{\ell+1,0} \leftarrow 0, \quad f^{\ell+1,0} \leftarrow f^{\ell+1}.
end for
Prolongation from coarse to fine level
for \ell = J - 1:1 do
     Coarse grid correction
              u^{\ell,\nu_\ell} \leftarrow u^{\ell,\nu_\ell} + P^\ell_{\ell+1} \star^2 u^{\ell+1,\nu_\ell}.
end for
\textbf{return} \ u = u^{1,\nu_1}.
```

The values of the convolution kernels of PDE-MgNet are learned from data by minimizing a loss function defined similarly as in [30]. Suppose the distribution of the right-and-side function f is F. If we sample the distribution F by M_{train} times and denote the training data set as \mathfrak{X}_F , then the empirical loss is given by

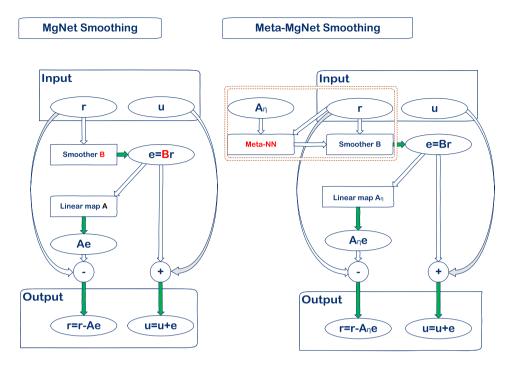


Fig. 7. Comparison between PDE-MgNet and Meta-MgNet in smoothing step. The red rectangle shows the major difference. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

$$Loss \approx \frac{1}{M_{\text{train}}} \sum_{f \in \mathcal{X}_F} \frac{||f - A \star u_T||^2}{||f||^2}.$$
 (17)

In our experiments, we choose T=1 following [27].

4. Meta-MgNet

The PDE-MgNet suffers from poor generalization when tested on a data set generated from the distribution far away from that of the training set, which significantly limits the practicality and utility of PDE-MgNet. This motivates us to improve PDE-MgNet with meta-learning by introducing a properly designed hypernetwork which infers specific components of the PDE-MgNet according to the parameters η of the parameterized PDE to achieve uniformly fast convergence. In this paper, the hypernetwork is introduced to make the smoother $\mathcal B$ in PDE-MgNet PDE-dependent. Now, we shall describe details of the design of such hypernetwork and the architecture of the entire Meta-MgNet.

4.1. Architecture of Meta-MgNet

The hypernetwork we introduce to the PDE-MgNet is called Meta-NN. The Meta-MgNet uses Meta-NN to infer an appropriate smoother (called meta-smoother) for each parameter η . The architecture of the meta-smoother in comparison with the smoother of the PDE-MgNet is presented in **Fig. 7**. The advantage of Meta-MgNet over PDE-MgNet is that the smoother of Meta-MgNet changes according to A_{η} and r, or we can write $\mathcal{B} = \mathcal{B}_{A_{\eta},r}$ which is realized by the Meta-NN. The entire architecture of the (iterative) Meta-MgNet solving (14) is given by

$$u_{t+1} = u_t + \text{Meta-MgNet}(f - A_{\eta} \star u_t, A_{\eta}), \tag{18}$$

where the Meta-MgNet (\cdot) is computed using **Algorithm 2** with (15) replaced by the meta-smoothing:

$$\mathbf{u}^{\ell,i} \leftarrow \mathbf{u}^{\ell,i-1} + \mathcal{B}_{\mathsf{A}_{\eta},\mathsf{r}}^{\ell,i-1}(\mathsf{r}^{\ell,i-1}).$$

In this paper, we consider two different methods to realize the meta-smoother $\mathcal{B}_{A_{\eta},r}$. The first one is based on the convolutional smoother, and its kernels are inferred from a convolutional hypernetwork. Thus, we call it the direct method, which is natural but mediocre. The second one is based on subspace correction smoother. Using subspace correction as a smoother is not as common as other smoothers such as Gauss-Seidel, but the numerical experiments show it performs better than direct methods.

4.1.1. Direct method

In the basic structure of PDE-MgNet, the smoothers are convolutions. Therefore, the direct method is to use Meta-NN to infer the value of these convolution kernels. We can use a vanilla DNN as the Meta-NN. **Algorithm 3** presents the details of this method. As for the structure of Meta-NN, we use a fully connect neural network with two hidden layers with 100

Algorithm 3 $\mathcal{B} = B_{d}(\mathsf{r}, \mathsf{A}_n; \mathcal{G}).$

Hyper-parameters: \mathcal{G} Inputs: \mathbf{r} , \mathbf{A}_{η} , Outputs: \mathcal{B}

1. Calculate subspace:

$$B \leftarrow \mathcal{G}(\frac{r}{||r||}, A_{\eta}).$$

2. Define the effect of ${\cal B}$ as

$$\mathcal{B}(\mathbf{r}) := \mathbf{B} \star \mathbf{r}.$$

return B.

neurons in each layer.

4.1.2. Subspace correction method

The subspace correction (SC) method is a classical numerical method for solving linear equations. The SC smoother has more flexible parameters than the convolutional smoother in PDE-MgNet. For a linear system $\mathbf{Au} = \mathbf{f}$, a subspace correction \mathcal{B} is determined by a subspace \mathbb{G} which is usually represented by the range of a matrix \mathbf{G} , i.e.

$$\mathbb{G} = \text{span}\{\mathbf{g}_1, \mathbf{g}_2, ..., \mathbf{g}_L\} = \text{range}(\mathbf{G}),$$

where $\mathbf{G} = (\mathbf{g}_1, ..., \mathbf{g}_L)$.

Notice that \mathbf{g}_i has the same dimension as \mathbf{f} and \mathbf{u} . Thus, we let Meta-NN export multi-channel tensors with the same shape as \mathbf{f} and then reshape each channel to form \mathbf{g}_i . Details are given by **Algorithm 4**. The architecture of Meta-NN \mathcal{G} in

Algorithm 4 $\mathcal{B} = B_{sc}(r, A_{\eta}; \mathcal{G}).$

Hyper-parameters: Meta-NN ${\cal G}$

Inputs: r, A_{η}

Output: \mathcal{B}

1. Calculate subspace:

$$G \leftarrow \mathcal{G}(r, A_n),$$

where G is a tensor with shape $L \times K \times J \times I$.

- 2. Reshape the tensor G to $L \times KJI$ matrix, and write its transpose as G, which is a $KJI \times L$ matrix.
- 3. Do subspace correction with the subspace $\mathbb{G} = \text{range}(\mathbf{G})$:

$$\mathbf{S}_n \leftarrow \mathbf{A}_n \mathbf{G}$$
.

$$\mathbf{e} = \mathbf{G}(\mathbf{G}^{\top}\mathbf{S}_n)^{-1}\mathbf{G}^{\top}\mathbf{r}.$$

4. Define the effect of ${\cal B}$ as

$$\mathcal{B}(\mathsf{r}, \mathbf{A}_n; \mathcal{G}) := \mathsf{Reshape}(\mathbf{e}).$$

where Reshape (\cdot) means to reshape ${\bf e}$ to the same shape as tensor r. return ${\cal B}.$

Algorithm 4 needs a more careful design than the one in **Algorithm 3**. To select an appropriate subspace for traditional SC method is also difficult. The most popular choice is the Krylov subspace, i.e. $\mathbb{G}_K = \{\mathbf{r}, \mathbf{Ar}, ..., \mathbf{A}^k \mathbf{r}\}$. If we write $\mathbb{G}_K = \{\mathbf{f}_0(\mathbf{A})\mathbf{r}, \mathbf{f}_1(\mathbf{A})\mathbf{r}, ..., \mathbf{f}_k(\mathbf{A})\mathbf{r}\}$, and $\mathbf{f}_i(\mathbf{A}) = \mathbf{A}^i$. Inspired by such formulation, we design the Meta-NN in **Algorithm 4** as

$$\mathcal{G}_{\theta}(\mathbf{r}, \mathbf{A}_{\eta}) = \mathcal{N}_{\mathsf{FC}_{\theta}(\mathbf{A}_{\eta})}(\mathbf{r}). \tag{19}$$

Here, \mathcal{N}_{γ} is a CNN used to convert r into a multi-channel tensor, and each channel of the output plays the role as a $\mathbf{f}_{l}(\mathbf{A})$ in \mathbb{G}_{K} . The weights γ of \mathcal{N}_{γ} are the output of the neural network FC_{θ} . In this paper, the CNN \mathcal{N}_{γ} is a 3-layer Dense-Net block [26] and FC_{θ} is a 2-layer fully connected neural network. For the Dense-Net block, we use l to represent the channel number. An illustration of this Meta-NN is given in **Fig. 8**.

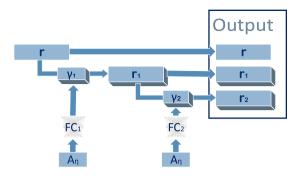


Fig. 8. The architecture of Meta-NN with 3 output channels as an example.

4.2. Training of Meta-MgNet

Suppose the distribution of η is Z and for each given η the distribution of f is F_{η} . Let the number of samples of η be $M_{\rm p}$. For each η , we sample f for $M_{\rm m-train}$ times and generate the data set $\mathfrak{X}_{F,\eta}$. Similar as the loss (17), we consider the following loss function

$$Loss = E_{\eta \sim Z, f \sim F_{\eta}} \frac{||f - A_{\eta} \star u_{T}||^{2}}{||f||^{2}} \approx \frac{1}{M_{p} M_{\text{m-train}}} \sum_{\eta \in \mathfrak{I}} \sum_{f \in \mathfrak{X}_{F, \eta}} \frac{||f - A_{\eta} \star u_{T}||^{2}}{||f||^{2}}.$$

In our experiments, we choose T = 1 following [27].

Remark 1. We can fine-tune the trained model when we are given a PDE with a new η , just like what meta-learning usually does. However, as shown in **Appendix I**, it brings little benefit and thus we shall omit fine-tuning.

4.3. Convergence analysis

This section analyzes the convergence of the proposed Meta-MgNet (16) with SC for Poisson's equation. We assume that the discretization schemes are either FDM with 7-point (or 9-point) stencil or FEM with P1 or Q1 elements. Now suppose the approximation of \mathbf{u} is \mathbf{u}_t , then the two-grid MG iteration can be written as

$$\begin{aligned} & \mathbf{r}_{t+\frac{1}{2}} = \mathbf{f} - \mathbf{A}\mathbf{u}_{t} = \mathbf{A}(\mathbf{u} - \mathbf{u}_{t}) \\ & \mathbf{u}_{t+\frac{1}{2}} = \mathbf{u}_{t} + \mathbf{P}\mathbf{A}_{c}^{-1}\mathbf{P}^{\top}\mathbf{r}_{t+\frac{1}{2}} \\ & \mathbf{r}_{t+1} = \mathbf{f} - \mathbf{A}\mathbf{u}_{t+\frac{1}{2}} = \mathbf{A}(\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \end{aligned}$$

$$\mathbf{u}_{t+1} = \mathbf{u}_{t+\frac{1}{2}} + \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{r}_{t+1}.$$

Then we have the recurrence relation of the error from t to t + 1:

$$||u-u_{t+1}||_A = ||(I-B_{r_{t+1}}A)(u-u_{t+\frac{1}{2}})||_A = ||(I-B_{r_{t+1}}A)(I-C)(u-u_t)||_A,$$

where $\mathbf{C} = \mathbf{P} \mathbf{A}_c^{-1} \mathbf{P}^{\top} \mathbf{A}$.

By [19], if the prolongation P is given by 7-point stencil or 9-point stencil, it is obviously that

$$||(\mathbf{I} - \mathbf{C})\mathbf{v}||_{\mathbf{A}}^2 < ||\mathbf{v}||_{\mathbf{A}}^2, \ \forall \mathbf{v} \in \mathbb{V}_h, \tag{20}$$

and there is a constant $c_a > 0$ s.t. $\forall \mathbf{v} \in \mathbb{V}_h$,

$$||(\mathbf{I} - \mathbf{C})\mathbf{v}||_{\mathbf{A}}^2 \le c_a \rho_{\mathbf{A}}^{-1} ||\mathbf{v}||_{\mathbf{A}^2}^2, \tag{21}$$

where $\rho_{\mathbf{A}}$ is the spectral radius of **A**.

Now, we introduce the following assumptions (see also [48,49]).

Assumption 4.1. For a given \mathbf{r} , we assume that the associated $\mathbf{B}_{\mathbf{r}}$ satisfies:

1. $\mathbf{B_r}$ is semi-symmetric positive defined (SSPD) and $\mathbf{B_r}\mathbf{AB_r} = \mathbf{B_r}$,

2. There is a constant $c_s > 0$ independent with \mathbf{r} , s.t.

$$||\mathbf{r}||^2 \le c_s \rho_{\mathbf{A}} \mathbf{r}^{\mathsf{T}} \mathbf{B}_{\mathbf{r}} \mathbf{r}. \tag{22}$$

Then, we have the following convergence theorem for Meta-MgNet (see also [2]).

Theorem 2. Let $\{\mathbf{u}_t\}$ be the sequence generated by the Meta-MgNet (16). Then, we have the convergence estimation

$$||\mathbf{u} - \mathbf{u}_t||_{\mathbf{A}} \le \delta^{\frac{t}{2}}||\mathbf{u} - \mathbf{u}_0||_{\mathbf{A}}, \quad with \ \delta = 1 - \frac{1}{c_a c_s}.$$

Proof. It suffices to show that

$$||\mathbf{u} - \mathbf{u}_{t+1}||_{\mathbf{A}} \le \delta^{\frac{1}{2}}||\mathbf{u} - \mathbf{u}_{t}||_{\mathbf{A}}. \tag{23}$$

If B_r satisfies **Assumption 4.1**, we have

$$\begin{split} ||\mathbf{u} - \mathbf{u}_{t+1}||_{\mathbf{A}}^2 &= ||(\mathbf{I} - \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A}) (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})||_{\mathbf{A}}^2 \\ &= (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})^\top (\mathbf{I} - \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}}) \mathbf{A} (\mathbf{I} - \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A}) (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \quad (\text{by } \mathbf{B}_{\mathbf{r}_{t+1}} \text{ is SSPD}) \\ &= (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})^\top (\mathbf{A} - 2 \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A} + \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A}) (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \\ &= (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})^\top (\mathbf{A} - \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A}) (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \quad (\text{by } \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} = \mathbf{B}_{\mathbf{r}_{t+1}}) \\ &= ||\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}||_{\mathbf{A}}^2 - (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})^\top \mathbf{A} \mathbf{B}_{\mathbf{r}_{t+1}} \mathbf{A} (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \\ &= ||\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}||_{\mathbf{A}}^2 - \mathbf{r}_{\mathbf{r}}^{-1} \mathbf{p}_{\mathbf{A}}^{-1} \mathbf{r}_{t+1}^{-1} \mathbf{r}_{t+1} \quad (\text{by } (22)) \\ &= ||\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}||_{\mathbf{A}}^2 - \mathbf{c}_{\mathbf{s}}^{-1} \rho_{\mathbf{A}}^{-1} (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})^\top \mathbf{A}^2 (\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}) \\ &= ||\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}}||_{\mathbf{A}}^2 - \mathbf{c}_{\mathbf{s}}^{-1} \rho_{\mathbf{A}}^{-1} ||(\mathbf{u} - \mathbf{u}_{t+\frac{1}{2}})||_{\mathbf{A}^2}^2 \end{aligned}$$

Write $\mathbf{w}_t = \mathbf{u} - \mathbf{u}_{t+\frac{1}{2}} = (\mathbf{I} - \mathbf{C})(\mathbf{u} - \mathbf{u}_t)$, with (20) and (21), we have

$$||\mathbf{u}-\mathbf{u}_{t+1}||_{\mathbf{A}}^2 \leq ||\mathbf{w}||_{\mathbf{A}}^2 - c_s^{-1}\rho_{\mathbf{A}}^{-1}||\mathbf{w}||_{\mathbf{A}^2} \leq ||\mathbf{u}-\mathbf{u}_t||_{\mathbf{A}}^2 - c_s^{-1}c_a^{-1}||\mathbf{u}-\mathbf{u}_t||_{\mathbf{A}}^2 = \delta||\mathbf{u}-\mathbf{u}_t||_{\mathbf{A}}^2$$

Thus, (23) is derived.

Now, we only need to verify that G = G(r) given in **Algorithm 4** stratifies **Assumption 4.1**. Indeed, we have

$$\mathbf{B}_{\mathbf{r}} = \mathbf{G}(\mathbf{G}^{\top}\mathbf{A}\mathbf{G})^{-1}\mathbf{G}^{\top}.$$

For the first assumption,

A is SPD
$$\Rightarrow$$
 $\mathbf{G}^{\top}\mathbf{A}\mathbf{G}$ is SPD \Rightarrow $(\mathbf{G}^{\top}\mathbf{A}\mathbf{G})^{-1}$ is SPD \Rightarrow $\mathbf{B_r} = \mathbf{G}(\mathbf{G}^{\top}\mathbf{A}\mathbf{G})^{-1}\mathbf{G}^{\top}$ is SSPD.

And

$$B_rAB_r = G(G^\top AG)^{-1}G^\top AG(G^\top AG)^{-1}G^\top = G(G^\top AG)^{-1}G^\top = B_r.$$

For the second assumption, without loss of generality, we assume $||\mathbf{r}|| = 1$, then we have $||\mathbf{r}||_{\mathbf{A}}^2 \le \rho_{\mathbf{A}} ||\mathbf{r}||^2 = \rho_{\mathbf{A}}$. Write $\mathbf{G} = [\mathbf{g}_1, ..., \mathbf{g}_L]$, and $\mathbf{g}_i, i = 1, 2, ..., L$ satisfy $\mathbf{g}_i^{\mathsf{T}} \mathbf{A} \mathbf{g}_j = \delta_{ij}$. Choose $\mathbf{g}_1 = \frac{\mathbf{r}}{||\mathbf{r}||_{\mathbf{A}}}$. Then, we obtain

$$\boldsymbol{r}^{\top}\boldsymbol{B}_{\boldsymbol{r}}\boldsymbol{r} \geq (\boldsymbol{r}^{\top}\boldsymbol{g}_{1})^{2} = \frac{1}{||\boldsymbol{r}||_{\boldsymbol{A}}^{2}} \geq \frac{1}{\rho_{\boldsymbol{A}}}.$$

This means that (22) is satisfied if $\mathbf{r} \in \text{range}(\mathbf{G})$, which is obvious from the design of the Meta-NN for B_{sc} .

5. Numerical experiments

In this section, we evaluate the performance of Meta-MgNet through a series of numerical experiments. In Section 5.1 and **5.2**, we apply Meta-MgNet to 2D and 3D anisotropic diffusion equations on domain $\Omega = [0, 1]^d$, d = 2, 3. In **Section 5.3**, we demonstrate why it is challenging to train a PDE-MgNet that generalizes well for all η . In **Section 5.4**, we include more

In the experiments, we compare Meta-MgNet with PDE-MgNet and MG method. Both Meta-MgNet and PDE-MgNet are trained on a data set created by a set of η , which will be later described in detail. Furthermore, we also train PDE-MgNet for each individual η and denote the trained model as PDE-MgNet- η . During testing, we only apply PDE-MgNet- η on the test data generated by the same η as during training. Thus, PDE-MgNet- η demonstrates the best accuracy of PDE-MgNet may achieve while it is impractical since it requires retraining of PDE-MgNet for each individual η . For the classical MG method, we choose the Krylov subspace smoother, GS smoother, line-Jacobi smoother, and damped Jacobi smoother. Since the GS and line-Jacobi smoother are challenging to implement efficiently with GPU, we use Matlab on CPU instead. All other algorithms are implemented in PyTorch on GPU.

5.1. 2D anisotropic diffusion equations

We consider the anisotropic diffusion equation

$$\begin{cases} -\nabla \cdot (C\nabla u) = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial \Omega, \end{cases}$$

$$\text{where } C = C(\epsilon, \theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \text{ is a } 2 \times 2 \text{ matrix, } \epsilon < 1, \theta \in [0, \pi]. \end{cases}$$

where
$$C = C(\epsilon, \theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$
 is a 2 × 2 matrix, $\epsilon < 1, \theta \in [0, \pi]$

5.1.1. Settings: training, testing and hyperparameter selection

- (1) For the training set, we randomly sample $M_p = 20$ different sets of parameters η from an interval I_{train} . The interval I_{train} is different for different PDEs. Thus, we will specify it later for each experiment. For each given η , we randomly sample $M_{\text{m-train}} = 100$ right-hand-side function, and each entry of f is sampled from the Gaussian distribution N(0, 1). We use ADAM method and the unsupervised learning loss (17) to train PDE-MgNet for 50 epochs and Meta-MgNet for 20 epochs. Since no matter if Meta-MgNet is fine-tuned, the result is almost the same, we skip the fine-tuning stage of Meta-MgNet. (See Appendix I for an ablation study on fine-tuning.) The learning rate for ADAM is 0.02 and the batch
- (2) For the test set, we choose some specific η from a set I_{test} , and for each η we randomly sample $M_{\text{m-test}} = 10$ righthand-side function. The stopping criterion for all compared algorithms is chosen as

$$\frac{||\mathbf{f} - \mathbf{A}_{\eta} \mathbf{u}_t||_2}{||\mathbf{f}||_2} < 10^{-6}.$$

Number of iterations and wall time are used as metrics to compare the performance of different algorithms. Furthermore, we use "mean±std" to show the average and standard deviation of the number of iteration and wall time over the $M_{\text{m-test}}$ samples.

(3) We use $N \times N$ rectangular mesh and Q_1 -element to discretize the PDEs. We select N = 256 and the number of layers J=5. The prolongations $\mathcal P$ and restrictions $\mathcal R$ are given by the traditional 9-point stencil, i.e. the kernel of $\mathcal P$ and $\mathcal R$ is

$$P = R = \begin{pmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{2} & 1 & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \end{pmatrix}.$$

The \mathcal{A}^{ℓ} on coarse grid is also given by the geometric multigrid method, and it is easy to verify that each A^{ℓ} is equal for Q_1 -element. We chose \-Cycle structure for all algorithms with $\nu_1 = 2, \nu_2 = \cdots = \nu_J = 1$. (We tried several different settings of v_I and found that $v_1 = 2$, $v_2 = \cdots = v_I = 1$ is the best option. See **Appendix II** for the corresponding ablation

(4) For PDE-MgNet, the kernel size of smoother at each layer is 7×7 . For Meta-MgNet, we choose SC method B_{SC} in **Algorithm 4** and use the Meta-NN in (19). The kernel size of the output is 7×7 and the increase channel number of Meta-NN is l=3. The size of the hidden layer of the fully connected neural network FC $_{\theta}$ is 200. The number of the parameters of one single Meta-NN is 122.2K. The Meta-MgNet in this numerical experiment contains 5 Meta-NNs, which means the total number of the parameters of Meta-MgNet is around 611K. Note that a comparison between B_d and B_{SC} is given in **Appendix III**, it shows that B_{SC} is the better option. Thus, we only present the numerical result of $B_{\rm sc}$ in the main body of this paper.

Table 2The number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method. "-" means the algorithm does not converge within 10⁴ iterations.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η	MG (Krylov)	MG (GS)	MG (line-Jacobi)	MG (Jacobi)
$\epsilon = 1$	4.0 ± 0.00	_	7.0 ± 0.00	4.0 ± 0.00	10.0 ± 0.00	14.0 ± 0.00	15.0 ± 0.00
$\epsilon = 10^{-1}$	7.5 ± 0.50	19.2 ± 0.40	21.2 ± 0.60	$\boldsymbol{7.9 \pm 0.30}$	33.7 ± 0.48	13.0 ± 0.00	90.2 ± 0.98
$\epsilon = 10^{-2}$	35.1 ± 1.04	178.9 ± 2.74	149.7 ± 3.44	52.5 ± 0.81	253.6 ± 4.19	13.0 ± 0.00	$\textbf{752.8} \pm \textbf{12.23}$
$\epsilon = 10^{-3}$	171.6 ± 6.34	$1.2e3 \pm 12.85$	910.9 ± 15.64	345.9 ± 3.88	$1.9e3 \pm 25.56$	13.0 ± 0.00	$5.6e3 \pm 119.42$
$\epsilon = 10^{-4}$	375.2 ± 5.88	_	$3.1e3 \pm 35.70$	$2.2\text{e}3 \pm 27.94$	_	11.0 ± 0.00	_
$\epsilon=10^{-5}$	$\textbf{797.8} \pm \textbf{12.76}$	-	$9.9e3 \pm 40.81$	$7.6\text{e}3 \pm 81.96$	-	11.0 ± 0.00	-
wall time							
$\epsilon = 1$	$}}}$	_	$\textbf{0.02} \pm \textbf{0.00}$	$\textbf{0.02} \pm \textbf{0.00}$	0.14 ± 0.01	0.34 ± 0.01	0.04 ± 0.00
$\epsilon = 10^{-1}$	$\boldsymbol{0.05 \pm 0.00}$	$\boldsymbol{0.05 \pm 0.00}$	$\boldsymbol{0.06 \pm 0.00}$	$\textbf{0.04} \pm \textbf{0.00}$	$\boldsymbol{0.48 \pm 0.02}$	$\boldsymbol{0.32 \pm 0.01}$	0.23 ± 0.00
$\epsilon = 10^{-2}$	$\textbf{0.22} \pm \textbf{0.01}$	$\boldsymbol{0.44 \pm 0.01}$	$\boldsymbol{0.37 \pm 0.01}$	$\boldsymbol{0.25 \pm 0.00}$	$\boldsymbol{3.47 \pm 0.15}$	$\boldsymbol{0.32 \pm 0.01}$	$\boldsymbol{1.85 \pm 0.03}$
$\epsilon=10^{-3}$	1.06 ± 0.04	3.04 ± 0.03	2.28 ± 0.05	$\boldsymbol{1.64 \pm 0.02}$	27.33 ± 0.68	$\textbf{0.32} \pm \textbf{0.01}$	13.95 ± 0.29
$\epsilon = 10^{-4}$	2.31 ± 0.03	_	$\boldsymbol{7.69 \pm 0.13}$	10.56 ± 0.14	_	$\textbf{0.27} \pm \textbf{0.01}$	_
$\epsilon = 10^{-5}$	4.91 ± 0.08	-	24.49 ± 0.14	35.67 ± 0.40	_	$\textbf{0.27} \pm \textbf{0.02}$	_

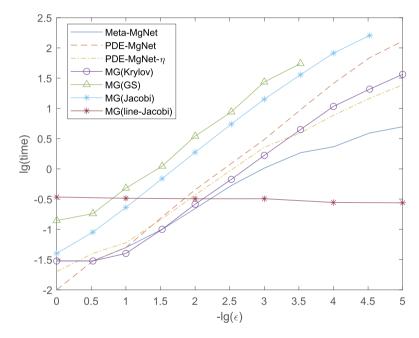


Fig. 9. Wall time of Meta-MgNet, PDE-MgNet, and the MG methods while $\theta=0$ and ϵ varies.

We compare number of iterations and wall time of Meta-MgNet, PDE-MgNet and MG method. We consider two different generalization scenarios that will be called "in-distribution generalization" and "out-of-distribution (OoD) transfer". For indistribution generalization, we have $I_{\text{test}} \subset I_{\text{train}}$, while for OoD transfer, the parameter η of the test data is entirely outside the interval of the training data, i.e. $I_{\text{test}} \subset I_{\text{train}}^c$.

5.1.2. In-distribution generalization

In this group of experiments, the training set of PDE-MgNet and Meta-MgNet is generated by fixing $\theta=0$ in **Table 2** and $\theta=0.1\pi$ in **Table 3** and randomly sampling ϵ with distribution $\lg\frac{1}{\epsilon}\sim U[0,5]$. The **Table 2**, **Table 3** and **Fig. 9** show Meta-MgNet has overall better performance than PDE-MgNet and MG methods, while the advantage is significant when ϵ is small. It is worth mentioning that the line-Jacobi smoother can only be applied to several specific θ , such as $0, \frac{\pi}{4}, \frac{\pi}{2}$, thus line-Jacobi smoother is limited in practical applications.

5.1.3. Out-of-distribution (OoD) transfer

The first group of experiments are with the fixed $\theta = 0$. The training set is generated with $\lg \frac{1}{\epsilon} \sim U[2,3]$ and the test set is generated with $\epsilon = 1, 10^{-1}, 10^{-4}, 10^{-5}$ (i.e. a test for OoD transfer with respect to ϵ). For the second group of

Table 3The number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method. "-" means the algorithm does not converge within 10⁴ iterations.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η	MG (Krylov)	MG (GS)	MG (line-Jacobi)	MG (Jacobi)
$\epsilon = 1$	4.0 ± 0.00	-	17.00 ± 0.00	4.0 ± 0.00	10.0 ± 0.00	-	15.0 ± 0.00
$\epsilon = 10^{-1}$	5.4 ± 0.49	136.7 ± 1.95	68.30 ± 2.69	8.0 ± 0.30	27.3 ± 0.48	_	64.70 ± 0.64
$\epsilon = 10^{-2}$	28.5 ± 0.50	$1.0e3 \pm 25.39$	861.40 ± 24.53	45.4 ± 0.49	187.0 ± 2.53	_	476.30 ± 4.78
$\epsilon = 10^{-3}$	94.2 ± 1.33	$3.8e3\pm183.22$	$1.9\text{e}3 \pm 60.41$	142.8 ± 1.89	707.6 ± 15.24	-	$1.8e3 \pm 44.47$
$\epsilon = 10^{-4}$	129.4 ± 2.42	$5.3e3 \pm 180.33$	$3.8e3 \pm 94.79$	187.8 ± 3.79	990.2 ± 19.03	-	$2.5\text{e}3 \pm 62.01$
$\epsilon=10^{-5}$	134.8 ± 2.86	$5.6e3 \pm 168.51$	$4.1\text{e}3\pm110.96$	195.9 ± 2.43	$1.0\text{e}3 \pm 37.45$	-	$2.6\text{e}3 \pm 79.24$
wall time							
$\epsilon = 1$	0.04 ± 0.00	-	0.05 ± 0.00	$\textbf{0.03} \pm \textbf{0.00}$	0.14 ± 0.01	-	0.04 ± 0.00
$\epsilon = 10^{-1}$	$\textbf{0.05} \pm \textbf{0.00}$	$\boldsymbol{0.36 \pm 0.01}$	$\boldsymbol{0.18 \pm 0.01}$	$\textbf{0.05} \pm \textbf{0.00}$	$\boldsymbol{0.29 \pm 0.02}$	_	$\boldsymbol{0.17 \pm 0.00}$
$\epsilon = 10^{-2}$	$\textbf{0.21} \pm \textbf{0.00}$	2.64 ± 0.06	2.23 ± 0.06	$\boldsymbol{0.24 \pm 0.00}$	2.05 ± 0.08	_	$\boldsymbol{1.26\pm0.01}$
$\epsilon=10^{-3}$	$\textbf{0.69} \pm \textbf{0.01}$	$\boldsymbol{9.89 \pm 0.53}$	4.90 ± 0.17	$\boldsymbol{0.75 \pm 0.01}$	$\boldsymbol{7.78 \pm 0.21}$	-	4.79 ± 0.11
$\epsilon = 10^{-4}$	$\textbf{0.94} \pm \textbf{0.02}$	13.87 ± 0.48	$\boldsymbol{9.90 \pm 0.24}$	$\boldsymbol{0.98 \pm 0.02}$	10.76 ± 0.25	-	6.66 ± 0.15
$\epsilon = 10^{-5}$	$\textbf{0.98} \pm \textbf{0.02}$	14.56 ± 0.44	10.69 ± 0.30	1.02 ± 0.01	11.36 ± 0.51	_	$\boldsymbol{7.03 \pm 0.21}$

Table 4The mean and std of the number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method on the testing set, "-" means the algorithm does not converge within 10⁴ iterations.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η	MG (Krylov)	MG (GS)	MG (line-Jacobi)	MG (Jacobi)
$\epsilon = 1$	7.0 ± 0.00	_	7.0 ± 0.00	4.0 ± 0.00	10.0 ± 0.00	14.0 ± 0.00	15.0 ± 0.00
$\epsilon = 10^{-1}$	10.0 ± 0.00	23.0 ± 0.00	21.2 ± 0.60	$\boldsymbol{7.9 \pm 0.30}$	33.7 ± 0.48	13.0 ± 0.00	90.2 ± 0.98
$\epsilon = 10^{-4}$	340.7 ± 3.52	$5.8e3 \pm 121.90$	$3.1e3 \pm 35.70$	$2.2\text{e}3 \pm 27.94$	_	11.0 ± 0.00	_
$\epsilon=10^{-5}$	817.2 ± 97.97	-	$9.9e3 \pm 40.81$	$7.6e3 \pm 81.96$	_	11.0 ± 0.00	_
wall time							
$\epsilon = 1$	$\boldsymbol{0.05 \pm 0.00}$	_	$\textbf{0.02} \pm \textbf{0.00}$	$\textbf{0.02} \pm \textbf{0.00}$	$\boldsymbol{0.14 \pm 0.01}$	$\boldsymbol{0.24 \pm 0.01}$	0.04 ± 0.00
$\epsilon = 10^{-1}$	$\boldsymbol{0.07 \pm 0.00}$	$\boldsymbol{0.06 \pm 0.00}$	$\boldsymbol{0.06 \pm 0.00}$	$\textbf{0.04} \pm \textbf{0.00}$	$\boldsymbol{0.48 \pm 0.02}$	$\boldsymbol{0.32 \pm 0.01}$	0.23 ± 0.00
$\epsilon=10^{-4}$	2.08 ± 0.02	14.38 ± 0.32	$\boldsymbol{7.69 \pm 0.13}$	10.56 ± 0.14	-	$\textbf{0.27} \pm \textbf{0.01}$	-
$\epsilon = 10^{-5}$	4.99 ± 0.59	_	24.49 ± 0.14	35.67 ± 0.40	_	$\textbf{0.27} \pm \textbf{0.02}$	_

experiments, the training set is generated with $\theta \sim U[0.125\pi, 0.375\pi]$ and $\lg \frac{1}{\epsilon} \sim U[0, 5]$, while the test set is generated with $\theta = 0.05\pi, 0.12\pi, 0.4\pi$ and 0.5π (i.e. a test for OoD transfer with respect to θ).

Table 4 shows that Meta-MgNet has superior ability of OoD transfer with respect to ϵ in comparison with PDE-MgNet, while **Table 5** shows that Meta-MgNet is noticeably superior in OoD transfer with respect to θ in comparison with PDE-MgNet. Note that classical methods are not learning-based, and hence they do not have the issue of OoD transfer. Thus, the results of MG methods in Table 4 are copied from Table 2.

5.1.4. Further discussions

It is also worth noting that, for both the scenarios of in-distribution generalization and OoD transfer, Meta-MgNet even outperforms PDE-MgNet- η which uses the exact same η for both training and testing. This shows the benefit of treating the problem of solving parameterized PDEs as a multi-task learning problem. From what it seems, the hypernetwork, i.e. Meta-NN, is able to extract certain common structure hidden within the tasks which helps with solving each individual task.

The training time of Meta-MgNet is around 45 minutes, while it is around 10 minutes for PDE-MgNet- η for each η . Although training Meta-MgNet takes more times than PDE-MgNet- η for each η . In practical application, especially multiquery scenarios, the utility of PDE-MgNet- η can be significantly reduced due to the constant retraining. Therefore, Meta-MgNet has a clear overall advantage.

5.2. 3D anisotropic diffusion equations

Consider the following 3D anisotropic diffusion equation

$$\begin{cases} -\nabla \cdot (C\nabla u) = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial \Omega, \end{cases} \quad \text{and} \quad C = \begin{pmatrix} \epsilon_0 \\ \epsilon_1 \\ \epsilon_2 \end{pmatrix},$$

Table 5The mean and std of the number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method on the testing set. "-" means the algorithm does not converge within 10⁴ iterations.

#iterations	Meta-MgNet, $\theta=0.05\pi$	PDE-MgNet, $\theta = 0.05\pi$	Meta-MgNet, $\theta=0.12\pi$	PDE-MgNet, $\theta=0.12\pi$
$\epsilon = 1$	3.0 ± 0.00	-	3.0 ± 0.00	-
$\epsilon = 10^{-1}$	10.6 ± 0.49	=	10.1 ± 0.30	132.6 ± 3.10
$\epsilon = 10^{-2}$	71.5 ± 1.57	=	72.0 ± 1.61	566.3 ± 11.36
$\epsilon = 10^{-3}$	322.4 ± 7.03	=	233.2 ± 7.49	$2.1e3 \pm 55.09$
$\epsilon = 10^{-4}$	526.7 ± 14.64	-	306.0 ± 7.78	$2.8e3 \pm 150.40$
$\epsilon=10^{-5}$	557.4 ± 14.72	-	314.0 ± 4.86	$2.8e3 \pm 33.31$
wall time				
$\epsilon = 1$	0.03 ± 0.00	-	0.03 ± 0.00	-
$\epsilon = 10^{-1}$	0.07 ± 0.00	-	0.07 ± 0.00	$\textbf{0.35} \pm \textbf{0.01}$
$\epsilon = 10^{-2}$	0.46 ± 0.01	_	0.46 ± 0.01	$\boldsymbol{1.49 \pm 0.04}$
$\epsilon=10^{-3}$	2.05 ± 0.04	_	1.48 ± 0.05	$\boldsymbol{5.42 \pm 0.18}$
$\epsilon=10^{-4}$	3.34 ± 0.09	-	1.95 ± 0.05	$\textbf{7.49} \pm \textbf{0.45}$
$\epsilon = 10^{-5}$	3.54 ± 0.09	-	$\boldsymbol{1.99 \pm 0.03}$	$\boldsymbol{7.49 \pm 0.15}$
#iterations	Meta-MgNet, $\theta=0.4\pi$	PDE-MgNet, $\theta = 0.4\pi$	Meta-MgNet, $\theta=0.5\pi$	PDE-MgNet, $\theta = 0.5\pi$
$\epsilon = 1$	3.0 ± 0.00	=	3.0 ± 0.0	-
$\epsilon = 10^{-1}$	9.0 ± 0.00	51.7 ± 1.27	8.9 ± 0.30	49.3 ± 1.42
$\epsilon = 10^{-2}$	65.2 ± 1.54	434.5 ± 7.75	53.5 ± 1.12	428.8 ± 12.83
$\epsilon=10^{-3}$	240.3 ± 3.41	$1.7e3 \pm 50.92$	262.9 ± 5.96	$2.8e3 \pm 16.81$
$\epsilon = 10^{-4}$	327.5 ± 5.33	$2.4e3 \pm 67.84$	526.4 ± 25.51	_
$\epsilon = 10^{-5}$	339.5 ± 6.92	$2.5e3 \pm 85.99$	908.7 ± 27.43	-
wall time				
$\epsilon = 1$	0.03 ± 0.00	-	0.03 ± 0.00	-
$\epsilon = 10^{-1}$	0.06 ± 0.00	0.14 ± 0.00	0.06 ± 0.00	$\boldsymbol{0.13 \pm 0.01}$
$\epsilon = 10^{-2}$	0.42 ± 0.01	1.15 ± 0.03	0.35 ± 0.01	$\boldsymbol{1.13\pm0.03}$
$\epsilon = 10^{-3}$	1.53 ± 0.02	4.48 ± 0.12	1.67 ± 0.04	7.34 ± 0.11
$\epsilon = 10^{-4}$	2.09 ± 0.03	6.38 ± 0.22	3.35 ± 0.16	_

with $\Omega = [0, 1]^3$ and $\epsilon_0, \epsilon_1, \epsilon_2 > 0$. Without loss of generality, we set $\epsilon_0 = 1$.

5.2.1. Settings: training, testing and hyperparameter selection

- (1) The settings of training and testing are the same as the 2D case.
- (2) We use $N \times N \times N$ rectangular mesh and use Q_1 -element to discretize the PDE. Let N=64 and the number of layers J=4. The prolongations $\mathcal P$ and restrictions $\mathcal R$ are given by the traditional 9-point stencil. The $\mathcal A^\ell$ on the coarse grid is given by the geometry MG method. It is easy to verify that each A^ℓ is equal for Q_1 -element. We chose \-Cycle structure of all algorithms with $\nu_1=2, \nu_2=\cdots=\nu_I=1$.
- (3) For PDE-MgNet, the size of the kernel for the convolution smoother is $3 \times 3 \times 3$. For Meta-MgNet, we choose the SC method B_{SC} in **Algorithm 4** and use the Meta-NN in (19). We simply set the convolution smoother to be one layer CNN without activation, and the size of the kernel for the output is $7 \times 7 \times 7$. We set the number of channels of Meta-NN to l = 3.

5.2.2. In-distribution generalization

In this group of experiments, the training data set is generated by sampling ϵ_1 and ϵ_2 from distribution $\lg \frac{1}{\epsilon_1} \sim U[0,5]$ and $\lg \frac{1}{\epsilon_2} \sim U[0,5]$ respectively. **Table 6** shows that Meta-MgNet is more efficient than classic MG methods.

5.2.3. Out-of-distribution (OoD) transfer

In this group of experiments, the training data set is generated by sampling ϵ_1 and ϵ_2 from distribution $\lg \frac{1}{\epsilon_1} \sim U[3,4]$ and $\lg \frac{1}{\epsilon_2} \sim U[3,4]$ respectively. **Table 7** shows that Meta-MgNet has an overall best performance, and PDE-MgNet has trouble in OoD transfer with respect to ϵ_1 and ϵ_2 .

Table 6The mean and std of the number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method on the testing set. "-" means the algorithm does not converge within 10⁴ iterations.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η	MG (Krylov)	MG (GS)	MG (Jacobi)
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-1})$	5.0 ± 0.00	11.0 ± 0.00	11.0 ± 0.00	7.0 ± 0.00	53.0 ± 1.05	-
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-2})$	13.0 ± 0.00	91.4 ± 1.20	46.6 ± 1.69	38.3 ± 1.00	159.9 ± 5.27	-
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-5})$	156.2 ± 3.57	$3.2\text{e}3 \pm 73.59$	475.0 ± 8.00	606.8 ± 30.08	-	-
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-2})$	10.3 ± 0.46	116.4 ± 0.49	54.3 ± 0.78	45.7 ± 0.78	178.4 ± 2.54	291.52 ± 9.14
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-5})$	73.4 ± 0.80	$3.5e3 \pm 77.82$	771.5 ± 9.39	631.70 ± 16.64	-	-
$(\epsilon_1, \epsilon_2) = (10^{-5}, 10^{-5})$	111.6 ± 0.92	$8.3\text{e}3 \pm 65.18$	$5.5\mathrm{e}3 \pm 99.85$	$1.8\text{e}3 \pm 20.23$	_	-
wall time						
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-1})$	0.13 ± 0.00	0.09 ± 0.00	0.09 ± 0.00	$\textbf{0.07} \pm \textbf{0.00}$	4.29 ± 0.10 -	_
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-2})$	$\textbf{0.30} \pm \textbf{0.00}$	$\boldsymbol{0.65 \pm 0.01}$	$\boldsymbol{0.34 \pm 0.02}$	$\boldsymbol{0.32 \pm 0.01}$	14.2 ± 0.64	_
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-5})$	3.40 ± 0.08	22.53 ± 0.51	$\textbf{3.34} \pm \textbf{0.06}$	4.93 ± 0.25	-	_
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-2})$	$\textbf{0.25} \pm \textbf{0.01}$	$\boldsymbol{0.82 \pm 0.00}$	$\boldsymbol{0.39 \pm 0.01}$	$\boldsymbol{0.38 \pm 0.01}$	14.67 ± 0.61	6.00 ± 0.00
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-5})$	$\boldsymbol{1.62 \pm 0.02}$	24.73 ± 0.55	$\boldsymbol{5.42 \pm 0.06}$	$\boldsymbol{5.12 \pm 0.13}$	-	-
$(\epsilon_1, \epsilon_2) = (10^{-5}, 10^{-5})$	$\textbf{2.45} \pm \textbf{0.02}$	58.87 ± 0.47	38.93 ± 0.71	14.90 ± 0.16	_	_

Table 7The mean and std of the number of iterations (when the stopping criteria is met) and the wall time of Meta-MgNet, PDE-MgNet and MG method on the testing set. "-" means the algorithm does not converge within 10⁴ iterations.

	=	_				
#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η	MG (Krylov)	MG (GS)	MG (Jacobi)
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-1})$	10.0 ± 0.00	-	11.0 ± 0.00	7.0 ± 0.00	53.0 ± 1.05	-
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-2})$	43.1 ± 0.54	_	46.6 ± 1.69	38.3 ± 1.00	159.9 ± 5.27	-
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-5})$	755.4 ± 55.88	_	475.0 ± 8.00	606.8 ± 30.08	-	-
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-2})$	11.0 ± 0.00	_	54.3 ± 0.78	45.7 ± 0.78	178.4 ± 2.54	291.52 ± 9.14
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-5})$	106.7 ± 1.35	-	771.5 ± 9.39	631.7 ± 16.64	-	-
$(\epsilon_1, \epsilon_2) = (10^{-5}, 10^{-5})$	125.7 ± 2.19	-	$5.5\text{e}3 \pm 99.85$	$1.8\text{e}3 \pm 20.23$	-	-
wall time						
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-1})$	0.24 ± 0.01	_	0.09 ± 0.00	$\textbf{0.07} \pm \textbf{0.00}$	4.29 ± 0.10 -	_
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-2})$	$\boldsymbol{0.95 \pm 0.01}$	_	$\boldsymbol{0.34 \pm 0.02}$	$\textbf{0.32} \pm \textbf{0.01}$	14.2 ± 0.64	_
$(\epsilon_1, \epsilon_2) = (10^{-1}, 10^{-5})$	16.37 ± 1.22	_	$\textbf{3.34} \pm \textbf{0.06}$	4.93 ± 0.25	-	-
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-2})$	$\textbf{0.26} \pm \textbf{0.00}$	_	$\boldsymbol{0.39 \pm 0.01}$	$\boldsymbol{0.38 \pm 0.01}$	14.67 ± 0.61	$\boldsymbol{6.00 \pm 0.00}$
$(\epsilon_1, \epsilon_2) = (10^{-2}, 10^{-5})$	$\textbf{2.34} \pm \textbf{0.03}$	_	$\boldsymbol{5.42 \pm 0.06}$	$\boldsymbol{5.12 \pm 0.13}$	_	_
$(\epsilon_1, \epsilon_2) = (10^{-5}, 10^{-5})$	$\boldsymbol{2.76 \pm 0.05}$	_	38.93 ± 0.71	14.90 ± 0.16		_

5.3. Why is it challenging to train a convergent PDE-MgNet for all η ?

For both the 2D and 3D anisotropic diffusion equations, we found it challenging to train appropriate weights for PDE-MgNet that can generalize beyond its training setting. This is, in fact, the motivation of viewing solving parameterized PDEs as multi-task learning and introducing Meta-MgNet. In this subsection, we conduct a simple experiment to demonstrate this issue with PDE-MgNet.

Consider the 3D anisotropic diffusion equation. We choose two different distributions for the parameters D_1 : $\lg_2 \epsilon_1 \sim U[-2,-1]$, $\lg_2 \epsilon_2 \sim U[1,2]$ and D_2 : $\lg_2 \epsilon_1 \sim U[1,2]$, $\lg_2 \epsilon_2 \sim U[-2,-1]$ to train two PDE-MgNet. We present the weights of the convolution smoothers at the finest level ($\ell=1$), namely K¹ in (25) and K² in (26):

$$\begin{split} \mathsf{K}^1 = \begin{bmatrix} \begin{bmatrix} -0.0170 & -0.0503 & -0.0169 \\ 0.0051 & -0.2370 & 0.0051 \\ -0.0171 & -0.0504 & -0.0170 \end{bmatrix} \begin{bmatrix} 0.0389 & 0.2431 & 0.0389 \\ -0.0696 & 1.1127 & -0.0695 \\ 0.0388 & 0.2431 & 0.0389 \end{bmatrix} \begin{bmatrix} -0.0171 & -0.0503 & -0.0170 \\ 0.0050 & -0.2370 & 0.0051 \\ -0.0170 & -0.0503 & -0.0169 \end{bmatrix} \end{bmatrix}, \\ \mathsf{K}^2 = \begin{bmatrix} \begin{bmatrix} -0.0127 & -0.0331 & -0.0126 \\ 0.0339 & 0.2175 & 0.0340 \\ -0.0126 & -0.0332 & -0.0127 \end{bmatrix} \begin{bmatrix} 0.0113 & -0.2041 & 0.0114 \\ -0.0766 & 1.0647 & -0.0765 \\ 0.0113 & -0.2041 & 0.0113 \end{bmatrix} \\ \times \begin{bmatrix} -0.0126 & -0.0332 & -0.0126 \\ 0.0340 & 0.2175 & 0.0340 \\ -0.0127 & -0.0331 & -0.0127 \end{bmatrix} \end{bmatrix}. \end{split}$$

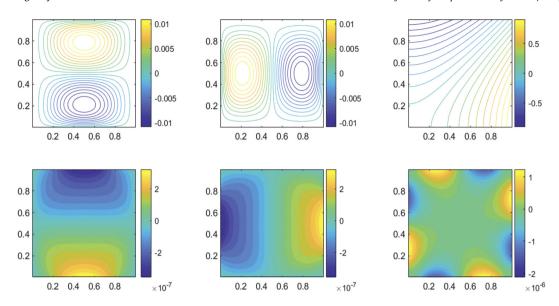


Fig. 10. The three pictures on the top row are the numerical solutions and the three at the bottom row are the error of u, v and p with $(a_x, a_y) = (0, 0)^{\top}$. The number of iterations is 367 and the wall time is 17.57 s for Meta-MgNet.

Noting the red numbers in K^1 and K^2 , we can see that $K^1_{-1,0,0}$, $K^1_{1,0,0} < 0$ while $K^2_{-1,0,0}$, $K^2_{1,0,0} > 0$; and $K^1_{0,-1,0}$, $K^1_{0,1,0} > 0$ while $K^2_{0,-1,0}$, $K^2_{0,1,0} < 0$. Furthermore, if we use K^1 to smooth the PDEs with parameters generated from D_2 , the error will increase, which means the convolution smoother with kernel K^1 is not fit for D_2 . We have the same issue with K_2 and D_1 . This phenomenon indicates that different distribution of parameters of the parameterized PDE may lead to weights of PDE-MgNet of contradictory effects. This is not only for PDE-MgNet, but rather an issue often occurs for supervised learning models. In contrast, Meta-MgNet handles this issue gracefully by adjusting the weights according to η in an adaptive fashion.

5.4. Ossen equations

The right-hand-side function \underline{f} in previous numerical experiments is randomly generated from the normal distribution. In this section, we include practical numerical examples to show the performance of Meta-MgNet. We adopt Ossen equations as an example:

$$\begin{cases} -\mu \, \Delta \underset{\sim}{u} + (\underset{\sim}{a} \cdot \nabla) \underset{\sim}{u} + \nabla p = \underset{\sim}{f}, & \text{in } \Omega, \\ -\text{div} \underset{\sim}{u} = \underset{\sim}{0}, & \text{in } \Omega, \\ \underset{\sim}{u} = \underset{\sim}{0}, & \text{on } \partial \Omega, \end{cases}$$

where $\mu=(u,v)^{\top}$, $\mu=\frac{1}{Re}$, Re is the Reynold number, and $\varrho=(a_x,a_y)^{\top}$. Without loss of generality, we let $\mu=1$. We choose $\mu=\left(\frac{-2xy^2(1-x)(1-2x)(1-y)^2}{2x^2y(1-y)(1-2y)(1-x)^2}\right)$ and $p=x^2-y^2$ as the solution and calculate the analytic form of the corresponding right-hand-side function f. The training data set is construct from sampling $a_x\sim U[0,200]$ and $a_y\sim U[0,200]$. We use the MAC scheme [38,6] to discretize Ossen equations and the mesh size is 512 \times 512. Except for the settings mentioned above, other settings are the same as **Section 5.1**. Numerical solutions and error maps for a few different ϱ are presented in **Fig. 10–13** with the number of iterations and the wall time of Meta-MgNet recorded in the captions.

6. Conclusions and future work

Solving parameterized PDEs is an essential and yet challenging task. In this paper, we provided a new perspective on the problem by viewing it as multi-task learning. With this, we proposed a new meta-learning based solver called Meta-MgNet by introducing a carefully designed hypernetwork (called Meta-NN) to the PDE-MgNet. Numerical experiments on 2D and 3D anisotropic diffusion equations showed that Meta-MgNet significantly outperforms the supervised learning-based PDE-MgNet and classical MG methods. Furthermore, Meta-MgNet manifested a clear advantage in training and generalization over PDE-MgNet, which demonstrated the feasibility of the proposed multi-task perspective and meta-learning approach to solving parameterized PDEs.

This paper only discussed using meta-learning to improve smoothers because the prolongations and restrictions in classic MG methods are efficient enough to solve the three PDEs considered in this paper. As for some other PDEs such as Helmholtz

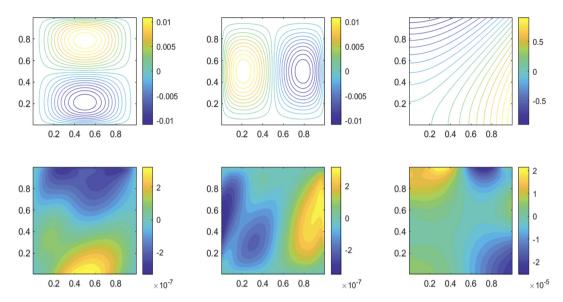


Fig. 11. The three pictures on the top row are the numerical solutions and the three at the bottom row are the error of u, v and p with $(a_x, a_y) = (50, 100)^{\top}$. The number of iterations is 281 and the wall time is 13.6 s for Meta-MgNet.

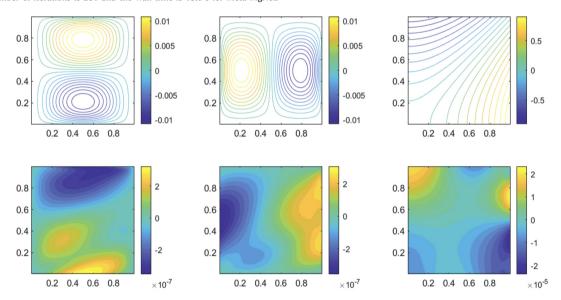


Fig. 12. The three pictures on the top row are the numerical solutions and the three at the bottom row are the error of u, v and p with $(a_x, a_y) = (100, 50)^{\top}$. The number of iterations is 275 and the wall time is 13.38 s for Meta-MgNet.

equations, the convolutional prolongations and restrictions may not work well. Therefore, it is worth exploring a data-driven approach to improve prolongations and restrictions as well. Furthermore, we only considered uniform mesh in this paper. We may consider generalizing MgNet or Meta-MgNet to nonuniform meshes, such as the triangular mesh, by exploiting tools from geometric deep learning, such as graph (convolutional) neural networks.

CRediT authorship contribution statement

Yuyan Chen: Conceptualization, Formal analysis, Methodology, Software, Validation, Writing – original draft. **Bin Dong:** Conceptualization, Funding acquisition, Methodology, Resources, Supervision, Writing – review & editing. **Jinchao Xu:** Funding acquisition, Resources, Validation, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

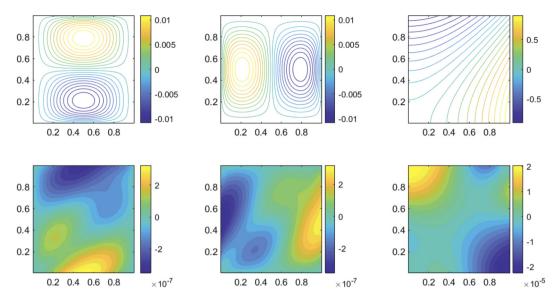


Fig. 13. The three pictures on the top row are the numerical solutions and the three at the bottom row are the error of u, v and p with $(a_x, a_y) = (100, 100)^{\top}$. The number of iterations is 235 and the wall time is 11.52 for Meta-MgNet.

Table 8The mean and std of the number of iterations and the wall time of Meta-MgNet. No matter if Meta-MgNet is fine-tuned or not, the results are almost the same.

#iterations	Meta-MgNet fine-tuning	Meta-MgNet
$\epsilon = 1$	3.0 ± 0.00	4.0 ± 0.00
$\epsilon = 10^{-1}$	7.0 ± 0.00	7.5 ± 0.50
$\epsilon = 10^{-2}$	32.7 ± 0.90	35.1 ± 1.04
$\epsilon = 10^{-3}$	192.7 ± 4.29	171.6 ± 6.34
$\epsilon = 10^{-4}$	352.2 ± 7.60	375.2 ± 5.88
wall time		
$\epsilon = 1$	0.03 ± 0.00	0.03 ± 0.00
$\epsilon = 10^{-1}$	0.05 ± 0.00	$\boldsymbol{0.05 \pm 0.00}$
$\epsilon = 10^{-2}$	0.21 ± 0.00	$\boldsymbol{0.22 \pm 0.01}$
$\epsilon = 10^{-3}$	1.18 ± 0.03	$\boldsymbol{1.06 \pm 0.04}$
$\epsilon = 10^{-4}$	2.16 ± 0.05	2.31 ± 0.03

Acknowledgement

Yuyan Chen was supported in part by the PSU-PKU Joint Center for Computational Mathematics and Applications, Bin Dong by the National Natural Science Foundation of China (grant No. 11831002), Beijing Natural Science Foundation (grant No. 180001), and Beijing Academy of Artificial Intelligence (BAAI), and Jinchao Xu by the Verne M. William Professorship Fund from the Pennsylvania State University and the National Science Foundation (grant No. DMS-1819157).

Appendix A

In the appendix, we add more numerical experiments to support some of our hyper-parameter choices. We use interpolation of 2D anisotropic diffusion equations as an example, and the setting of these experiments is the same as section 5.1.

A.1. Appendix I

The **Table 8** shows the efficiency of fine-tuning for Meta-MgNet. We can find that the parameters inferred by meta smoother B_{sc} are good enough so that the result is almost the same after fine-tuning. Therefore, we can skip the fine-tuning stage.

Table 9 The mean and std of the numbers of iteration and the wall time of different choices of $v_1, ..., v_J$. The parameters of the 2D anisotropic diffusion equation are $\epsilon = 10^{-2}, \theta = 0$.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η
$(v_1,, v_J) = (1, 1, 1, 1, 1)$	42.1 ± 1.58	356.6 ± 1.36	230.60 ± 6.99
$(v_1,, v_J) = (2, 1, 1, 1, 1)$	35.1 ± 1.04	178.9 ± 2.74	149.7 ± 3.44
$(v_1,, v_J) = (3, 1, 1, 1, 1)$	21.9 ± 1.81	210.90 ± 4.23	198.60 ± 6.89
$(v_1,, v_J) = (2, 2, 1, 1, 1)$	22.90 ± 0.70	214.50 ± 2.42	199.00 ± 4.75
wall time			
$(v_1,, v_J) = (1, 1, 1, 1, 1)$	0.23 ± 0.01	0.91 ± 0.01	0.59 ± 0.02
$(v_1,, v_J) = (2, 1, 1, 1, 1)$	0.22 ± 0.01	$\textbf{0.44} \pm \textbf{0.01}$	$\textbf{0.37} \pm \textbf{0.01}$
$(v_1,, v_J) = (3, 1, 1, 1, 1)$	$\textbf{0.16} \pm 0.01$	0.56 ± 0.01	$\boldsymbol{0.52 \pm 0.02}$
$(v_1,, v_J) = (2, 2, 1, 1, 1)$	0.17 ± 0.00	0.56 ± 0.01	0.53 ± 0.01

Table 10 The mean and std of the number of iterations and the wall time of different choices of $v_1, ..., v_J$. The parameters of the 2D anisotropic diffusion equation are $\epsilon = 10^{-4}, \theta = 0$.

#iterations	Meta-MgNet	PDE-MgNet	PDE-MgNet- η
$(v_1,, v_J) = (1, 1, 1, 1, 1)$	$1.2e3 \pm 22.51$	-	$6.9e3 \pm 140.05$
$(v_1,, v_J) = (2, 1, 1, 1, 1)$	375.2 ± 5.88	_	$3.1e3 \pm 35.70$
$(v_1,, v_J) = (3, 1, 1, 1, 1)$	355.30 ± 10.51	_	$3.6e3 \pm 40.44$
$(v_1,, v_J) = (2, 2, 1, 1, 1)$	339.00 ± 12.70	-	$2.8\text{e}3 \pm 33.22$
wall time			
$(v_1,, v_J) = (1, 1, 1, 1, 1)$	6.36 ± 0.11	-	17.52 ± 0.34
$(v_1,, v_J) = (2, 1, 1, 1, 1)$	$\textbf{2.31} \pm \textbf{0.03}$	_	$\textbf{7.69} \pm \textbf{0.13}$
$(v_1,, v_J) = (3, 1, 1, 1, 1)$	2.50 ± 0.07	_	$\boldsymbol{9.42 \pm 0.11}$
$(v_1,, v_J) = (2, 2, 1, 1, 1)$	2.34 ± 0.09	-	8.45 ± 0.10

Table 11 The mean and std of the number of iterations and the wall time for the convolutional smoother $B_{\rm d}$ and the SC method $B_{\rm SC}$.

#iterations	B_{d}	B_{sc}
$\epsilon = 1$	-	4.0 ± 0.00
$\epsilon = 10^{-1}$	58.9 ± 0.30	7.5 ± 0.50
$\epsilon = 10^{-2}$	597.8 ± 8.08	35.1 ± 1.04
$\epsilon = 10^{-3}$	$5.5e3 \pm 60.89$	171.6 ± 6.34
$\epsilon = 10^{-4}$	_	375.2 ± 5.88
wall time		
$\epsilon = 1$	-	0.03 ± 0.00
$\epsilon = 10^{-1}$	0.21 ± 0.00	$\boldsymbol{0.05 \pm 0.00}$
$\epsilon = 10^{-2}$	2.07 ± 0.03	$\boldsymbol{0.22 \pm 0.01}$
$\epsilon = 10^{-3}$	19.04 ± 0.20	$\boldsymbol{1.06 \pm 0.04}$
$\epsilon = 10^{-4}$	_	2.31 ± 0.03

A.2. Appendix II

In section 5.1, we choose $(v_1, ..., v_J) = (2, 1, 1, 1, 1)$. Now, we compare the result of several groups of $v_1, ..., v_J$. Since it is easier to estimate error on coarse grid, to set $v_i > 1$, i = 3, 4, 5 is not necessary. Thus, we only test some pair of v_1 and v_2 , and set $v_i = 1$ for i = 3, 4, 5. The numerical experiments in **Tables 9 and 10** show $(v_1, ..., v_J) = (2, 1, 1, 1, 1)$ is a better choice.

A.3. Appendix III

We compare the efficiency between B_d in **Algorithm 3** and B_{sc} in **Table 11**. As we can see, B_{sc} is much better than B_d . Therefore, we choose B_{sc} as the meta smoother in other numerical experiments in this paper.

References

- [1] C. Audouze, F. De Vuyst, P. Nair, Reduced-order modeling of parameterized pdes using time-space-parameter principal component analysis, Int. J. Numer. Methods Eng. 80 (2009).
- [2] R.E. Bank, C.C. Douglas, Sharp estimates for multigrid rates of convergence with general smoothing and acceleration, SIAM J. Numer. Anal. 22 (1985) 617–633.
- [3] C. Beck, W. E, A. Jentzen, Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations, J. Nonlinear Sci. 29 (2019) 1563–1619, https://doi.org/10.1007/s00332-018-9525-3.
- [4] P. Benner, S. Gugercin, K. Willcox, A survey of projection-based model reduction methods for parametric dynamical systems, SIAM Rev. 57 (2015) 483–531, https://doi.org/10.1137/130932715.
- [5] A. Brock, T. Lim, J.M. Ritchie, N. Weston, SMASH: one-shot model architecture search through hypernetworks, CoRR abs/1708.05344, http://arxiv.org/abs/1708.05344, arXiv:1708.05344, 2017.
- [6] L. Chen, M. Wang, L. Zhong, Convergence analysis of triangular mac schemes for two dimensional Stokes equations, J. Sci. Comput. 63 (2015) 716-744.
- [7] T.Q. Chen, Y. Rubanova, J. Bettencourt, D.K. Duvenaud, Neural ordinary differential equations, in: Advances in Neural Information Processing Systems, 2018, pp. 6571–6583.
- [8] Y. Chen, W. Yu, T. Pock, On learning optimized reaction diffusion processes for effective image restoration, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015.
- [9] N. Discacciati, J.S. Hesthaven, D. Ray, Mathicse technical report: controlling oscillations in high-order discontinuous Galerkin schemes using artificial viscosity tuned by neural networks, http://infoscience.epfl.ch/record/263616, 2019, 10.5075/epfl-MATHICSE-263616.
- [10] W. E, A proposal on machine learning via dynamical systems, Commun. Math. Stat. 5 (2017) 1-11.
- [11] W. E, B. Yu, The deep Ritz method: a deep learning-based numerical algorithm for solving variational problems, Commun. Math. Stat. 6 (2018) 1–12, https://doi.org/10.1007/s40304-018-0127-z.
- [12] J. Feliu-Faba, Y. Fan, L. Ying, Meta-learning pseudo-differential operators with deep neural networks, arXiv e-prints, arXiv:1906.06782, 2019.
- [13] C. Finn, P. Abbeel, S. Levine, Model-agnostic meta-learning for fast adaptation of deep networks, in: Proceedings of the 34th International Conference on Machine Learning Volume 70, JMLR.org, 2017, pp. 1126–1135.
- [14] S. Fresca, L. Dede, A. Manzoni, A comprehensive deep learning-based approach to reduced order modeling of nonlinear time-dependent parametrized pdes, arXiv:2001.04001, 2020.
- [15] D. Greenfeld, M. Galun, R. Kimmel, I. Yavneh, R. Basri, Learning to optimize multigrid pde solvers, arXiv:1902.10248, 2019.
- [16] K. Gregor, Y. LeCun, Learning fast approximations of sparse coding, in: Proceedings of the 27th International Conference on International Conference on Machine Learning, Omnipress, 2010, pp. 399–406.
- [17] D. Ha, A.M. Dai, Q.V. Le, Hypernetworks, CoRR abs/1609.09106, http://arxiv.org/abs/1609.09106, arXiv:1609.09106, 2016.
- [18] E. Haber, L. Ruthotto, Stable architectures for deep neural networks, Inverse Probl. 34 (2017) 014004.
- [19] W. Hackbusch, Multi-Grid Methods and Applications, vol. 4, Springer Science & Business Media, 2013,
- [20] J. Han, A. Jentzen, W. E. Solving high-dimensional partial differential equations using deep learning, Proc. Natl. Acad. Sci 115 (2018) 8505–8510, https://doi.org/10.1073/pnas.1718942115, https://www.pnas.org/content/115/34/8505, arXiv: https://www.pnas.org/content/115/34/8505.full.pdf.
- [21] J. He, J. Xu, Mgnet: a unified framework of multigrid and convolutional neural network, CoRR abs/1901.10415, http://arxiv.org/abs/1901.10415, arXiv: 1901.10415, 2019
- [22] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.
- [23] G. Hinton, L. Deng, D. Yu, G. Dahl, A.r.Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, B. Kingsbury, T. Sainath, Deep neural networks for acoustic modeling in speech recognition, IEEE Signal Process. Mag. 29 (2012) 82–97, https://www.microsoft.com/en-us/research/publication/deep-neural-networks-for-acoustic-modeling-in-speech-recognition/.
- [24] T. Hospedales, A. Antoniou, P. Micaelli, A. Storkey, Meta-learning in neural networks: a survey, arXiv:2004.05439, 2020.
- [25] J.T. Hsieh, S. Zhao, S. Eismann, L. Mirabella, S. Ermon, Learning neural PDE solvers with convergence guarantees, arXiv e-prints, arXiv:1906.01200, 2019.
- [26] G. Huang, Z. Liu, L. Van Der Maaten, K.Q. Weinberger, Densely connected convolutional networks, in: CVPR, 2017, p. 3.
- [27] A. Katrutsa, T. Daulbaev, I. Oseledets, Deep multigrid: learning prolongation and restriction matrices, arXiv preprint, arXiv:1711.03825, 2017.
- [28] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, Nature 521 (2015) 436.
- [29] K. Lee, K.T. Carlberg, Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders, J. Comput. Phys. 404 (2020) 108973, https://doi.org/10.1016/ji.jcp.2019.108973, http://www.sciencedirect.com/science/article/pii/S0021999119306783.
- [30] Y. Li, J. Lu, A. Mao, Variational training of neural network approximations of solution maps for physical models, arXiv e-prints, arXiv:1905.02789, 2019.
- [31] Z. Liu, W. Cai, Z.Q.J. Xu, Multi-scale deep neural network (mscalednn) for solving Poisson-Boltzmann equation in complex domains, arXiv preprint, arXiv:2007.11207, 2020.
- [32] Z. Liu, H. Mu, X. Zhang, Z. Guo, X. Yang, K.T. Cheng, J. Sun, Metapruning: meta learning for automatic neural network channel pruning, in: The IEEE International Conference on Computer Vision (ICCV), 2019.
- [33] J. Lorraine, D. Duvenaud, Stochastic hyperparameter optimization through hypernetworks, CoRR abs/1802.09419, http://arxiv.org/abs/1802.09419, arXiv: 1802.09419, 2018.
- [34] Y. Lu, A. Zhong, Q. Li, B. Dong, Beyond finite layer neural networks: bridging deep architectures and numerical differential equations, in: Proceedings of the 35th International Conference on Machine Learning, PMLR, 2018.
- [35] I. Luz, M. Galun, H. Maron, R. Basri, I. Yavneh, Learning algebraic multigrid using graph neural networks, arXiv:2003.05744, 2020.
- [36] J. Magiera, D. Ray, J.S. Hesthaven, C. Rohde, Constraint-aware neural networks for Riemann problems, arXiv e-prints, arXiv:1904.12794, 2019.
- [37] A. Nichol, J. Achiam, J. Schulman, On first-order meta-learning algorithms, CoRR abs/1803.02999, http://arxiv.org/abs/1803.02999, arXiv:1803.02999, 2018
- [38] R. Nicolaides, X. Wu, Analysis and convergence of the mac scheme. ii. Navier-Stokes equations, Math. Comput. 65 (1996) 29-44.
- [39] M. Raissi, P. Perdikaris, G.E. Karniadakis, Physics informed deep learning (part I): data-driven solutions of nonlinear partial differential equations, arXiv e-prints, arXiv:1711.10561, 2017.
- [40] M. Raissi, P. Perdikaris, G.E. Karniadakis, Physics informed deep learning (part II): data-driven discovery of nonlinear partial differential equations, arXiv e-prints, arXiv:1711.10566, 2017.
- [41] D. Ray, J.S. Hesthaven, Detecting troubled-cells on two-dimensional unstructured grids using a neural network, J. Comput. Phys. 397 (2019) 108845, https://doi.org/10.1016/j.jcp.2019.07.043, http://www.sciencedirect.com/science/article/pii/S0021999119305297.
- [42] H. Sheng, C. Yang, Pfnn: a penalty-free neural network method for solving a class of second-order boundary-value problems on complex geometries, arXiv preprint, arXiv:2004.06490, 2020.
- [43] D. Silver, A. Huang, C.J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al., Mastering the game of go with deep neural networks and tree search, Nature 529 (2016) 484.
- [44] J. Sirignano, K. Spiliopoulos, Dgm: a deep learning algorithm for solving partial differential equations, J. Comput. Phys. 375 (2018) 1339–1364, https://doi.org/10.1016/j.jcp.2018.08.029, http://www.sciencedirect.com/science/article/pii/S0021999118305527.

- [45] S. Thrun, L. Pratt, Learning to learn: introduction and overview, in: Learning to Learn, 1998.
- [46] Z. Wang, Z. Zhang, A mesh-free method for interface problems using the deep learning approach, J. Comput. Phys. 400 (2020) 108963, https://doi.org/10.1016/j.jcp.2019.108963, http://www.sciencedirect.com/science/article/pii/S0021999119306680.
- [47] H. Xu, H. Chang, D. Zhang, Dlga-pde: discovery of pdes with incomplete candidate library via combination of deep learning and genetic algorithm, J. Comput. Phys. 109584 (2020).
- [48] J. Xu, Theory of Multilevel Methods, volume 8924558 Cornell University, Ithaca, NY, 1989.
- [49] J. Xu, Iterative methods by space decomposition and subspace correction, SIAM Rev. 34 (1992) 581-613.
- [50] J. Xu, L. Zikatanov, The method of alternating projections and the method of subspace corrections in Hilbert space, J. Am. Math. Soc. 15 (2002) 573-597.
- [51] J. Xu, L. Zikatanov, Algebraic multigrid methods, Acta Numer. 26 (2017) 591-721.
- [52] G. Yu, J. Xu, L.T. Zikatanov, Analysis of a two-level method for anisotropic diffusion equations on aligned and nonaligned grids, Numer. Linear Algebra Appl. 20 (2013) 832–851.
- [53] D. Zhang, L. Guo, G.E. Karniadakis, Learning in modal space: solving time-dependent stochastic pdes using physics-informed neural networks, SIAM J. Sci. Comput. 42 (2020) A639–A665, https://doi.org/10.1137/19M1260141.
- [54] H. Zhang, B. Liu, H. Yu, B. Dong, Metainv-net: meta inversion network for sparse view ct image reconstruction, arXiv:2006.00171, 2020.