# Data driven approximation of parametrized PDEs by Reduced Basis and Neural Networks

Niccolò Dal Santo[1], Simone Deparis[1], Luca Pegolotti[1]

July 2, 2019

[1] SCI-SB-SD, École Polytechnique Fédérale de Lausanne (EPFL),
Station 8, 1015 Lausanne, Switzerland.

## Abstract

We are interested in the approximation of partial differential equations with a data-driven approach based on the reduced basis method and machine learning. We suppose that the phenomenon of interest can be modeled by a parametrized partial differential equation, but that the value of the physical parameters is unknown or difficult to be directly measured. Our method allows to estimate fields of interest, for instance temperature of a sample of material or velocity of a fluid, given data at a handful of points in the domain. We propose to accomplish this task with a neural network embedding a reduced basis solver as exotic activation function in the last layer. The reduced basis solver accounts for the underlying physical phenomenon and it is constructed from snapshots obtained from randomly selected values of the physical parameters during an expensive offline phase. The same full order solutions are then employed for the training of the neural network. As a matter of fact, the chosen architecture resembles an asymmetric autoencoder in which the decoder is the reduced basis solver and as such it does not contain trainable parameters. The resulting latent space of our autoencoder includes parameter-dependent quantities feeding the reduced basis solver, which – depending on the considered partial differential equation – are the values of the physical parameters themselves or the affine decomposition coefficients of the differential operators.

## 1 Introduction

During the last decade, machine learning (ML) techniques have gained considerable attention for their ability to exploit data abundance in problems that are typically difficult to be solved through classic algorithmic paradigms. This success is partly motivated by the recent exponential growth in the size of available data, which is due to the technological advancements in consumer devices, and by the increase in computational power of parallel architectures, such as supercomputers and graphics processing units (GPUs). Deep learning (DL) [17, 4] currently represents one of the most active and promising areas of research in machine learning. Notable applications of DL include classification tasks, for instance image recognition [27, 34, 46], text categorization [28, 25] and natural language processing [48, 8], and regression tasks [42], which play an important role e.g. in computer vision [29], object detection [43], stock prediction [13], and cancer detection [24].

In the context of the numerical solution of partial differential equations (PDEs), DL has only lately become an alternative to standard approximation techniques. This is arguably due to the reluctancy of the community to adopt these algorithms because of the lack of a solid theoretical foundation, which instead characterizes classical numerical methods – e.g. the finite element (FE) method – and provides error bounds and stability conditions. In fact, the approximation properties of neural networks (NNs), which are the basis of every DL algorithm, have been well-known for at least thirty years; for example, in [9, 10, 33] the authors show that feed-forward neural networks with one hidden layer are universal approximators. However, these results do not address fundamental aspects regarding the application of neural networks in practice; for example, they do not provide

guidelines on how to design optimal architectures retaining a sufficient level of expressiveness with a minimal number of parameters and they do not focus on the *trainability* of such algorithms. Some of the recent efforts oriented towards defining a mathematical theory of neural networks and which give insights on the aforementioned critcal issues are e.g. [14, 5, 36, 18, 19, 20].

In this paper, we address the problem of approximating the solution of parametrized PDEs by combining the reduced basis (RB) method [37, 22] and NNs. We refer the reader to [30, 23, 47] for other applications of DL to reduced order modelling with the RB method. We set ourselves in the scenario in which we are given data at some input physical points – e.g. the value of the variable of interest – and we are interested in determining the solution itself (or functions thereof) at output locations or even in the whole domain. Moreover, we suppose that the values of the physical parameters characterizing the PDE be either unknown or difficult to be measured (e.g. with non-invasive procedures). The key idea of our approach is to design a NN to learn the input-output mapping by exploiting the knowledge of the underlying physical phenomenon. This concept has already been introduced in the works by E et al. [15], Raissi et al. [38, 39], and Schwab and Zech [41]. In these cases, the authors propose to include the knowledge of the PDE in the functional that is minimized during training: in the former work, the loss function to be minimized is identified as the energy functional of the corresponding differential problem, whereas in the latter the classic mean squared error loss function is modified to account for the residual of the PDE (in strong form) and the initial and boundary conditions. The novelty of this paper consists in including the physical knowledge directly in the NN by considering as output layer a RB solver that, given the input of the previous layers, assembles and solves the reduced system, projects the reduced solution onto the original full-order space, and computes the solution at the desired output locations. We refer to this type of architecture as PDE-aware deep neural network (PDE-DNN). This can be interpreted as a standard multi-layer perceptron (MLP) in which the output layer is equipped with a non-linear activation function – i.e. the RB solver – in the sense that it does not contain any trainable parameters. Contrarily to standard activation functions, however, the RB solver maps vectors – which in our case contain either the physical parameters of the PDE or the coefficients of the affine decomposition – to (typically larger) vectors, i.e. the values of the solution at the output locations. Alternatively, PDE-DNNs can be regarded as asymmetric autoencoders in which the MLP and the RB solver play the roles of the encoder and the decoder, respectively. In this view, the *latent space* coincides with either the space of physical parameters or of the space of coefficients of the affine decomposition, depending on the structure of the underlying PDE.

This paper is structured as follows. In Section 2, we recall basic concepts of NNs and set the notation and the terminology that will be extensively adopted throughout the rest of the paper. In Section 3, we provide a (non-exhaustive, but sufficiently detailed for our purposes) introduction to the reduced order modeling of parametrized PDEs through the RB method. In Section 4, we introduce the concept of PDE-DNN. In Section 5, we present numerical experiments on a three-dimensional example of advection-diffusion allowing for an affine decomposition (Section 5.1) and on two-dimentional non-affine examples, namely a diffusion problem with a non-affine diffusion coefficient (Section 5.2) and the steady Navier-Stokes equations in a bifurcation with a resistive term modelling a stenosis (Section 5.3). Finally, in Section 6 some conclusions are drawn.

## 2   Deep Neural Networks

Artificial neural networks, often simply denoted neural networks (NNs) in the literature, are biologically inspired ML algorithms which are able to *learn* from observational data [21]. Similarly to other ML techniques, NN models are described by means of a collection of trainable parameters. These parameters are inferred during a process whose goal is to find a locally (in the parameter space) optimal choice that leads to an "accurate" input-output mapping of the network on a dataset – usually called training dataset – for which the expected result is known.

Deep Neural Networks (DNNs) are a class of NNs in which simple nonlinear modules are composed together in multiple layers. The first layer is the *input layer* and processes the raw data which
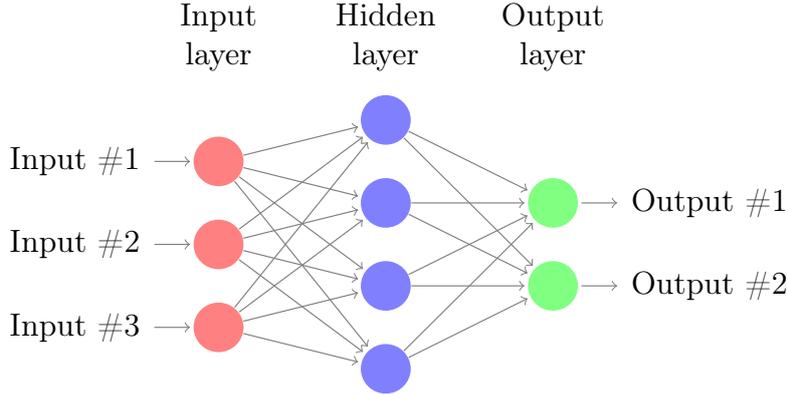
Figure 1: Example of DNN with one fully connected hidden layer

feeds the network, while the last layer is the *output layer* which provides the output of the network; the ones located between the two are called *hidden layers*, see Figure 1. A common choice for the architecture of the DNN is given by ordered fully connected layers of perceptrons, i.e. multi-layer perceptrons (MLPs). Perceptrons are simple computational units containing a set of nodes (the neurons) processing the data coming from the previous layer, and a non-linear *activation function*. MLPs are often preferred over other architectures for statistical regression, because of their ability to approximate non-trivial functions. It has been shown that MLPs with at least one hidden layer and differentiable activation functions are universal function approximators; in particular, MLPs with one hidden layer are able to approximate any continuous function, and MLPs with two hidden layers or more are able to approximate any function [9, 10].

Let us denote with $L$ the number of trainable layers in a MLP (i.e. to the hidden and output layers) and let us use the index 0 to chatacterize the input layer. The input layer performs a trivial mapping of the input of the neural network $\mathbf{x}^{(0)} \in \mathbb{R}^{N_{\mathrm{in}}}$ into itself. The $l^{\mathrm{th}}$ layer of a MLP, with $l > 0$, is a mapping of the form

$$\mathbf{y}^{(l)} = \sigma^{(l)}(W^{(l)}\mathbf{x}^{(l)} + \mathbf{b}^{(l)}), \tag{1}$$

where $\mathbf{x}^{(l)} \in \mathbb{R}^{N_{\mathrm{in}}^{(l)}}$ is the input, $\mathbf{y}^{(l)} \in \mathbb{R}^{N_{\mathrm{out}}^{(l)}}$ is the output, $W^{(l)} \in \mathbb{R}^{N_{\mathrm{out}}^{(l)} \times N_{\mathrm{in}}^{(l)}}$ is a matrix of parameters (weights), $\mathbf{b}^{(l)} \in \mathbb{R}^{N_{\mathrm{out}}^{(l)}}$ is a vector of parameters (biases), and $\sigma^{(l)} : \mathbb{R}^{N_{\mathrm{out}}^{(l)}} \to \mathbb{R}^{N_{\mathrm{out}}^{(l)}}$ is a non-linear activation function; note that $\mathbf{x}^{(l)} = \mathbf{y}^{(l-1)}$. In this work, we restrict ourselves to two of the most commonly used activation functions: $\mathrm{ReLU}(x) = \max(0, x)$ and $\mathrm{sigmoid}(x) = (1 + \exp(-x))^{-1}$; applied to multidimensional inputs, these must be intended as component-wise functions. Equation 1 can be trivially extended to the case in which $N_{\mathrm{s}}$ samples are processed at the same time; in such scenario, the input and output of the layer are matrices – that is $\mathbf{X}^{(l)} \in \mathbb{R}^{N_{\mathrm{s}} \times N_{\mathrm{in}}^{(l)}}$ and $\mathbf{Y}^{(l)} \in \mathbb{R}^{N_{\mathrm{s}} \times N_{\mathrm{out}}^{(l)}}$ – where each sample is stored row-wise.

We indicate $\theta^{(l)} = (W^{(l)}, \mathbf{b}^{(l)})$ the weights-biases pair of the $l^{\mathrm{th}}$ layer, and $\Theta = \{\theta^{(1)}, \ldots, \theta^{(L)}\}$ the set of all the parameters of a MLP. Moreover, we denote $N_{\mathrm{out}} = N_{\mathrm{out}}^{(L)}$ the size of the output of the neural network. The goal of the training process for a MLP is to find suitable values of the parameters $W^{(l)}$ and $\mathbf{b}^{(l)}$ for $l = 1, \ldots, L$ such that the function modeled by the network, $f : \mathbb{R}^{N_{\mathrm{in}}} \to \mathbb{R}^{N_{\mathrm{out}}}$ defined as $f(\mathbf{x}^{(0)}) := \mathbf{y}^{(L)}$, well approximates the actual input-output relationship. During the training process, the discrepancy between the real output $\mathbf{y}$ and the one of the network $\mathbf{y}^{(L)}$ is measured by a *loss function* $\mathcal{L}_{\Theta}(\mathbf{y}, \mathbf{y}^{(L)})$, which clearly depends on the parameters and is minimized during the training process. A common choice of loss function, which we also consider in this work, is given by the mean squared error (MSE). In the case of multiple samples, this is defined as

$$\mathrm{MSE}_{\Theta}(\mathbf{Y}, \mathbf{Y}^{(L)}) = \sum_{i=1}^{N_{\mathrm{s}}} \sum_{j=1}^{N_{\mathrm{out}}} \left(Y_{ij} - Y_{ij}^{(L)}\right)^2, \tag{2}$$

where $\mathbf{Y}$ is the matrix storing the desired outputs for all the samples. The minimization is typically performed via (stochastic) gradient descent or variations such as the Adam algorithm [26], which all require the computation of the gradients of the loss function with respect to the parameters. In modern DL frameworks, these quantities are automatically computed by *backpropagation*, which is a specialization of the chain rule for derivatives to the functionals modeled by such architectures.

# 3 Parametrized partial differential equations

In this section we introduce the general framework of parametrized PDEs and we describe the fundamental principles of the reduced basis method.

## 3.1 Parametrized PDEs

Let us consider a parameter space $\mathcal{D} \subset \mathbb{R}^p$, $p \geq 1$, and denote by $\boldsymbol{\mu} \in \mathcal{D}$ a parameter vector encoding physical and/or geometrical properties of the problem. Furthermore, let us introduce an open and bounded domain $\Omega \subset \mathbb{R}^d$, $d = 2, 3$, and denote by $\partial\Omega$ its boundary. We consider the following equation

$$\mathcal{N}[u; \boldsymbol{\mu}] = f(\boldsymbol{\mu}) \qquad \text{in } \Omega \tag{3}$$

where $\mathcal{N}[\cdot; \boldsymbol{\mu}]$ is a $\boldsymbol{\mu}$-differential operator which models a physical phenomenon. For any $\boldsymbol{\mu} \in \mathcal{D}$, we are interested in the solution $u = u(\boldsymbol{\mu})$. Equation (3) is provided with proper boundary conditions on $\partial\Omega$.

A numerical approximation to the solution $u$ of (3) can be obtained by means of a full order model (FOM). The latter is derived, e. g., starting from the variational formulation of (3) and employing a (Petrov-)Galerkin projection onto a finite-dimensional space. Examples of such FOMs are the spectral element method, the finite volume method and the finite element (FE) method. In our numerical experiments we consider the FE approximation as our reference *true* solution. Solving the FE problem is equivalent to solving a nonlinear algebraic system of (generally) large dimension

$$\mathbf{N}[\mathbf{u}, \boldsymbol{\mu}]\mathbf{u} = \mathbf{f}(\boldsymbol{\mu}), \tag{4}$$

where $\mathbf{N}[\mathbf{u}, \boldsymbol{\mu}] \in \mathbb{R}^{N_h \times N_h}$ is the parametrized nonlinear matrix, $\mathbf{f}(\boldsymbol{\mu}) \in \mathbb{R}^{N_h}$ the right hand side and $\mathbf{u} = \mathbf{u}(\boldsymbol{\mu}) \in \mathbb{R}^{N_h}$ is the FE vector containing the value of the solution at the mesh nodes. The dimension $N_h$ is in general very large when facing real-world applications, and can easily reach millions. When facing parameter dependent problem, i.e. when the solution of the PDE (3) is sought for many instances of the physical parameter $\boldsymbol{\mu}$, the computational load entailed by repeatedly solving the FE problem (4) can be unbearable. This issue can be overcome by Reduced Order Modelling (ROM) techniques, which allow to exploit the parameter dependence to boost the solution of the parametrized-problem at hand.

## 3.2 The reduced basis method for parametrized PDEs

Among all the ROM approaches, the RB method has been one of the most studied and exploited in the field of parametrized PDEs, and represents a convenient framework for the modeling reduction of differential problems. In the following, we briefly recall the basic ideas of the RB method, we refer to [37, 22] for full reviews and additional details.

The RB method relies on the idea that the $\boldsymbol{\mu}$-dependent solution of the $N_h \times N_h$ FE problem can be well approximated by a linear combination of $N$ basis functions obtained as solution of the same problem for (suitably chosen) parameter values, with $N \ll N_h$. The computation is usually based on an *offline/online* splitting. The offline phase computes a *reduced space* (the RB space), which is algebraically represented by the matrix $\mathbf{V} \in \mathbb{R}^{N_h \times N}$, $\mathbf{V} = [\boldsymbol{\xi}_1 | \ldots | \boldsymbol{\xi}_N]$. The goal of the RB method is to find a vector of weights $\mathbf{u}_N \in \mathbb{R}^N$ such that $\mathbf{V}\mathbf{u}_N$ is an accurate approximation

to the FE solution, that is $\mathbf{u} \approx \mathbf{V}\mathbf{u}_N$. This is done by solving an empirical algebraic RB problem of (small) dimension $N$ for any new instance of the parameter $\boldsymbol{\mu}$

$$\mathbf{N}_N[\mathbf{u}_N, \boldsymbol{\mu}]\mathbf{u}_N = \mathbf{f}_N(\boldsymbol{\mu}), \tag{5}$$

where $\mathbf{N}_N \in \mathbb{R}^{N \times N}$ is the RB matrix and $\mathbf{f}_N(\boldsymbol{\mu}) \in \mathbb{R}^N$ the RB right hand side, respectively obtained as Galerkin projection of the orginal FE arrays, where the nonlinear operator is computed at the approximation $\mathbf{V}\mathbf{u}_N$

$$\mathbf{N}_N[\mathbf{u}_N, \boldsymbol{\mu}] = \mathbf{V}^T\mathbf{N}[\mathbf{V}\mathbf{u}_N; \boldsymbol{\mu}]\mathbf{V}, \qquad \mathbf{f}_N(\boldsymbol{\mu}) = \mathbf{V}^T\mathbf{f}(\boldsymbol{\mu}). \tag{6}$$

We notice that Petrov-Galerkin projections are also viable ways to obtain a robust RB approximation, especially for noncoercive problems, see e.g. [1, 11].

### 3.2.1 Constructing the RB space

The RB matrix $\mathbf{V}$ can be obtained by employing either a greedy or proper orthogonal decomposition (POD) technique. We briefly remind the latter as we will employ it in the numerical experiments, for further reading we refer to [45]. Let us consider a set of $n_s$ parameter values $\{\boldsymbol{\mu}_i\}_{i=1}^{n_s}$ and FE vectors $\{\mathbf{u}(\boldsymbol{\mu}_i)\}_{i=1}^{n_s} \subset \mathbb{R}^{N_h}$ (called *snapshots*) collected as columns of a matrix $\mathbf{S} = [\mathbf{u}_1 | \dots | \mathbf{u}_{n_s}]$, $\mathbf{S} \in \mathbb{R}^{N_h \times n_s}$. For any prescribed dimension $N$, the POD allows to find an orthonormal basis $\{\boldsymbol{\xi}_i\}_{i=1}^{n_s}$ and the corresponding $N$-dimensional subspace, spanned by the columns of the matrix $\mathbf{V} = [\boldsymbol{\xi}_1 | \dots | \boldsymbol{\xi}_N]$ which best approximates $\{\mathbf{u}(\boldsymbol{\mu}_i)\}_{i=1}^{n_s}$ up to a tolerance $\varepsilon_{\text{POD}}$ with respect to the Euclidean norm. POD takes advantage of the singular value decomposition (SVD) of $\mathbf{S}$

$$\mathbf{S} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{Z}^T,$$

with $\mathbf{U} \in \mathbb{R}^{N_h \times N_h}$ and $\mathbf{Z} \in \mathbb{R}^{n_s \times n_s}$ orthogonal matrices and $\boldsymbol{\Sigma} = diag(\sigma_1, \dots \sigma_{n_s})$, $\boldsymbol{\Sigma} \in \mathbb{R}^{N_h \times n_s}$, containing the singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{n_s} \geq 0$. Then, $\mathbf{V}$ is provided by the first $N$ columns of $\mathbf{U}$, which form by construction an orthonormal basis for the best $N$-dimensional approximation subspace. In practical applications, $N$ is usually not set a priori by the user but for a fixed tolerance $\varepsilon_{\text{POD}}$, $N$ is found as the minimum $j$ such that

$$1 - \frac{\sum_{i=1}^{j} \sigma_i^2}{\sum_{i=1}^{n_s} \sigma_i^2} \leq \varepsilon_{\text{POD}}^2. \tag{7}$$

### 3.2.2 Affine decompositions of RB arrays

It is clear that the solution of the RB problem (5) yields a number of operations which depends on the number of RB degrees of freedom $N$, hence independent of the FE dimension $N_h$. However, the assembly (6) of problem (5) requires the projection of the FE matrix and right hand side on the RB space, thus entailing operations with a complexity dependent on $N_h$ and preventing an efficient *offline/online* splitting. To overcome this bottleneck, the RB method relies on the affine dependence of the RB arrays, that is $\mathbf{N}[\mathbf{u}, \boldsymbol{\mu}]$ and $\mathbf{f}(\boldsymbol{\mu})$ can be written as sum of $Q_n$ and $Q_f$ terms independent of $\mathbf{u}_N$ and $\boldsymbol{\mu}$, as follows

$$\mathbf{N}[\mathbf{V}\mathbf{u}_N, \boldsymbol{\mu}] = \sum_{q=1}^{Q_n} \theta_n^q(\mathbf{u}_N; \boldsymbol{\mu})\mathbf{N}^q, \qquad \mathbf{f}(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \theta_f^q(\boldsymbol{\mu})\mathbf{f}^q, \tag{8}$$

with $\mathbf{N}^q \in \mathbb{R}^{N_h \times N_h}$, $q = 1, \dots, Q_n$ and $\mathbf{f}^q \in \mathbb{R}^{N_h}$, $q = 1, \dots, Q_f$. From (8) it easily follows that

$$\mathbf{N}_N[\mathbf{u}_N, \boldsymbol{\mu}] = \sum_{q=1}^{Q_n} \theta_n^q(\mathbf{u}_N; \boldsymbol{\mu})\mathbf{N}_N^q, \qquad \mathbf{f}_N(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \theta_f^q(\boldsymbol{\mu})\mathbf{f}_N^q, \tag{9}$$

with $\mathbf{N}_N^q = \mathbf{V}^T \mathbf{N}^q \mathbf{V} \in \mathbb{R}^{N \times N}$, $q = 1, \ldots, Q_n$, and $\mathbf{f}_N^q = \mathbf{V}^T \mathbf{f}^q$, $q = 1, \ldots, Q_f$. Being the RB affine arrays in (8) $\mathbf{u}_N$- and $\boldsymbol{\mu}$-independent, they can be preassembled in the offline phase. Then the assembly of $\mathbf{N}_N[\mathbf{u}_N, \boldsymbol{\mu}]$ and $\mathbf{f}_N(\boldsymbol{\mu})$ in the online phase entails only the sums in (9), whose operations are $N_h$-independent.

When we are facing a PDE problem which does not feature by construction an affine decomposition as in (8), we cannot directly split the computation in an *offline* and *online* phase. However, an approximated affine decomposition can be computed by using the Empirical Interpolation Method (EIM) [3] or its discrete variants Discrete EIM (DEIM) and Matrix DEIM (MDEIM), where the latter is specific for matrices [7, 35]. During the offline phase, these algorithms are provided with a set of vector (DEIM) or matrix (MDEIM) snapshots and a tolerance which encodes the accuracy of the resulting affine approximation and return a basis of affine components which is computed through an internal POD. Then, during the online phase an interpolation problem is solved to compute the coefficients $\theta_n^q(\mathbf{u}_N; \boldsymbol{\mu})$ and $\theta_f^q(\boldsymbol{\mu})$ for each new instance of the PDE problem.

Depending on the problem at hand and the chosen accuracy, the number of affine components can largely vary, from tens to hundreds or even thousands, which may largely affect the speedup provided by the RB approximation compared with the FE one. This issue is prominent when facing nonlinear unsteady problems where the time-dynamics plays a relevant role for the creation of the RB approximation, see e.g. [12] in the case of the Navier-Stokes equations, where the number of affine components is large even for a relatively modest Reynolds number (about 400).

(M)(D)EIM techniques provide a satisfactory tool to deal with nonlinear and nonaffine PDE problems, however they present two significant bottlenecks:

- the underlying FE mesh must always be available both in the offline and online phase; in particular, in the latter the mesh is essential for assembling the right hand side of the interpolation problem for computing $\theta_n^q(\mathbf{u}_N; \boldsymbol{\mu})$ and $\theta_f^q(\boldsymbol{\mu})$. When the RB method is used to speedup the solution of large FE problems, the storage of the mesh can be significantly demanding and can prevent from using the RB approximation outside of a HPC environment.

- Given the (M)(D)EIM basis, the interpolation problem solved to produce the affine decomposition does not necessarily provides the best combination of basis functions, and in practice, for complex problems, this happens only if a significantly large number of basis is employed.

# 4  PDE-aware deep neural networks

In this section we introduce an original framework for approximating parametrized PDEs by coupling a DNN and a RB solver.

We focus on the following class of problems. Let us consider a physical system, mathematically modeled by a parametrized PDE as that in Equation (3), for which we know the value of the solution $u$ at certain points $\mathcal{P}_{\text{in}} = \{\vec{p}_i\}_{i=1}^{N_{\text{in}}}$; in realistic scenarios, these measurements could be obtained from multiple sensors positioned in $\Omega$ or on its boundary $\partial \Omega$. In this work, we restrict ourselves to steady physical phenomena, since unsteady cases need further developments. The PDE is parametrized by a vector of parameters $\boldsymbol{\mu} \in \mathcal{D} \subset \mathbb{R}^p$, $p \geq 1$ (following the notation introduced in Section 3.1). We assume that the values of the parameters are not known, or that their direct measurements is either complex, expensive or invasive. We want to define a neural network which, starting from measurements, is able not only to predict values on different parts of the domain, but also to provide the related parameter vector $\boldsymbol{\mu}$ or the coefficients of the affine decomposition of the reduced system in Equation (9). The train samples are generated by varying the parameters $\boldsymbol{\mu}$, whose value is not included in the dataset. We are interested in predicting the values of the solution at a set of points $\mathcal{Q}_{\text{out}} = \{\vec{q}_i\}_{i=1}^{N_{\text{out}}}$, given $u(\vec{p}_i; \boldsymbol{\mu})$, $\vec{p}_i \in \mathcal{P}_{\text{in}}$ for $i = 1, \ldots, N_{\text{in}}$, and at the same time to infer the value of the underlying parameter vector $\boldsymbol{\mu}$. In our numerical applications, we will explore the cases in which $\mathcal{P}_{\text{in}} \cap \mathcal{Q}_{\text{out}} = \emptyset$ and $\mathcal{P}_{\text{in}} \equiv \mathcal{Q}_{\text{out}}$.

We attempt the solution of such problems with PDE-aware deep neural networks (PDE-DNNs). The key idea of PDE-DNNs lies on the use of a PDE solver as building block of a DNN. This is
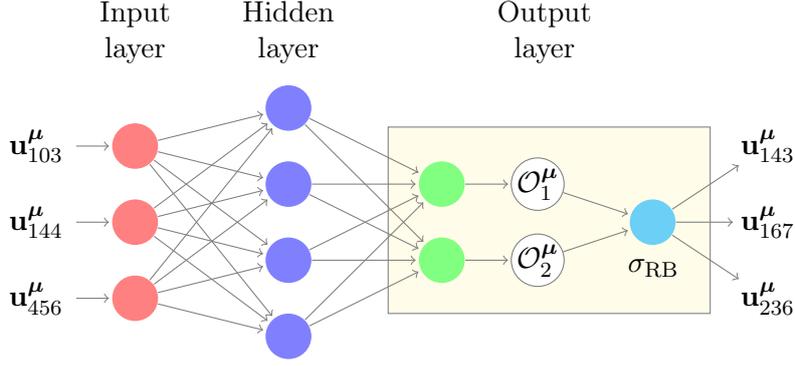
Figure 2: Example of PDE-DNN with one hidden layer in the MLP.

motivated by the fact that standard DNNs employed for solving PDE problems do not exploit the underlying physics, and, to the best of our knowledge, the PDE plays a relevant role only in the definition of the loss function, as in [38, 39]. In particular, our proposition is to consider a MLP with $L$ trainable layers in which the output layer encodes the discrete mathematical model.

The first $L - 1$ trainable layers of the network are mappings of the form of Equation (1) in which $\sigma^{(l)} = \text{ReLU}$ for $l = 1, \ldots, L - 1$. The mapping of the output layer takes the form

$$\mathbf{y}^{(L)} = \sigma_{\text{RB}}(W^{(L)}\mathbf{x}^{(L)} + \mathbf{b}^{(L)}), \tag{10}$$

where $\sigma_{\text{RB}} : \mathbb{R}^s \to \mathbb{R}^{N_{\text{out}}}$ represents the function of a RB solver acting from a $s$-dimensional representation of the solution space to the space of the solution values at the output locations $\mathcal{Q}_{\text{out}}$. The RB solver $\sigma_{\text{RB}}$ constructs such mapping as follows. Let us assume without loss of generality that $s = Q_n + Q_f$ and that $\boldsymbol{\xi} = [\theta_n^1, \ldots, \theta_n^{Q_n}, \theta_f^1, \ldots, \theta_f^{Q_f}] \in \mathbb{R}^s$ is the vector of the coefficients of the affine decomposition in Equation (9). As a matter of fact, in the case the PDE allow for an affine decomposition, $\boldsymbol{\xi}$ could contain the physical parameters themselves. Given $\boldsymbol{\xi}$, the reduced solution is simply obtained by solving the linear system efficiently assembled as in Equation (9); we remark that the matrices $\mathbf{N}_N^q$ and the vectors $\mathbf{f}_N^q$ appearing in Equation (9) must be computed a priori in the offline phase. Then, the desired output is found by projecting the reduced solution onto the FE space, i.e. $\mathbf{u} = \mathbf{V}\mathbf{u}_N$, and by interpolating the FE solution at the points of interest in $\mathcal{Q}_{\text{out}}$. In our tests, for simplicity, the input and the outputs are taken as values of the FE solution at the nodes of the mesh; since we employ standard Lagrangian basis functions, no actual interpolation of the solution is needed, and, eventually, we can write the PDE-DNN activation function as follows

$$\sigma_{\text{RB}}(\boldsymbol{\xi}) = \mathbf{R}_{\text{out}}^T \mathbf{V}\mathbf{u}_N(\boldsymbol{\xi}) = \mathbf{R}_{\text{out}}^T \mathbf{V}\mathbf{N}_N^{-1}(\boldsymbol{\xi})\mathbf{f}_N(\boldsymbol{\xi}), \tag{11}$$

where $\mathbf{R}_{\text{out}} \in \mathbb{R}^{N_h \times N_{\text{out}}}$ restricts the FE representation of the RB approximation to the output locations.

We remark that, as we observed in our numerical simulations, taking $\xi \in [0,1]^s$ considerably improves the convergence rate of the optimization through gradient descent. For this reason, in our numerical simulations we actually consider $\widetilde{\sigma}_{\text{RB}} = \sigma_{\text{RB}} \circ \text{sigmoid}$ as activation function in the output layer. The choices of ReLU and sigmoid as activation functions for the hidden layers and as preprocessing of the RB solver in the output layer respectively has been determined adequate after empirical observation. We aknowledge, however, that other architectures could lead to different and possibly better results than the ones reported in this paper.

Figure 2 shows the schematic representation of a PDE-DNN with two trainable layers ($L = 2$) and with $s = 2$.

**Remark 4.1.** *Activation functions play a relevant role when training a network, since both the forward and backward propagation stages depend on them. In particular, in the latter the derivatives of the activation function with respect to the latent space $\boldsymbol{\xi} = [\theta_n^1, \ldots, \theta_n^{Q_n}, \theta_f^1, \ldots, \theta_f^{Q_f}]$, that is $\frac{\partial \sigma_{\text{RB}}}{\partial \theta_n^q}$*

and $\frac{\partial \sigma_{\mathrm{RB}}}{\partial \theta_f^q}$, enter into play. Following the definition of the affine decompositions (9), we can compute them as follows

$$\frac{\partial \sigma_{\mathrm{RB}}}{\partial \theta_n^q} = \mathbf{R}_{\mathrm{out}}^T \mathbf{V} \frac{\partial \left( \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{f}_N(\boldsymbol{\xi}) \right)}{\partial \theta_n^q} = \mathbf{R}_{\mathrm{out}}^T \mathbf{V} \frac{\partial (\mathbf{N}_N^{-1}(\boldsymbol{\xi}))}{\partial \theta_n^q} \mathbf{f}_N(\boldsymbol{\xi})$$

$$= -\mathbf{R}_{\mathrm{out}}^T \mathbf{V} \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \frac{\partial (\mathbf{N}_N(\boldsymbol{\xi}))}{\partial \theta_n^q} \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{f}_N = -\mathbf{R}_{\mathrm{out}}^T \mathbf{V} \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{N}_N^q \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{f}_N(\boldsymbol{\xi}), \quad (12)$$

for $q = 1, \ldots, Q_n$, and, for $q = 1, \ldots, Q_f$,

$$\frac{\partial \sigma_{\mathrm{RB}}}{\partial \theta_f^q} = \mathbf{R}_{\mathrm{out}}^T \mathbf{V} \frac{\partial \left( \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{f}_N(\boldsymbol{\xi}) \right)}{\partial \theta_f^q} = \mathbf{R}_{\mathrm{out}}^T \mathbf{V} \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \frac{\partial \left( \mathbf{f}_N(\boldsymbol{\xi}) \right)}{\partial \theta_f^q} = \mathbf{R}_{\mathrm{out}}^T \mathbf{V} \mathbf{N}_N^{-1}(\boldsymbol{\xi}) \mathbf{f}_N^q.$$

Let us assume that we are able to build, either from direct measurements or – as we consider in this paper – from numerical simulations, a training dataset composed of $N_s$ input vectors of the form $\mathbf{x}_i = [u(\vec{p}_1; \boldsymbol{\mu}_i), \ldots, u(\vec{p}_{N_{\mathrm{in}}}; \boldsymbol{\mu}_i)]$ and $N_s$ output vectors of the form $\mathbf{y}_i = [u(\vec{q}_1; \boldsymbol{\mu}_i), \ldots, u(\vec{q}_{N_{\mathrm{out}}}; \boldsymbol{\mu}_i)]$, possibly arranged in matrix form in $\mathbf{X} \in \mathbb{R}^{N_s \times N_{\mathrm{in}}}$ and $\mathbf{Y} \in \mathbb{R}^{N_s \times N_{\mathrm{out}}}$ respectively; the parameters $\boldsymbol{\mu}_i \in \mathcal{D}$, for $i = 1, \ldots, N_s$, are sampled in $\mathcal{D}$ to well represent the whole set. In this regard, we note that in the case the training dataset be obtained from numerical simulation, the computational burden of the RB offline phase is partially mitigated by the fact that the numerical solutions corresponding to each $\boldsymbol{\mu}_i$ for $i = 1, \ldots, N_s$ can be used as snapshots to build the RB basis by POD. For the training of the network, we consider the MSE loss function defined as in Equation (2). For each sample, our loss function does not consider at all the values of the corresponding parameters $\boldsymbol{\mu}_i$. However, $\boldsymbol{\mu}_i$ or the coefficients of the affine decomposition are internally estimated by the MLP. Such additional information is extracted in our code by combining in the output tensor the solution by the RB solver with the result of the linear operation in the output layer (i.e. the input required by $\sigma_{\mathrm{RB}}$). We stress that the latter is not used in the computation of the loss function.

## 5 Numerical experiments

### 5.1 Affinely parametrized elliptic problems

Let us consider the domain $\Omega = (0,1)^3$ and the parametrized advection-diffusion problem

$$\begin{cases} -\nabla \cdot (\nu \nabla u) + \vec{b}(\vec{x}; \alpha) \cdot \nabla u = 0 & \text{in } \Omega, \\ +\text{b.c.} \end{cases} \tag{13}$$

with parameters $\boldsymbol{\mu} = (\nu, \alpha) \in \mathcal{D} = (0.5, 10) \times (0, \pi/6) \subset \mathbb{R}^2$. The former defines the diffusivity of the problem and the latter is the angle defining the direction of the advection field $\vec{b}(\vec{x}; \alpha) = \sin(\alpha)\vec{b}_1(\vec{x}) + \cos(\alpha)\vec{b}_2(\vec{x})$, which is obtained as linear combination of two divergence free $\alpha$-independent vector fields $\vec{b}_1(\vec{x})$ and $\vec{b}_2(\vec{x})$. These are computed as solutions of the Navier-Stokes equations with viscosity 1 and equipped with non-homogeneous Dirichlet boundary conditions on one face of the cube – specifically, $\vec{b}_1(\vec{x}) = [100, 0, 0]$ and $\vec{b}_2(\vec{x}) = [0, 100, 0]$ for every $\vec{x} \in \partial\Omega$ having $x_3 = 1$ – and homogeneous Dirichlet boundary conditions everywhere else. Fig. 3 (left) shows the streamlines of the vector field obtained by choosing $\alpha = 0.48$. The boundary of the domain $\partial\Omega$ is partitioned in $\Gamma_D = \{\vec{x} \in \partial\Omega : x_1 = 0 \text{ or } x_1 = 1\}$ and $\Gamma_N$, such that $\partial\Omega = \Gamma_D \cup \Gamma_N$ with $\mathring{\Gamma}_D \cap \mathring{\Gamma}_N = \emptyset$. We impose non-homogeneous Dirichlet boundary conditions on $\Gamma_D$, in particular $u = 1$ for every $\vec{x} \in \Gamma_D$ with $x_1 = 1$ and $u = 0$ otherwise, and homogeneous Neumann conditions on $\Gamma_N$. Such problem could model, for instance, the temperature of a fluid in steady regime in which the values at the two Dirichlet boundaries are kept fixed and heat is freely exchanged at the Neumann boundaries. Fig. 3 (left) shows the solution of Eq. (13) over the streamlines of the advecting vector field.
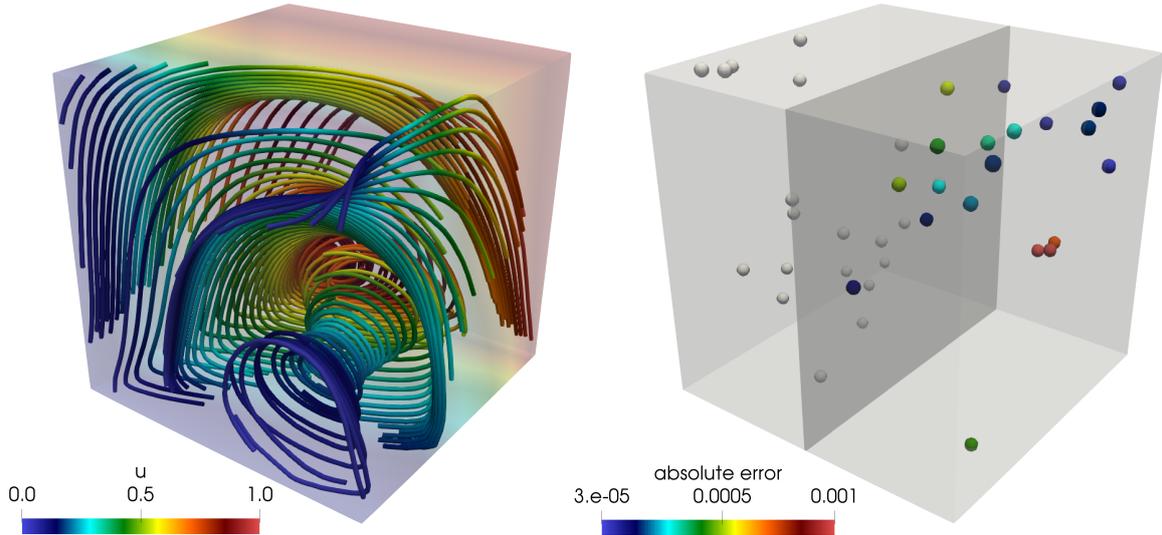
Figure 3: On the left, streamlines of the vector field $\vec{b}(\vec{x}; \alpha)$, colored with respect to the solution of Eq. 13. On the right, absolute error on the prediction of a PDE-DNN with 20 input points (white spheres) and 20 output points (colored spheres). Both figures refer to $\nu = 2.22$ and $\alpha = 0.48$.

After the discretization of problem 13 by the FE method (which we consider for all of the high fidelity models presented in this paper), the system is rewritten in algebraic form as

$$\mathbf{Au} = \mathbf{F}, \tag{14}$$

where $\mathbf{A} = \nu \mathbf{A}_1 + \sin(\alpha)\mathbf{A}_2 + \cos(\alpha)\mathbf{A}_3 \in \mathbb{R}^{N_h \times N_h}$ is the affine decomposition of the finite element matrix; $\mathbf{F} \in \mathbb{R}^{N_h}$ accounts for the boundary conditions and admits a similar affine decomposition.

### 5.1.1 PDE-DNNs for affinely parametrized PDEs

We aim at evaluating the performance of a PDE-DNN in predicting solutions of Eq. (13) for several instances of $\boldsymbol{\mu}$ in a set of output points in terms of the following *hyperparameters*: *i)* the number of input points $N_{\mathrm{in}}$, *ii)* the number of output points $N_{\mathrm{out}}$, *iii)* the tolerance for the POD of the RB solver $\varepsilon_{\mathrm{POD}}$, and *iv)* the number of samples in the training dataset $N_{\mathrm{s}}$.

For this test case we randomly sample the input and output points in the partitions of $\Omega$ corresponding to $x_2 \leq 1/2$ and $x_2 > 1/2$ respectively, as depicted in Fig. 13 (right); this choice is motivated by the fact that, if two sets $\mathcal{P}_{\mathrm{in}}$ and $\mathcal{Q}_{\mathrm{out}}$ were obtained from sampling over the entire $\Omega$, then also simple interpolation between points in the input set could potentially provide a rough approximation of the solution in the output points. As mentioned in Section 4, the sampling is done over the nodes of the mesh which is used in the computations by taking care of not selecting nodes on the Dirichlet boundaries; we consider a computational mesh composed of tetrahedra with $\mathbb{P}^1$ finite elements and a total of $N_h = 12416$ degrees of freedom.

In order to build a PDE-DNN, we first need to perform the offline phase of the RB method. This entails the collection of a predetermined number of snapshots $n_s$, which are obtained by solving the problem by FOM – in our case, the FE method with linear Lagrangian basis functions – for $\boldsymbol{\mu}_i \in \mathcal{D}$, $i = 1, \ldots, n_s$. In this example, $n_s = 350$ and the parameters are sampled from a uniform distribution in $\mathcal{D}$. Performing the POD of the snapshots matrix $\mathbf{S}$ with tolerances $\varepsilon_{\mathrm{POD}} = 10^{-4}$, $\varepsilon_{\mathrm{POD}} = 10^{-5}$, $\varepsilon_{\mathrm{POD}} = 10^{-6}$ and $\varepsilon_{\mathrm{POD}} = 10^{-7}$ yields bases of sizes $N = 12$, $N = 19$, $N = 28$ and $N = 38$ respectively.

The $n_s = 350$ snapshots generated during the RB offline phase are also used as samples in the training dataset, which is performed through Adam algorithm with learning rate $10^{-3}$. We decide to enlarge such dataset by exploiting the ability of the RB method to generate solutions corresponding to new parameters in $\mathcal{D}$ at a negligible cost. The largest training dataset we consider
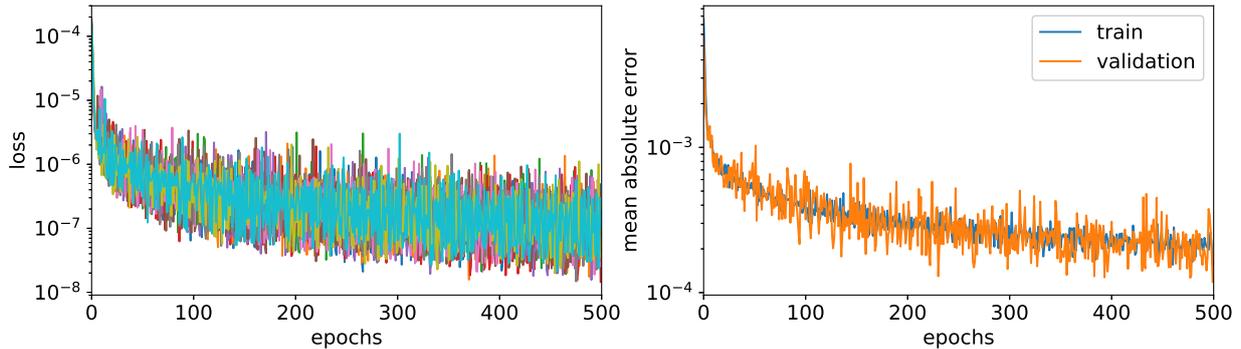
Figure 4: On the left, decaying of the MSE function during 10 trainings of PDE-DNN model with choice of hyperparameters $N_{\text{in}} = 20$, $N_{\text{out}} = 20$, $\varepsilon_{\text{POD}} = 10^{-5}$ and $N_{\text{s}} = 10000$. On average, the training has taken 1366 s. On the right, decaying of the mean absolute of error (averaged on the 10 runs) for train and validation datasets. The validation is obtained from the training dataset by selecting 20% of the samples.

is composed of 20000 samples: 350 full order solutions and 19650 RB solutions obtained with POD tolerance $\varepsilon_{\text{POD}} = 10^{-7}$. As we wish to investigate the role of the number of samples on the training process and performance of the network, we also consider smaller training datasets by selecting the first $N_{\text{s}}$ samples of the largest training dataset. We remark that in our tests we actually select 20% of the training samples for the validation of the model, that is, the training process is based on the remaining 80% samples and the loss and other performance metrics are computed also on the validation dataset. In the deep learning community, such practice is adopted to prevent common issues such as *over-* or *underfitting*; we refer the reader to [17] for more information on these topics. In the following, we will refer to training dataset as the subset of the total number of samples that is used during training. The test dataset is composed of 600 FE solutions at different random parameter values.

The PDE-DNN consists of 4 hidden layers. In this test case, each hidden layer consists of 256 neurons. The linear part of the output layer is in charge of estimating a two dimensional vector, which is interpreted by the RB solver as an estimation of the physical parameters of the problem. The training of the PDE-DNN is performed through minimization of the loss in Eq. (2) by gradient descent for 500 epochs; such number has been determined as adequate after empirical observations. Fig. 4 depicts the decaying of the MSE loss over 10 trainings of a PDE-DNN; we remark that the evolution of the loss function is not deterministic because the initialization of the trainable parameters of the MLP is random. In Fig. 4 (right), we focus on the mean absolute error – averaged over the 10 runs – as metric of the performance of the network. Since the mean absolute error reaches a plateau at around 500 epochs, we safely conclude that underfitting does not affect the training and this number of epochs is sufficient. Moreover, as the mean average error on the training and the validation dataset follow the same trend, overfitting does not occur during training. This is also true for the other choices of hyperparameters that we have considered. Our interpretation of this pheonomenon is that, as the train samples are obtained from numerical simulations, there are no "hidden parameters" that could characterize only the training dataset and penalize the performance over the validation. In other words, well representing the training data automatically leads to good approximation of the validation dataset.

Fig. 5 displays the estimation of the normalized parameters $\nu$ and $\alpha$ performed by the PDE-NN for different values of the POD tolerance and different number of training samples. The findings are obtained by averaging the performance of 30 networks (trained over 500 epochs) on the test dataset. The effect of the two hyperparameters $\varepsilon_{\text{POD}}$ and $N_{\text{s}}$ is the following. Given a sufficiently high number of input samples, for instance $N_{\text{s}} = 20000$, the RB tolerance does not significantly affect the accuracy of the estimation for neither of the two physical parameters (the tolerance regions corresponding to $N_{\text{s}} = 20000$, indeed, do not change significantly across the two rows of
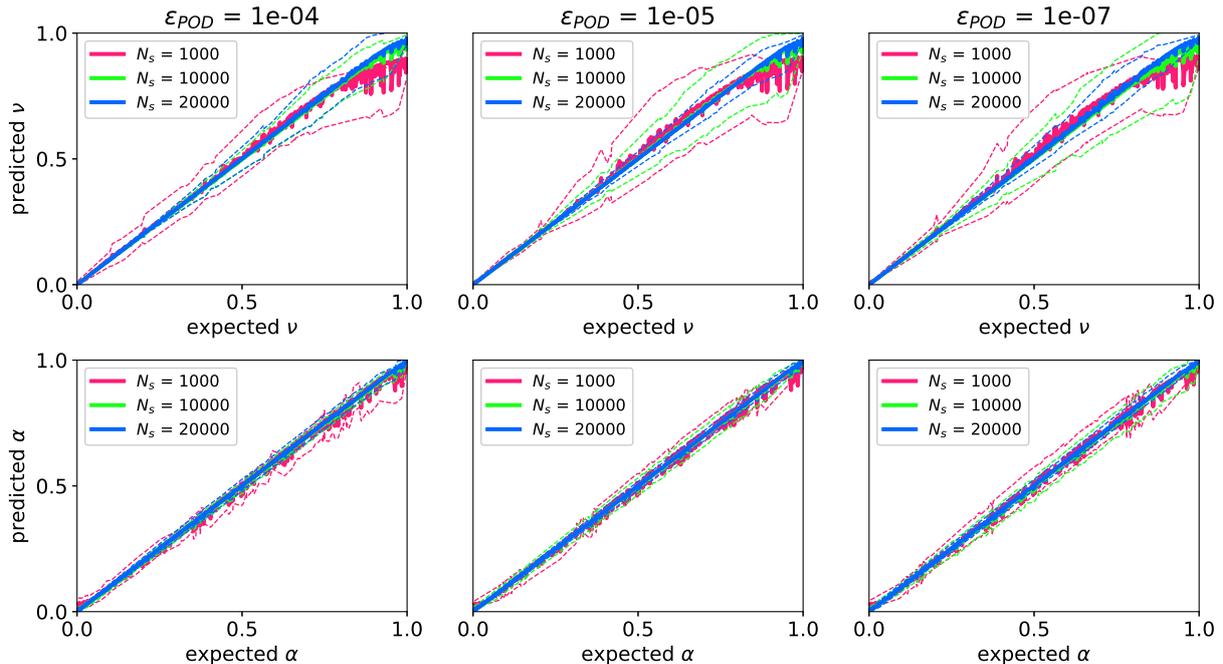
Figure 5: Estimation of the two (normalized over the respective ranges) physical parameters of the problem, $\nu$ (top row) and $\alpha$ (bottom row), for different POD tolerances of the RB problem and different number of samples. The results refer to the test dataset, which is composed of 600 data points, and are average over 30 trainings of the PDE-DNN network. The regions between dashed lines display the 95% confidence intervals corresponding to the average lines with the matching colors.

plots of Fig. 5). If, however, the training dataset is scarce, decreasing the tolerance of the POD leads to more uncertain results (i.e. larger confidence regions) without an appreciable increase of the precision on average. In other words, based on these results we conclude that the choice $\varepsilon_{\text{POD}} = 10^{-4}$ or $\varepsilon_{\text{POD}} = 10^{-5}$ is sufficient to get accurate results regardless of the size of the training dataset.

The estimation of the diffusion coefficient $\nu$ appears more challenging than the one of the angle $\alpha$: for values approaching $\nu = 10$ the networks are not able to distinguish accurately diffusivity coefficients close to each other, and in particular the predicted values are on average underestimating the correct ones (this is particularly true for $N_{\text{s}} = 1000$). The estimation of the parameter $\alpha$ is quite accurate for all the considered configurations.

Table 1 and 2 show the errors on the prediction of $\nu$ and $\alpha$ respectively for different configurations of the PDE-DNN, averaged over 10 trainings of the network; all the results are obtained with $\varepsilon_{\text{POD}} = 10^{-5}$.

As a comparison, we consider three neural networks, denoted $\text{MLP}_\mu$, $\text{MLP}_{\text{out}}$ and MLP. These are are designed and trained to compute the parameter $\boldsymbol{\mu}$ ($\text{MLP}_\mu$), or the value of $u$ at the given output locations ($\text{MLP}_{\text{out}}$) or both (MLP). The hidden layers of these networks are the same as those in the considered PDE-DNN, whereas the output layer is a perceptron equipped with sigmoid and with dimension $\dim(\boldsymbol{\mu})$, $N_{\text{out}}$ and $N_{\text{out}} + \dim(\boldsymbol{\mu})$ for $\text{MLP}_\mu$, $\text{MLP}_{\text{out}}$ and MLP, respectively.

As we noted in Fig. 5, the number of samples plays a relevant role in the approximation of the physical parameters, whereas the number of input and output locations does not affect significantly the results. We remark that the PDE-DNN is able to obtain the approximately same accuracy of the other two neural networks, even though the parameters are only found as a byproduct (i.e. it is not necessary to train the model by providing the value of the parameters).

Table 3 shows the error on the output normalized with respect to the norm of the expected output (to take into account the different output sizes) in the same configurations considered in

|  | $N_{\mathrm{s}}$ | $(N_{\mathrm{in}}, N_{\mathrm{out}})$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
|  |  | $(20, 20)$ | $(40, 40)$ | $(100, 100)$ | $(20, 40)$ | $(20, 100)$ | $(40, 20)$ | $(100, 20)$ |
| PDE-NN | 1000 | 2.18e-02 | 2.64e-02 | 2.00e-02 | 2.03e-02 | 2.31e-02 | 2.29e-02 | 7.05e-02 |
|  | 10000 | 3.62e-03 | 5.28e-03 | 5.00e-03 | 1.11e-02 | 9.99e-03 | 4.60e-03 | 3.87e-03 |
|  | 20000 | 2.87e-03 | 7.33e-03 | 5.69e-03 | 2.24e-03 | 4.23e-03 | 4.81e-03 | 3.50e-03 |
| MLP | 1000 | 2.37e-02 | 1.91e-02 | 1.78e-02 | 2.24e-02 | 1.62e-02 | 1.75e-02 | 1.56e-02 |
|  | 10000 | 5.81e-03 | 7.49e-03 | 3.69e-03 | 6.84e-03 | 4.46e-03 | 4.70e-03 | 5.32e-03 |
|  | 20000 | 4.71e-03 | 3.96e-03 | 3.53e-03 | 4.17e-03 | 3.55e-03 | 4.07e-03 | 4.08e-03 |
| $\mathrm{MLP}_\mu$ | 1000 | 2.39e-02 | 2.51e-02 | 3.04e-02 | 1.87e-02 | 2.01e-02 | 2.61e-02 | 2.02e-02 |
|  | 10000 | 7.99e-03 | 4.13e-03 | 5.80e-03 | 4.93e-03 | 6.50e-03 | 9.62e-03 | 6.12e-03 |
|  | 20000 | 2.97e-03 | 1.30e-02 | 2.82e-03 | 5.56e-03 | 8.53e-03 | 5.95e-03 | 3.17e-03 |

Table 1: Average error on $\nu$ over 600 test samples and over the outputs of 10 independently trained networks.

|  | $N_{\mathrm{s}}$ | $(N_{\mathrm{in}}, N_{\mathrm{out}})$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
|  |  | $(20, 20)$ | $(40, 40)$ | $(100, 100)$ | $(20, 40)$ | $(20, 100)$ | $(40, 20)$ | $(100, 20)$ |
| PDE-NN | 1000 | 8.42e-03 | 1.58e-02 | 9.16e-03 | 9.52e-03 | 1.09e-02 | 1.30e-02 | 3.93e-02 |
|  | 10000 | 3.05e-03 | 2.36e-03 | 2.91e-03 | 4.71e-03 | 3.30e-03 | 1.69e-03 | 2.07e-03 |
|  | 20000 | 2.66e-03 | 4.28e-03 | 2.30e-03 | 1.02e-03 | 1.71e-03 | 4.54e-03 | 1.41e-03 |
| MLP | 1000 | 8.44e-03 | 7.62e-03 | 1.46e-02 | 1.02e-02 | 8.10e-03 | 7.03e-03 | 8.61e-03 |
|  | 10000 | 7.10e-03 | 3.02e-03 | 4.49e-03 | 4.22e-03 | 1.81e-03 | 4.41e-03 | 2.43e-03 |
|  | 20000 | 2.31e-03 | 1.67e-03 | 1.34e-03 | 1.85e-03 | 2.46e-03 | 2.27e-03 | 1.85e-03 |
| $\mathrm{MLP}_\mu$ | 1000 | 1.57e-02 | 1.07e-02 | 2.20e-02 | 8.79e-03 | 1.04e-02 | 1.21e-02 | 1.49e-02 |
|  | 10000 | 3.03e-03 | 1.95e-03 | 4.83e-03 | 2.08e-03 | 2.09e-03 | 2.67e-03 | 6.47e-03 |
|  | 20000 | 2.09e-03 | 3.64e-03 | 1.96e-03 | 2.11e-03 | 3.88e-03 | 4.26e-03 | 3.03e-03 |

Table 2: Average error on $\alpha$ over 600 test samples and over the outputs of 10 independently trained networks.

| | $N_{\mathrm{s}}$ | $(N_{\mathrm{in}}, N_{\mathrm{out}})$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | $(20,20)$ | $(40,40)$ | $(100,100)$ | $(20,40)$ | $(20,100)$ | $(40,20)$ | $(100,20)$ |
| PDE-NN | 1000 | 7.59e-04 | 8.72e-04 | 5.75e-04 | 5.74e-04 | 6.35e-04 | 1.57e-03 | 4.84e-03 |
| | 10000 | 2.08e-04 | 2.70e-04 | 3.77e-04 | 3.18e-04 | 3.30e-04 | 1.80e-04 | 1.28e-04 |
| | 20000 | 1.63e-04 | 2.21e-04 | 1.78e-04 | 7.88e-05 | 1.24e-04 | 2.66e-04 | 1.40e-04 |
| MLP | 1000 | 1.54e-03 | 1.37e-03 | 2.56e-03 | 1.17e-03 | 9.85e-04 | 1.03e-03 | 1.41e-03 |
| | 10000 | 8.93e-04 | 5.58e-04 | 5.34e-04 | 5.47e-04 | 3.48e-04 | 5.66e-04 | 4.34e-04 |
| | 20000 | 2.98e-04 | 3.39e-04 | 2.89e-04 | 2.43e-04 | 2.28e-04 | 3.51e-04 | 3.93e-04 |
| $\mathrm{MLP_{out}}$ | 1000 | 1.05e-03 | 8.47e-04 | 1.45e-03 | 1.01e-03 | 1.34e-03 | 1.22e-03 | 1.12e-03 |
| | 10000 | 5.99e-04 | 4.51e-04 | 1.81e-04 | 2.82e-04 | 1.56e-04 | 2.80e-04 | 1.99e-04 |
| | 20000 | 1.63e-04 | 8.50e-05 | 1.78e-04 | 3.68e-04 | 1.33e-04 | 1.37e-04 | 1.11e-04 |

Table 3: Average error – normalized with respect to the output norm – on $\mathbf{y}_i$ over 600 test samples and over the outputs of 10 independently trained networks.



Figure 6: $H^1$ error against the 600 full order solutions in the test dataset when employing 20, 40 and 100 input locations in the dataset. Each point has $x$-component equal to the estimated value of the normalized physical parameter $\nu$ and it is colored with respect to the estimated value of the normalized physical parameter $\beta$; these and the errors have been obtained as averages over 10 trainings (500 epochs) of PDE-DNNs corresponding to $\varepsilon_{\mathrm{POD}} = 10^{-5}$ and $Ns = 20000$.

the previous tables. Also in the case of $\mathrm{MLP_{out}}$, the three networks perform similarly for any choice of the hyperparameters.

Finally, we consider the case in which the physical input and output location coincide. In other words, we train the network such that the RB solution coming from the output layer is as close as possible to the full order solution at the input locations($x \leq 1/2$). In this scenario, the PDE-NN is trained to model the identity function and the architecture is effectively an autoencoder in which the encoder is the MLP, the decoder is the RB solver, and the latent space is the space of physical parameters. In order to evaluate the performance of the network, we consider as metrics the $H^1$ error against the full order solution. We recall that this is possible because, given the estimation of the physical parameters by the MLP, the RB solver is able to approximate the solution at any point in the domain.

Fig. 6 shows the $H^1$ error averaged over 10 trainings of PDE-DNNs when employing 20, 40 and 100 physical points. The average $H^1$ errors over the whole test dataset in the three cases is $3.0 \times 10^{-4}$, $2.3 \times 10^{-4}$ and $2.4 \times 10^{-4}$. As a comparison, we report that the average error achieved with the RB method alone with the same POD accuracy $\varepsilon_{\mathrm{POD}} = 10^{-5}$ is $1.1 \times 10^{-5}$. We observe that increasing the number of input locations has the effect of slightly decreasing the average error
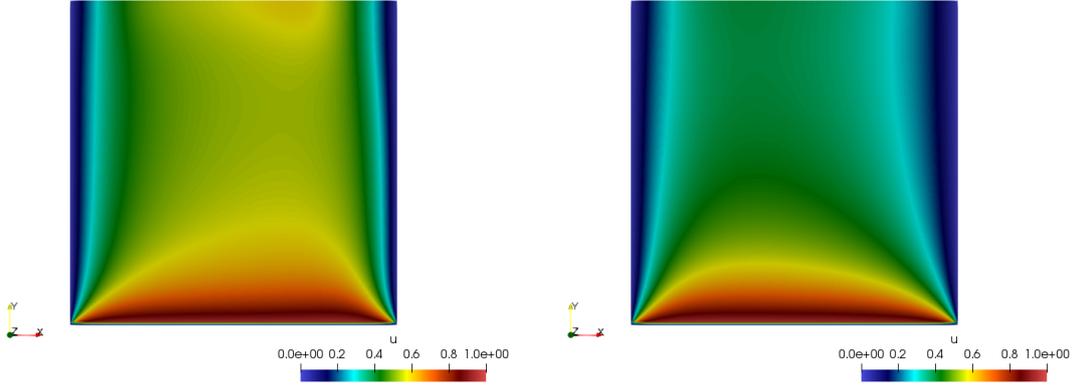
Figure 7: Examples of solutions of problem (15) for the parameters $\boldsymbol{\mu}_1 = (0.42, 0.42, 0.06)$ (left) and $\boldsymbol{\mu}_2 = (0.45, 0.59, 0.09)$ (right).

performed over the test dataset. Moreover, Fig. 6 highlights that the prediction of the networks is considerably worse when the normalized $\nu$ approaches 0. This phenomenon and the fact that the loss function we consider is only concerned with the accuracy on the solution could explain the tendency of the networks to be considerably more precise in the estimation of $\nu$ when it is small (as displayed in Fig 5).

## 5.2 Nonaffinely parametrized elliptic problems

In the second test case, we consider a second-order diffusion problem where the diffusion coefficient is nonaffinely parametrized. Let us consider the domain $\Omega = (0,1)^2$ describing a square beam and a diffusion problem of the form

$$\begin{cases} -\nabla \cdot (\alpha(\boldsymbol{\mu})\nabla T) = 1 & \text{in } \Omega, \\ +\text{b.c.} \end{cases} \tag{15}$$

where $T$ is the temperature of the beam. We define the following boundaries

$$\Gamma_N = \left\{ \vec{x} \in \bar{\Omega} : y = 1 \right\}, \qquad \Gamma_D^1 = \left\{ \vec{x} \in \bar{\Omega} : y = 0 \right\}, \qquad \Gamma_D^0 = \partial\Omega \backslash \Gamma_N \backslash \Gamma_D^1,$$

such that $\partial\Omega = \bar{\Gamma}_0^D \cup \bar{\Gamma}_1^D \cup \bar{\Gamma}_N$, and we set homogeneous Neumann boundary conditions on $\Gamma_N$, homogeneous Dirichlet conditions on $\Gamma_D^0$ and the solution to be equal to 1 on $\Gamma_D^1$. We introduce the parameter vector

$$\boldsymbol{\mu} = (x_0^1, x_0^2, \sigma) \in [0.4, 0.6]^2 \times [0.05, 0.1] \subset \mathbb{R}^3$$

and the coefficients

$$\alpha(\boldsymbol{\mu}) = \sigma + \frac{1}{\sigma} \exp\left( -\frac{\|\vec{x} - \vec{x}_0(\boldsymbol{\mu})\|^2}{\sigma} \right), \qquad \vec{x}_0(\boldsymbol{\mu}) = (x_0^1, x_0^2).$$

Examples of solutions for two different physical parameters $\boldsymbol{\mu}$ are reported in Figure 7.

We employ a lifting function $T_l = T_l(\vec{x})$ to deal with the nonhomogeneous Dirichlet boundary conditions, and we use a structured mesh with $N_h = 10201$ vertices and the FE method with first order polynomial piecewise basis functions to discretize the variational problem corresponding to (15). This leads to the following parametrized linear system of dimension $N_h$

$$\mathbf{A}(\boldsymbol{\mu})\mathbf{T} = \mathbf{F}(\boldsymbol{\mu}), \tag{16}$$

14

| MLP architecture | layers description |
|:---:|:---:|
| $\mathcal{A}_{\text{na}}^1$ | 4 hidden layers with sizes 1024, 512, 256, 128 |
| $\mathcal{A}_{\text{na}}^2$ | 4 hidden layers with all sizes equal to 256 |
| $\mathcal{A}_{\text{na}}^3$ | 4 hidden layers with all sizes equal to 64 |

Table 4: Architectures employed for Nonaffine test case.

where $\mathbf{A}(\boldsymbol{\mu}) \in \mathbb{R}^{N_h \times N_h}$ and $\mathbf{T} = \mathbf{T}(\boldsymbol{\mu})$, $\mathbf{F}(\boldsymbol{\mu}) \in \mathbb{R}^{N_h}$. A significant difference with respect to the previous test case is the nonaffine dependence of the FE matrix $\mathbf{A}(\boldsymbol{\mu})$ and right hand side $\mathbf{F}(\boldsymbol{\mu})$ (the latter due to the lifting function) with respect to the parameters, i.e. assumptions (8) do not hold. Hence, to build a RB model independent of dimension $N_h$, we employ MDEIM and DEIM leading to $Q_a$ sand $Q_f$ affine basis for the FE matrix and right hand side, respectively, that is

$$\mathbf{A}(\boldsymbol{\mu}) \approx \sum_{q=1}^{Q_a} \theta_q^a(\boldsymbol{\mu}) \mathbf{A}_q, \qquad \mathbf{F}(\boldsymbol{\mu}) \approx \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) \mathbf{F}_q. \qquad (17)$$

We recall that given a new parameter $\boldsymbol{\mu}$, the coefficients $\{\theta_q^a(\boldsymbol{\mu})\}_{q=1}^{Q_a}$ and $\{\theta_q^f(\boldsymbol{\mu})\}_{q=1}^{Q_f}$ are computed by solving an interpolation problem which enforces the true values of some preselected entries of the matrix (or right hand side) in its approximation. Clearly, the higher the number of affine components $Q_a$ and $Q_f$, the more accurate the resulting affine approximation and the corresponding RB model. Notice that the underlying parametrization of $\alpha(\boldsymbol{\mu})$ results in a complex parameter dependence, especially related to the approximate affine approximation of $\mathbf{A}(\boldsymbol{\mu})$. This is due to the Gaussian-like shape of the diffusion coefficient, which changes its centre and amplitude according to the value of the parameter $\boldsymbol{\mu}$ and requires MDEIM to compute a large number of basis functions to obtain an accurate affine approximation. This fact can hamper the online efficiency of the RB method and lead to a less competitive approximation with respect to compute the full FE solution. Finally, we highlight that being the problem linear, the coefficients $\theta_q^a(\boldsymbol{\mu})$ do not depend on $\mathbf{u}_N(\boldsymbol{\mu})$.

### 5.2.1 PDE-DNNs for nonaffinely parametrized PDEs

Given the measurements of the FE solution at $N_{\text{in}} = 40$ physical points located on the boundary $\Gamma_N$ (green in Figure 8a), we are interested in predicting the value of the same solution at $N_{\text{out}} = 100$ physical points located in the top left quadrant of the domain (red in Figure 8a). To this aim, we employ a PDE-DNN with 4 hidden layers, for which we investigate three different architectures, reported in Table 4. The fifth and last layer combines the outputs of the forth layer into the coefficients $\{\theta_q^a(\boldsymbol{\mu})\}_{q=1}^{Q_a}$ and $\{\theta_q^f(\boldsymbol{\mu})\}_{q=1}^{Q_f}$ of the affine approximations (17). These coefficients are used to assemble the RB matrix and right hand side. After solving the small dense linear problem, the layer recovers the values of the temperature at the output locations.

With respect to the PDE-DNN in Figure 2, $\mathcal{O}_i^{\boldsymbol{\mu}}$ represent the following values:

$$\mathcal{O}_i^{\boldsymbol{\mu}} = \begin{cases} \theta_i^a(\boldsymbol{\mu}) & \text{if } i = 1, \ldots, Q_a \\ \theta_{i-Q_a}^f(\boldsymbol{\mu}) & \text{if } i = Q_a + 1, \ldots, Q_a + Q_f. \end{cases}$$

The number of employed affine components is fixed to $Q_f = 10$ for the right hand side, yielding an accuracy of $\approx 10^{-6}$ for the approximation of the right hand side. On the other hand, the considered parametrization leads to a large number of terms for the affine approximation of $\mathbf{A}(\boldsymbol{\mu})$ to obtain its accurate representation (up to about 40 for an accuracy of $\approx 10^{-5}$), leading to a significant overhead when employing the standard RB method. In the latter, the values $\{\theta_i^a(\boldsymbol{\mu})\}_{i=1}^{Q_a}$ are obtained by solving an interpolation problem, and a large number of affine components is required to obtain an accurate RB solution. In the PDE-NN framework we are proposing, the values of the functions $\{\theta_i^a(\boldsymbol{\mu})\}_{i=1}^{Q_a}$ (and $\{\theta_i^f(\boldsymbol{\mu})\}_{i=1}^{Q_f}$) are instead optimized by the encoding part of the network, the MLP.

(a) Input location (green) and output location (red)

(b) $\mathcal{A}_{\mathrm{na}}^1$ and $Q_a = 5$

(c) $\mathcal{A}_{\mathrm{na}}^1$ and $Q_a = 40$
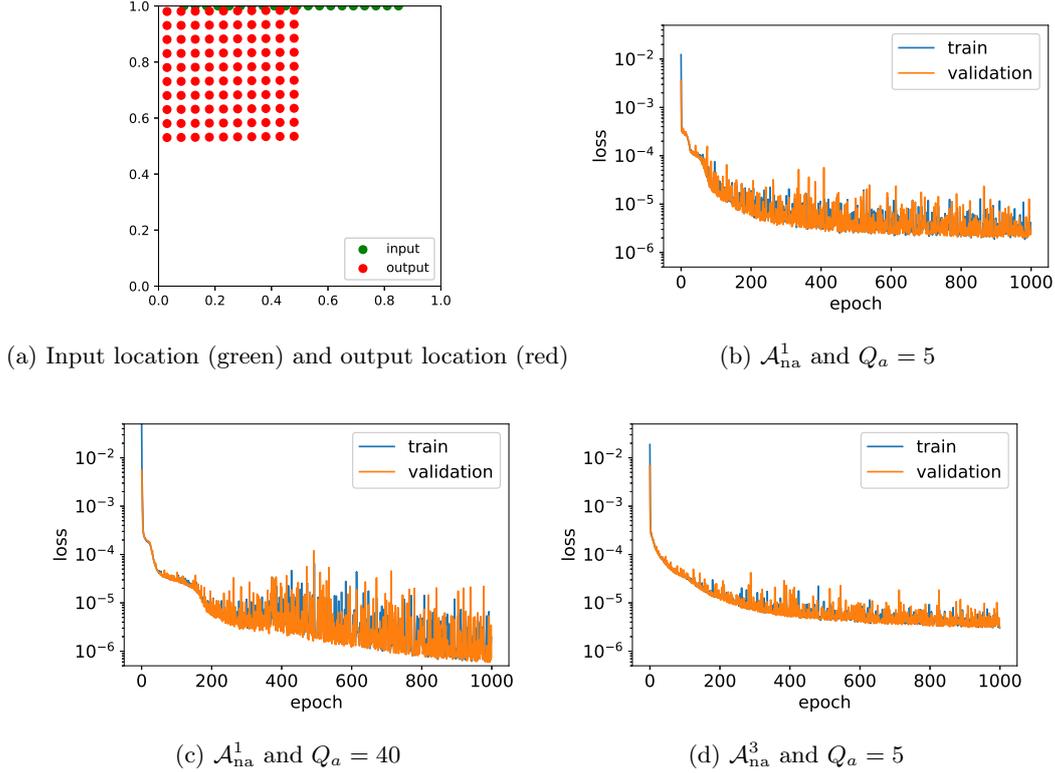
(d) $\mathcal{A}_{\mathrm{na}}^3$ and $Q_a = 5$

Figure 8: Location of input and output locations and decay of loss for the training and validation set different architectures and number of affine terms $Q_a$.

It is important to notice in the case of PDE-NNs the values $\{\theta_i^a(\boldsymbol{\mu})\}_{i=1}^{Q_a}$ are not trained to minimize the error between the original FE matrix $\mathbf{A}(\boldsymbol{\mu})$ and its affine approximation, but rather to minimize the error between the final output of the RB solver and the training output. We analyze the behavior of our network by varying the number of matrix affine components $Q_a$, which is set to 1, 2, 3, 4, 5, 10, 20, 40, and used to feed the RB solver. These numbers of affine components lead to an affine approximation of $\mathbf{A}(\boldsymbol{\mu})$ that by the standard computation of $\{\theta_i^a(\boldsymbol{\mu})\}_{i=1}^{Q_a}$ would lead to a relative error ranging from 1 to $10^{-5}$. We stress that for our PDE-DNNs we employ MDEIM for building an affine basis, but not for computing the corresponding coefficients $\{\theta_i^a(\boldsymbol{\mu})\}_{i=1}^{Q_a}$ of the affine approximation.

The training of the PDE-DNN is performed through minimization of the loss in Eq. (2) by Adam algorithm (with learning rate set to $10^{-3}$) for 1000 epochs; the parameters of the optimization method were set after empirical observations. For training the network, we employ $N_{\mathrm{s}} = 10000$ samples obtained by computing the FE full solutions of problem (16) corresponding to randomly sampled values of $\boldsymbol{\mu} \in \mathcal{D}$. We use this set also to construct a RB solver with $\varepsilon_{\mathrm{POD}} = 10^{-4}$. We considered this setting adequate after the analysis carried out in the test case in Section 5.1; 80% of the samples are used as training set, whereas the remainder serves as validation set. Examples of loss decay for the training validation sets corresponding to $\mathcal{A}_{\mathrm{na}}^1$ and $\mathcal{A}_{\mathrm{na}}^3$ architectures and $Q_a = 5$ and $Q_a = 40$ affine components are reported in Figure 8. As expected, the lower number of trainable parameters in the case of $\mathcal{A}_{\mathrm{na}}^3$ results in a significantly smoother decay of the loss function.

The trained network is then evaluated on a test set with 1000 samples corresponding to values of the parameters randomly sampled in $\mathcal{D}$ and different from the ones used during the training phase, and we evaluate the error achieved on the value of the solution at the output location by the three proposed architecture as function of the number of affine components predicted by the MLP. We also compare this result by the one computed by a standard RB solver, where the number of affine components is $Q_a$ and, given the physical parameter $\boldsymbol{\mu}$, the coefficients $\{\theta_q^a(\boldsymbol{\mu})\}_{q=1}^{Q_a}$ are computed online by solving an interpolation problem in MDEIM. As a matter of fact, the RB method requires

| $Q_a$ | 1 | 2 | 3 | 4 | 5 | 10 | 20 | 40 |
|---|---|---|---|---|---|---|---|---|
| $\mathcal{A}_{\mathrm{na}}^1$ | 2.19e-02 | 3.59e-02 | 2.55e-02 | 7.31e-03 | 5.15e-03 | 6.16e-03 | 5.58e-03 | 2.59e-03 |
| $\mathcal{A}_{\mathrm{na}}^2$ | 2.22e-02 | 8.58e-03 | 1.18e-02 | 5.46e-03 | 4.87e-03 | 3.11e-03 | 4.55e-03 | 2.36e-03 |
| $\mathcal{A}_{\mathrm{na}}^3$ | 2.22e-02 | 2.97e-02 | 6.26e-02 | 7.79e-03 | 5.69e-03 | 4.42e-03 | 4.52e-03 | 3.45e-03 |
| RB | 6.84e-01 | 7.63e-01 | 3.43e+00 | 1.34e+00 | 6.52e-01 | 4.98e-02 | 7.27e-03 | 3.32e-04 |

Table 5: Errors on the output for the three architectures and the RB method over 1000 test samples randomly selected and different from the ones using in the training phase.



Figure 9: Accuracy obtained by the PDE-DNN with three different architectures of the MLP and the standalone RB method as function the number of matrix affine components $Q_a$.

a fine affine decomposition of $\mathbf{A}(\boldsymbol{\mu})$ in order to provide a satisfactory result. As observed in the literature, the accuracy of the RB solution strictly depends on the accuracy of the matrix affine approximation: the lower $Q_a$ is, the less accurate the RB model becomes. On the contrary, the accuracy of the PDE-DNNs is mildly dependent of the number of affine components $Q_a$: a minimal amount of $Q_a = 4$ affine basis is required, where a plateau is reached. By employing a larger $Q_a$ no gain in accuracy is observed. This different behavior is ascribed to the fact that the PDE-DNN does not necessarily provide the values of $\{\theta_q^a(\boldsymbol{\mu})\}_{q=1}^{Q_a}$ that would be computed by the interpolation problem solved by the MDEIM algorithm, and that would provide a poor approximation in the case of small $Q_a$. Instead, it computes the coefficients $\{\theta_q^a(\boldsymbol{\mu})\}_{q=1}^{Q_a}$ which yield an accurate RB approximation with the given affine basis. Eventually, this results in the minimization of the loss function. On average, the PDE-DNN approximation is more accurate in the 100% of cases when using $Q_a = 1, 2, 3, 4, 5, 10$ and in about the 90% of cases if $Q_a = 20$. The standalone RB method always results in a better approximation only if $Q_a = 40$. The detailed errors can be found in Table 5.

## 5.3 Steady Navier-Stokes equations

In this section we introduce the steady Navier-Stokes (NS) equations for a viscous Newtonian incompressible fluid, which model the blood dynamics in artery bifurcations. Given an open bounded domain $\Omega \subset \mathbb{R}^2$, shown in Figure (10a), such that $\partial\Omega = \Gamma_{\mathrm{in}} \cup \Gamma_{\mathrm{out}} \cup \Gamma_{\mathrm{w}}$ and $\mathring{\Gamma}_w \cap \mathring{\Gamma}_{in} = \mathring{\Gamma}_{out} \cap \mathring{\Gamma}_{in} = \mathring{\Gamma}_w \cap \mathring{\Gamma}_{out} = \emptyset$ the steady NS equations read as follows:

$$\begin{cases} \vec{u} \cdot \nabla \vec{u} - \nu \Delta \vec{u} + \nabla p = \vec{0} & \text{in } \Omega \\ \nabla \cdot \vec{u} = 0 & \text{in } \Omega \\ +\text{b.c.} & \text{on } \partial\Omega \end{cases} \tag{18}$$

Here $\vec{u} = \vec{u}(\boldsymbol{\mu})$ and $p = p(\boldsymbol{\mu})$ are the velocity and the pressure fields describing the dynamics of the fluid and $\nu$ denotes the kinematic viscosity, which in our formulation is taken as physical parameter.

We employ a nonhomogeneous Dirichlet inlet condition $u = u(\vec{x}) = 4U(0.4 - x)^2$, with $U = 10$, on $\Gamma_{\text{in}}$, homogeneous Dirichlet conditions on the wall $\Gamma_{\text{w}}$ and homogeneneous Neumann conditions on the outlet boundary $\Gamma_{\text{out}}$. We define the Reynolds number $Re = L\bar{U}/\nu$ as the non-dimensional ratio of convection to diffusion, where $L$ and $\bar{U}$ are the characteristic length of the domain and velocity of the flow, respectively; here we deal with laminar flows, featuring $Re \in [1, 10^3]$. We are interested in the velocity dynamics in the presence of stenosis, that is the partial occlusion of the artery before the bifurcation. In this work, this is modelled by summing an additional reaction term $c(\vec{x}; r)\vec{u}$ to the left hand side of the first equation in (18), where

$$c(\vec{x}; r) = 10^3 \mathbb{I}_{\Omega_s(r)} + 10^{-10} \mathbb{I}_{\Omega_c(r)}$$

where $\mathbb{I}_A$ is the characteristic function on a set $A \subset \mathbb{R}^2$ and the regions $\Omega_s(r)$ and $\Omega_c(r)$ are defined as

$$\Omega_s = \Omega_s(r) = \{\vec{x} \in \Omega : 10(\vec{x}_1 - 1.5).^2 + (\vec{x}_2 - 0.46)^2 < r^2\}$$
$$\Omega_c = \Omega_c(r) = \{\vec{x} \in \Omega \backslash \Omega_s(r) : 10(\vec{x}_1 - 1.5).^2 + (\vec{x}_2 - 0.46)^2 < r_{\max}^2\},$$

where $r_{\max} = 0.25$ is fixed. The function $c(\vec{x}; r)$ defines an ellipsis modelling the presence of an obstacle which represents the stenosis in the region $\Omega_s$. The volume of the stenosis increases by increasing the value of $r$, which is the second parameter for this test case. A non-zero value is assigned also in the region $\Omega_c$, in order to guarantee that the support of $c(\vec{x}; r)$ be not parameter dependent. This ensures that the sparsity pattern of the nonaffinely parametrized FE matrix obtained by discretizing the additional reaction term does not change by changing the value of the parameter $r$. Should this property not hold, we could not apply MDEIM to affinely approximate it. Eventually, we define the parameter vector as

$$\boldsymbol{\mu} = (\nu, r) \in \mathcal{D} = [0.01, 0.1] \times [0, 0.25].$$

From (18), we can derive the corresponding variational formulation and subsequently the FE discretization by employing continuous piecewise second order polynomials for the the velocity and continuous piecewise first order polynomials for the pressure. This choice, also known as Taylor-Hood FE basis functions, is motivated by the saddle-point nature of problem (18) and guarantees the well-posedness of the resulting FE nonlinear system. We refer to, e.g., [16, 44, 6] for further details on the NS equations and their FE discretization. Solving the FE formulation of the NS equations is equivalent to solving a nonlinear system as (4) of the following form, where the matrix at the left hand side has a saddle-point structure

$$\begin{bmatrix} \mathbf{D}(\boldsymbol{\mu}) + \mathbf{C}(\mathbf{u}; \boldsymbol{\mu}) + \mathbf{K}(\boldsymbol{\mu}) & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} \mathbf{f}^u(\boldsymbol{\mu}) \\ \mathbf{0} \end{bmatrix}. \tag{19}$$

Here $\mathbf{u} \in \mathbb{R}^{N_h^u}$ and $\mathbf{p} \in \mathbb{R}^{N_h^p}$ are the vector representations of the velocity and pressure FE solutions, $\mathbf{f}^u$ corresponds to the discretization of the nonhomogeneous Dirichlet condition, $\mathbf{D}(\boldsymbol{\mu})$, $\mathbf{C}(\mathbf{u}; \boldsymbol{\mu})$, $\mathbf{K}(\boldsymbol{\mu}) \in \mathbb{R}^{N_h^u \times N_h^u}$ arise from the discretization of the second order differential operator, the nonlinear term and the additional obstruction term in the first equation of (18), respectively, and $\mathbf{B} \in \mathbb{R}^{N_h^p \times N_h^u}$ from the discretization of the incompressibility constraint given by the second equation in (18). We stress that $\mathbf{B}$ is not parameter dependent, whereas the nonlinear matrix $\mathbf{C}(\mathbf{u}; \boldsymbol{\mu})$ depends on the parameter by means of the solution $\mathbf{u}$.

**Remark 5.1.** *Let us introduce, for $\vec{w}, \vec{u}, \vec{v}$ in a suitable functional space, the trilinear form*

$$c(\vec{w}, \vec{u}, \vec{v}) = \int_\Omega (\vec{w} \cdot \nabla)\vec{u} \cdot \vec{v} d\Omega. \tag{20}$$

*The nonlinear matrix $\mathbf{C}(\mathbf{u}; \boldsymbol{\mu})$ arises from the FE discretization of (20) where $\vec{w} = \vec{u}$.*
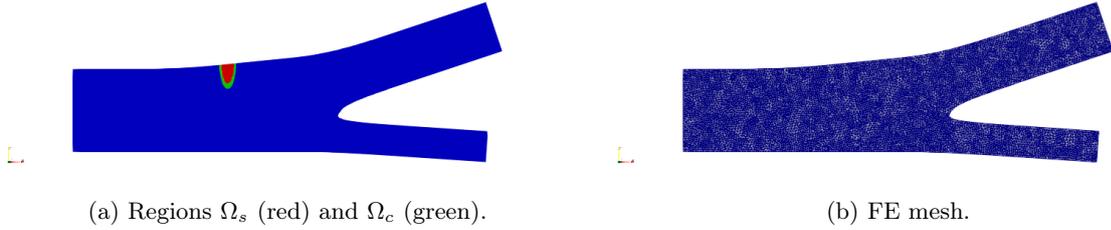
(a) Regions $\Omega_s$ (red) and $\Omega_c$ (green).



(b) FE mesh.

Figure 10: On the left, domain $\Omega$, boundaries, regions $\Omega_s(r)$, $\Omega_c(r)$ for $r = 0.2$. On the right, the FE mesh.



(a) $\boldsymbol{\mu} = (0.05, 0.15)$
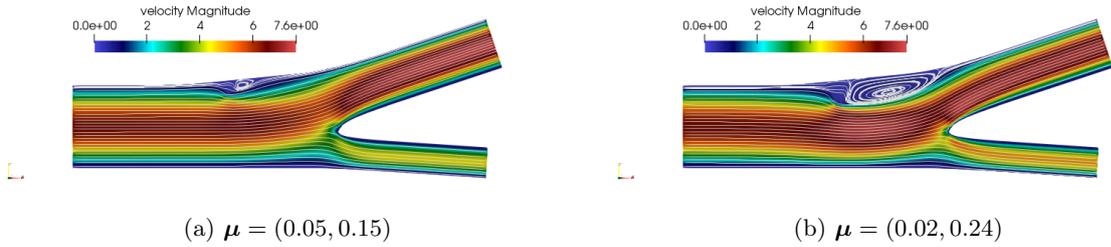


(b) $\boldsymbol{\mu} = (0.02, 0.24)$

Figure 11: Velocity solutions of the modified NS problem for $\boldsymbol{\mu} = (0.05, 0.15)$ (left) and $\boldsymbol{\mu} = (0.02, 0.24)$ (right), where a stenosis is simulated, yielding an obstruction and the subsequent creation of vortices.

In the numerical experiment we use a mesh with 15182 triangles, shown in Figure 10b, leading to $N_h^u = 59686$ dofs for the velocity and $N_h^p = 7592$ for the pressure. The FE system is then solved with the Newton method (with a tolerance of $10^{-8}$) in about 165.2 seconds, where at each iteration a linear system is solved with a direct method. An example of regions $\Omega_s$ and $\Omega_c$ is displayed in Figure 10a and the velocity obtained for two choices of the parameter, simulating the absence and the presence of the stenosis, are shown in Figures 11a and 11b, respectively. The larger the value of $r$, the smaller is the volume of central arterial lumen and the larger the corresponding velocity, leading also to the creation of vortices, as it is shown in Figure 11b.

The RB method (and in particular POD) has been largely employed to accelerate the solution of the parametrized steady NS equations, see e.g. [2]. To this aim, as explained in Section 3.2, we collect velocity solutions of (18) for a large enough amount of randomly chosen parameter in $\mathcal{D}$, we perform POD to build the RB projection matrix for velocity $\mathbf{V}_u \in \mathbb{R}^{N_h^u \times N_u}$, and approximate $\mathbf{u}$ in (19) with $\mathbf{V}_u \mathbf{u}_N$, with $\mathbf{u}_N = \mathbf{u}_N(\boldsymbol{\mu}) \in \mathbb{R}^{N_u}$. Being the columns of $\mathbf{V}_u$ obtained through POD, i.e. as linear combinations of solutions of (19) for some parameter instances, they are divergence free, that is $\mathbf{B}\mathbf{V}_u = \mathbf{0}$, meaning that the second equation in (19) is automatically satisfied when we substitute $\mathbf{u}$ with $\mathbf{V}_u \mathbf{u}_N$. Next, we left-multiply the remaining first equation by $\mathbf{V}_u^T$ to obtain the Galerkin-projected problem. Similarly, $\mathbf{V}_u^T \mathbf{B}^T = \mathbf{0}$, providing us with the following nonlinear RB system where the reduced velocity $\mathbf{u}_N$ is the only unknown

$$\left( \mathbf{D}_N(\boldsymbol{\mu}) + \mathbf{C}_N(\mathbf{V}_u \mathbf{u}_N; \boldsymbol{\mu}) + \mathbf{K}_N(\boldsymbol{\mu}) \right) \mathbf{u}_N = \mathbf{f}_N^u(\boldsymbol{\mu}). \tag{21}$$

The definition of the linear reduced operators in (21) take advantage of the (approximated) affine decomposition of the corresponding FE arrays, that is,

$$\mathbf{D}_N(\boldsymbol{\mu}) = \nu \mathbf{V}_u^T \mathbf{D} \mathbf{V}_u \qquad \mathbf{K}_N(\boldsymbol{\mu}) = \sum_{q=1}^{Q_k} \theta_q^k(\boldsymbol{\mu}) \mathbf{V}_u^T \mathbf{K}_q \mathbf{V}_u, \qquad \mathbf{f}^u(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f^u} \theta_q^{f^u}(\boldsymbol{\mu}) \mathbf{V}_u^T \mathbf{f}_q^u, \tag{22}$$

where $\mathbf{D} \in \mathbb{R}^{N_h \times N_h}$ is the FE matrix corresponding to the viscosity-independent Laplace operator, $\{\mathbf{K}_q\}_{q=1}^{Q_k} \subset \mathbb{R}^{N_h \times N_h}$ and $\{\mathbf{f}_q^u\}_{q=1}^{Q_f^u} \subset \mathbb{R}^{N_h}$ are the (M)DEIM basis for the affine approximation of

$\mathbf{K}(\boldsymbol{\mu})$ and $\mathbf{f}^u(\boldsymbol{\mu})$, respectively. For what concerns the nonlinear term $\mathbf{C}_N(\mathbf{V}_u\mathbf{u}_N;\boldsymbol{\mu})$, it admits the following (exact) affine decomposition

$$\mathbf{C}_N(\mathbf{V}_u\mathbf{u}_N;\boldsymbol{\mu}) = \sum_{j=1}^{N_u} u_N^j \mathbf{V}_u^T \mathbf{C}_j \mathbf{V}_u, \tag{23}$$

where $u_N^j$ is the $j$-th component of the RB solution vector, i.e. $u_N^j = [\mathbf{u}_N]_j$, and $\mathbf{C}_j \in \mathbb{R}^{N_h \times N_h}$, $j = 1, \ldots, N_u$, is the FE matrix corresponding to the nonlinear term (20) where $\vec{w} = \vec{\xi}_j$, with $\vec{\xi}_j$ the $j$-th RB function. The RB problem (21) for the steady NS equations is nonlinear and requires to employ a Newton solver and the affine decomposition of the Jacobian and the residual, which are inherited from the ones of the linear and nonlinear operators (22) and (23). The need to assemble and solve a nonlinear problem can clearly slow down the efficiency of the RB method.

There are several techniques to obtain a RB problem coupling both the velocity and pressure, which involve the use of the so-called supremizing functions, the interested reader can find more details in [2, 40, 32]. In this work are not interested to the value of pressure, however, in case of need, it can be reconstructed from the RB velocity through an additional pressure equation, as in [31].

### 5.3.1 PDE-DNNs for parametrized Navier-Stokes PDEs

We use a setting similar to the one employed in the previous test case and construct a PDE-DNN which has the measurements of the velocity as input and output, with an architecture as (2), where the RB solver interprets $\mathcal{O}_i^{\boldsymbol{\mu}}$ as the viscosity $\nu$, the coefficients $\{\theta_q^k(\boldsymbol{\mu})\}_{q=1}^{Q_k}$ and $\{\theta_q^{f^u}(\boldsymbol{\mu})\}_{q=1}^{Q_f^u}$ of the affine approximations (22) and eventually the coefficients $u_N^j$ of the affine expansion (23). Since the coefficients $u_N^j$ have an exponential decay, we decide to truncate them to the first $Q_n$ coefficients, leading to the approximation

$$\mathbf{C}_N(\mathbf{V}_u\mathbf{u}_N;\boldsymbol{\mu}) \approx \sum_{j=1}^{Q_n} u_N^j \mathbf{V}_u^T \mathbf{C}_j \mathbf{V}_u. \tag{24}$$

To summarize, the last layer processes

$$\mathcal{O}_i^{\boldsymbol{\mu}} = \begin{cases} \nu & \text{if } i = 1 \\ \theta_{i-1}^k(\boldsymbol{\mu}) & \text{if } i = 2, \ldots, Q_k + 1 \\ \theta_{i-Q_k-1}^f(\boldsymbol{\mu}) & \text{if } i = Q_k + 2, \ldots, Q_k + Q_f^u + 1 \\ u_N^j & \text{if } i = Q_k + Q_f^u + 2, \ldots, Q_k + Q_f + Q_n, \ j = i - Q_k - Q_f^u - 1. \end{cases}$$

The choice of predicting the coefficients $u_N^j$ has evident advantages, since the final layer of our PDE-DNN is called to solve a linear RB problem rather than a non-linear one. Indeed, we predict the coefficients $u_N^j$ in the affine decomposition of $\mathbf{C}_N(\mathbf{V}_u\mathbf{u}_N;\boldsymbol{\mu})$, instead of considering them as the unknown solution, which is a remarkable advantage when compared to the standard RB method for the NS equations. As a consequence, also the backpropagation (12) is easily computable.

In the numerical experiments, we have randomly selected a varying number of $N_{\text{in}}$ input locations and $N_{\text{out}}$ locations across the full domain $\Omega$.

We construct the RB problem (21) starting from $n_s = 19137$ snapshots for the velocity and with a tolerance $\varepsilon_{\text{POD}} = 10^{-4}$, leading to $N_u = 245$ reduced basis functions. Similarly, we employ (M)DEIM with 1000 vector or matrix snapshots to build the affine approximations of $\mathbf{f}_N^u(\boldsymbol{\mu})$ and $\mathbf{K}_N(\boldsymbol{\mu})$; in the experiments we will use $Q_k = Q_f^u = 2, 4, 8, 16$, which correspond to an affine approximation with a relative accuracy ranging from about 1 to $10^{-2}$. For the MLP component, we take as architecture the ones reported in Table 6 and we employ the same $n_s = 19137$ samples obtained by solving the FE problem corresponding to randomly selected instances of $\boldsymbol{\mu}$. We train our
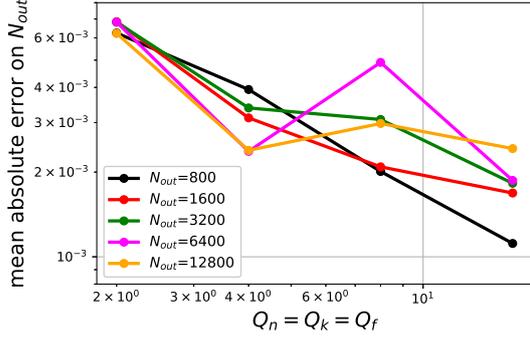
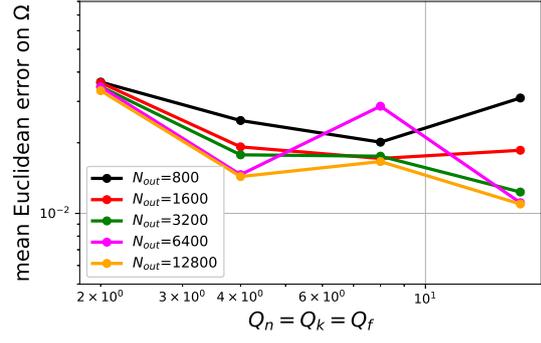| MLP architecture | Number of layers $L$ | Number of nodes per layer |
|:---:|:---:|:---:|
| $\mathcal{A}_{\text{ns}}^{1}$ | 4 | 64 |
| $\mathcal{A}_{\text{ns}}^{2}$ | 4 | 256 |
| $\mathcal{A}_{\text{ns}}^{3}$ | 4 | 1024 |
| $\mathcal{A}_{\text{ns}}^{4}$ | 8 | 1024 |

Table 6: Architectures employed for NS equations.

PDE-DNN with the Adam algorithm and 500 epochs and test against 1000 FE solutions, obtained for randomly selected instances of $\boldsymbol{\mu}$, different from the ones employed during the training. For simplicity, in the following we set the number of predicted affine components for the nonlinear term equal to the one $Q_k$ for the linear part, thus resulting in $Q_n = Q_k = Q_f^u = 2, 4, 8, 16$. The errors after training as function of the number of output locations $N_{\text{out}}$ are reported in Figure 12 for the case $N_{\text{in}} = 100$ for many architectures.

As a first consideration, we can see that the error has a decreasing trend with the number of affine components used for the nonaffine and nonlinear term, which is however interrupted by few bumps which are motivated by the training with the Adam algorithm with a mini-batch size equal to 128 and the random initialization of the trainable parameters. Then, we see that employing an architecture with a smaller amount of trainable parameters performs better when also a small number of output locations is employed; as an example we refer to Figure 12a, where a better result on the test set is achieved with $N_{\text{out}} = 800$. The behavior can be ascribed to the fact that a smaller number of trainable parameter is not able to capture the complex dynamics considered in this test case, not being able to create a map with too many output points. Yet, an overall better accuracy is achieved when more output points are employed, since in this case the error on a larger portion of the domain is penalized in the loss function used for the training (Figure 12b). As one could expect, the better results are provided when the architecture with a larger number of trainable parameter, $\mathcal{A}_{\text{ns}}^{4}$, is taken into account; this situation corresponds to Figure 12g, and achieves the best results in term of accuracy both on the $N_{\text{out}}$ test output locations and in the full domain. As expected, the larger the number of testing output locations $N_{\text{out}}$, the better the results on the full domain $\Omega$ (Figure 12h). Finally, we observe that by using a larger number of output locations, which consists in providing more information to the network during the training phase, yields a more significant use of a larger number of affine operators. Indeed, with $N_{\text{out}} = 800$, the accuracy does not significantly depend on the number of affine components; instead, the larger $N_{\text{out}}$, the better are the improvements given by a larger number of affine components. The architectures with an intermediate number of trainable parameters, present intermediate results in term of accuracy on both the test set and the full solution.
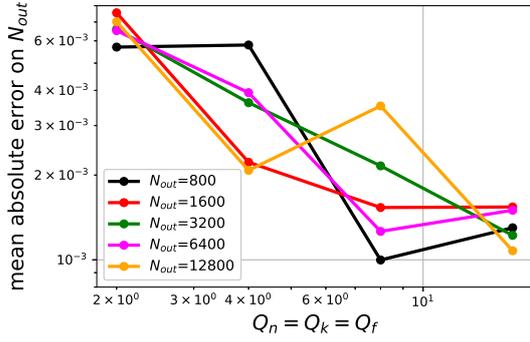
We consider in Figure 13 the cases $N_{\text{in}} = 200, 400$, both with the architecture $\mathcal{A}_{\text{ns}}^{4}$; we see that a larger number of input locations does not significantly affect the network in terms of accuracy and we observe the same trends with respect to the number of affine components and the output points as in the test with $N_{\text{in}} = 100$. In Figure 14 we report the error for the test parameter values, each point corresponds to a single instance of the parameter and it is colored according to the error committed when reconstructing the solution on the full domain (and compared with the true FE solution). The results consider a number of affine component $Q_n = Q_k = Q_f^u = 2, 16$ (respectively top and bottom line) and a number of output locations $N_{\text{out}} = 800, 12800$ (respectively left and right). When using a small number of affine components, i.e. $Q_n = Q_k = Q_f^u = 2$, the error is significantly larger than when using 16 affine components, since the complex dynamics is not captured by such a small number of affine components. This is confirmed for example when the fluid is less viscous (small $\nu$) and the occlusion is larger (large $r$), which yield the creation of larger vortices. Instead, the larger the number of affine components $Q_n = Q_k = Q_f^u = 16$, the smaller the errors and the more effective the use of a larger amount of output locations, up to a maximum error of about 1 % across the parameter space.
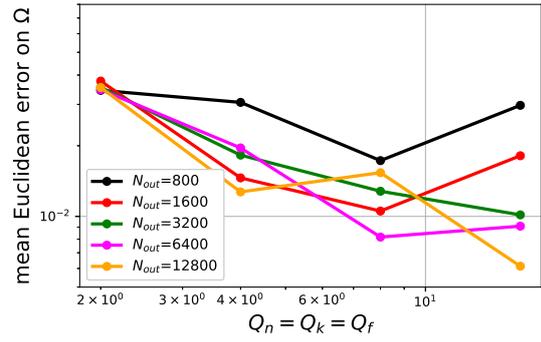
(a) NS errors at the test locations, $\mathcal{A}_{\mathrm{ns}}^1$ architecture.
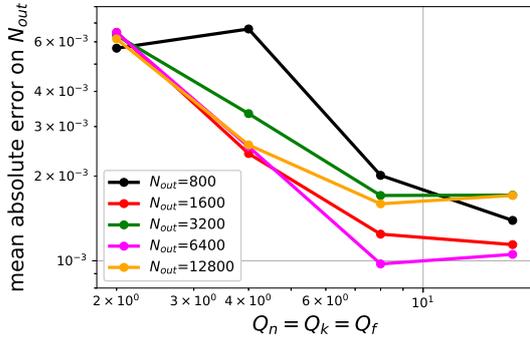
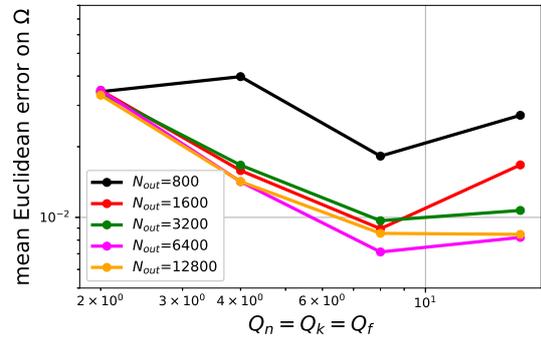(b) NS errors in $\Omega$, $\mathcal{A}_{\mathrm{ns}}^1$ architecture.

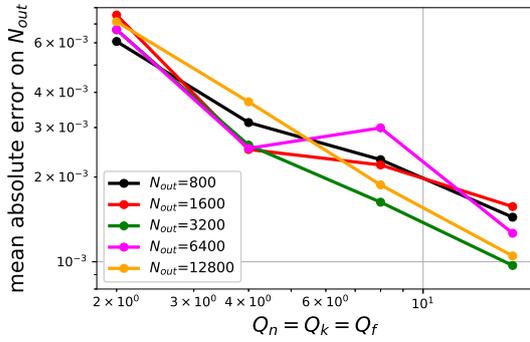(c) NS errors at the test locations, $\mathcal{A}_{\mathrm{ns}}^2$ architecture.

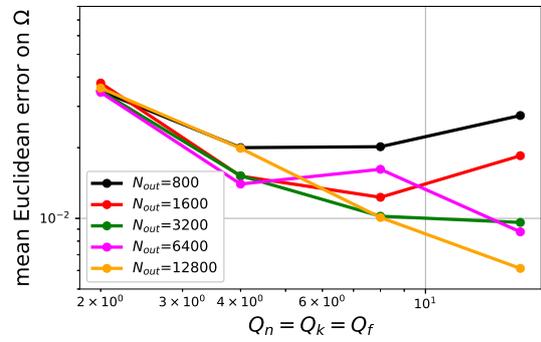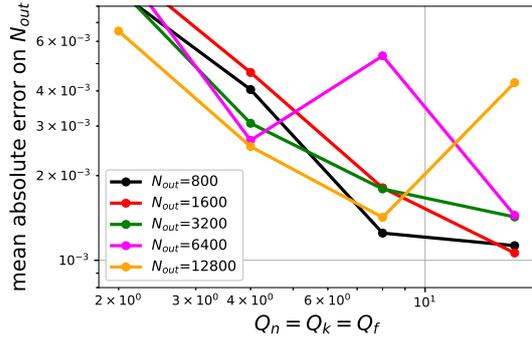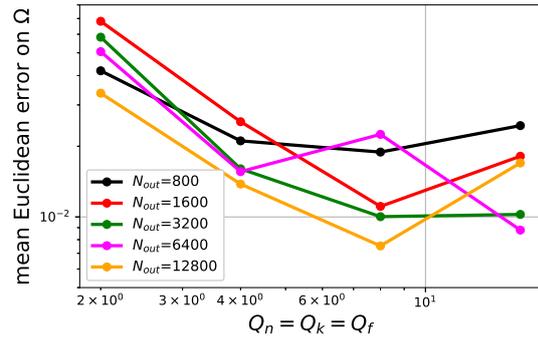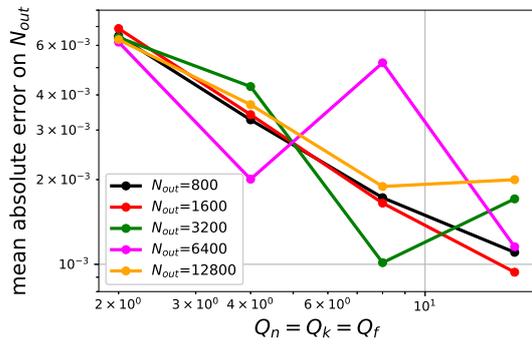(d) NS errors in $\Omega$, $\mathcal{A}_{\mathrm{ns}}^2$ architecture.

(e) NS errors at the test locations, $\mathcal{A}_{\mathrm{ns}}^3$ architecture.

(f) NS errors in $\Omega$, $\mathcal{A}_{\mathrm{ns}}^3$ architecture.

(g) NS errors at the test locations, $\mathcal{A}_{\mathrm{ns}}^4$ architecture.

(h) NS errors in $\Omega$, $\mathcal{A}_{\mathrm{ns}}^4$ architecture.

Figure 12: On the left, mean error on the test set FE functions at the output locations as function of the number of affine components for the linear (nonaffine) and nonlinear terms. On the right, mean error on the test set FE functions in the full $\Omega$. Different lines correspond to different architectures. All results are provided with $N_{\mathrm{in}} = 100$.
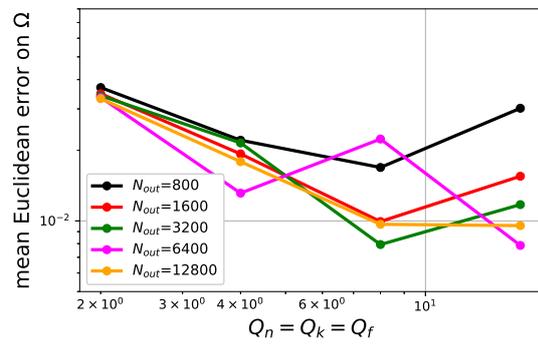
(a) NS errors at the test locations, $N_{\text{in}} = 200$.
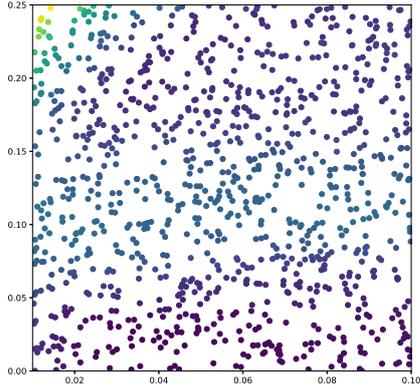
(b) NS errors in $\Omega$, $N_{\text{in}} = 200$.

(c) NS errors at the test locations, $N_{\text{in}} = 400$.
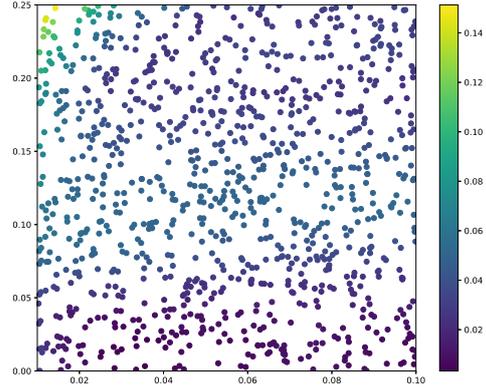
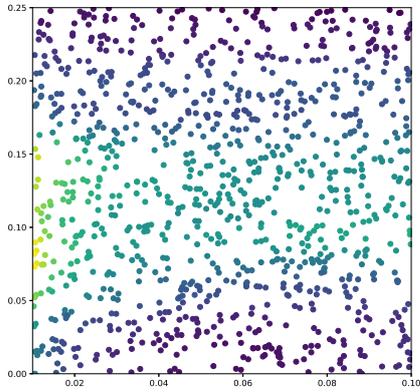(d) NS errors in $\Omega$, $N_{\text{in}} = 400$.

Figure 13: On the left, mean error on the test set FE functions at the output locations as function of the number of affine components for the linear (nonaffine) and nonlinear terms. On the right, mean error on the test set FE functions in the full $\Omega$. The top line refers to $N_{\text{in}} = 200$ and the bottom line refers to $N_{\text{in}} = 400$. All results are computed with $\mathcal{A}_{\text{ns}}^4$ architecture.
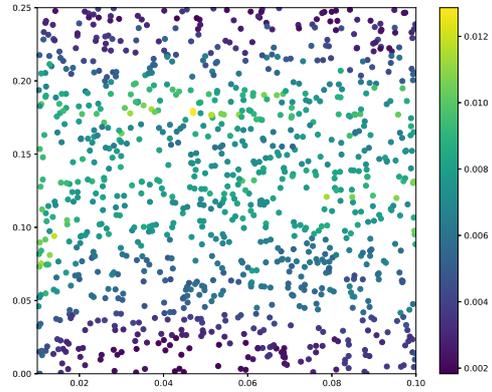
(a) $N_{\text{out}} = 800$, $Q_n = Q_k = Q_f^u = 2$.

(b) $N_{\text{out}} = 12800$, $Q_n = Q_k = Q_f^u = 2$.

(c) $N_{\text{out}} = 800$, $Q_n = Q_k = Q_f^u = 16$.

(d) $N_{\text{out}} = 12800$, $Q_n = Q_k = Q_f^u = 16$.

Figure 14: Scattered errors (on full $\Omega$) as function of the physical parameters computed with $\mathcal{A}_{\text{ns}}^4$ architecture, $N_{\text{in}} = 100$ and for different configurations of output locations (800 on the left, 12800 on the right) and number of affine components (2 on the top, 16 on the bottom).

# 6  Conclusions

In this paper we have proposed a novel way to integrate data and PDE numerical simulation by combining DNNs and RB solvers for the prediction of the solution of a parametrized PDE at some physical points. The proposed architecture is composed of a MLP followed by a RB solver, where the latter acts as a final nonlinear activation function interpreting the output of the MLP as prediction of parameter dependent quantities - physical parameters (affine case), theta functions of the approximated affine decomposition (nonaffine case) and approximated RB solutions (nonlinear case). Our architecture recovers an autoencoder structure when the input and output match, in this case the MLP acts as encoder and the RB solver acts as decoder. We have shown with a wide range of numerical experiments the features of the proposed methodology. Compared to standard DNN, we can obtain as byproduct the solution in the full physical space and, for affine dependencies, the value of the parameter. In the nonaffine and nonlinear case, we obtain accurate solutions by employing only a small amount of information compared to what needed by the standard RB method. Indeed, for the nonaffine case, only a small number of theta functions is needed by the affine approximation and, in the nonlinear case, we were able to turn the nonlinear problem to a linear one, yet recovering an accurate solution on the full computational domain.

Eventually, the advantages of the methods proposed in this paper can be leveraged for tackling even more complex problem classes. These include, e.g., multiphysics problems, where multiple PDE solvers can be weakly coupled by means of the quantities predicted by the MLP, FE error correctors, to improve the solution given by the considered PDE solver, and, finally, time dependent PDEs, where the (possibly expensive) time marching scheme can take advantage of an extrapolation given by the MLP.

# Acknowledgements

# References

[1] A. Abdulle and O. Budáč. A Petrov–Galerkin reduced basis approximation of the Stokes equation in parameterized geometries. *C. R. Math. Acad. Sci. Paris*, 353(7):641–645, 2015.

[2] F. Ballarin, A. Manzoni, A. Quarteroni, and G. Rozza. Supremizer stabilization of POD–Galerkin approximation of parametrized steady incompressible Navier–Stokes equations. *International Journal for Numerical Methods in Engineering*, 102(5):1136–1161, 2015.

[3] M. Barrault, Y. Maday, N. C. Nguyen, and A. T. Patera. An 'empirical interpolation' method: application to efficient reduced-basis discretization of partial differential equations. *Comptes Rendus Mathematique Académie des Sciences Paris*, 339(9):667–672, 2004.

[4] Y. Bengio et al. Learning deep architectures for ai. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.

[5] H. Bölcskei, P. Grohs, G. Kutyniok, and P. Petersen. Optimal approximation with sparsely connected deep neural networks. *SIAM Journal on Mathematics of Data Science*, 1(1):8–45, 2019.

[6] F. Brezzi. On the existence, uniqueness and approximation of saddle-point problems arising from lagrangian multipliers. *ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique*, 8(R2):129–151, 1974.

[7] S. Chaturantabut and D. C. Sorensen. Nonlinear model reduction via discrete empirical interpolation. *SIAM Journal on Scientific Computing*, 32(5):2737–2764, 2010.

[8] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. *Journal of machine learning research*, 12(Aug):2493–2537, 2011.

[9] G. Cybenko. Continuous valued neural networks with two hidden layers are sufficient. Technical report, Department of Computer Science, Tufts University, 1988.

[10] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.

[11] N. Dal Santo, S. Deparis, A. Manzoni, and A. Quarteroni. An algebraic least squares reduced basis method for the solution of nonaffinely parametrized stokes equations. *Computer Methods in Applied Mechanics and Engineering*, 344:186 – 208, 2019.

[12] N. Dal Santo and Manzoni. Hyper-reduced order models for parametrized unsteady navier-stokes equations on domains with variable shape. 2018.

[13] X. Ding, Y. Zhang, T. Liu, and J. Duan. Deep learning for event-driven stock prediction. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.

[14] W. E, J. Han, and Q. Li. A mean-field optimal control formulation of deep learning. *arXiv preprint arXiv:1807.01083*, 2018.

[15] W. E and B. Yu. The deep ritz method: A deep learning-based numerical algorithm for solving variational problems. *Communications in Mathematics and Statistics*, 6(1):1–12, Mar 2018.

[16] H. C. Elman, D. J. Silvester, and A. J. Wathen. Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics, 2005.

[17] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.

[18] B. Hanin. Universal function approximation by deep neural nets with bounded width and relu activations. *arXiv preprint arXiv:1708.02691*, 2017.

[19] B. Hanin. Which neural net architectures give rise to exploding and vanishing gradients? In *Advances in Neural Information Processing Systems*, pages 580–589, 2018.

[20] B. Hanin and D. Rolnick. How to start training: The effect of initialization and architecture. In *Advances in Neural Information Processing Systems*, pages 569–579, 2018.

[21] S. Haykin. *Neural networks: a comprehensive foundation*. Prentice Hall, Upper Saddle River, NJ, 2004.

[22] J. S. Hesthaven, G. Rozza, and B. Stamm. Certified reduced basis methods for parametrized partial differential equations. *SpringerBriefs in Mathematics*, 2016.

[23] J. S. Hesthaven and S. Ubbiali. Non-intrusive reduced order modeling of nonlinear problems using neural networks. *Journal of Computational Physics*, 363:55–78, 2018.

[24] Z. Hu, J. Tang, Z. Wang, K. Zhang, L. Zhang, and Q. Sun. Deep learning for image-based cancer detection and diagnosis- a survey. *Pattern Recognition*, 83:134–149, 2018.

[25] Y. Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.

[26] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[27] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

[28] S. Lai, L. Xu, K. Liu, and J. Zhao. Recurrent convolutional neural networks for text classification. In *Twenty-ninth AAAI conference on artificial intelligence*, 2015.

[29] S. Lathuilière, P. Mesejo, X. Alameda-Pineda, and R. Horaud. A comprehensive analysis of deep regression. *arXiv preprint arXiv:1803.08450*, 2018.

[30] K. Lee and K. Carlberg. Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders. *arXiv preprint arXiv:1812.08373*, 2018.

[31] A. E. Løvgren, Y. Maday, and E. M. Rønquist. A reduced basis element method for the steady Stokes problem. *ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique*, 40(3):529–552, 2006.

[32] A. Manzoni. An efficient computational framework for reduced basis approximation and a posteriori error estimation of parametrized Navier-Stokes flows. *ESAIM: Mathematical Modelling and Numerical Analysis*, 48(4):1199–1226, 2014.

[33] H. N. Mhaskar. Approximation properties of a multilayered feedforward artificial neural network. *Advances in Computational Mathematics*, 1(1):61–80, 1993.

[34] A. Mordvintsev, C. Olah, and M. Tyka. Inceptionism: Going deeper into neural networks, 2015.

[35] F. Negri, A. Manzoni, and D. Amsallem. Efficient model reduction of parametrized systems by matrix discrete empirical interpolation. *Journal of Computational Physics*, 303:431–454, 2015.

[36] P. Petersen, M. Raslan, and F. Voigtlaender. Topological properties of the set of functions generated by neural networks of fixed size. *arXiv preprint arXiv:1806.08459*, 2018.

[37] A. Quarteroni, A. Manzoni, and F. Negri. *Reduced Basis Methods for Partial Differential Equations: An Introduction.* Springer, 2016.

[38] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*, 2017.

[39] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics informed deep learning (part ii): data-driven discovery of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10566*, 2017.

[40] G. Rozza and K. Veroy. On the stability of the reduced basis method for Stokes equations in parametrized domains. *Computer Methods in Applied Mechanics and Engineering*, 196(7):1244–1260, 2007.

[41] C. Schwab and J. Zech. Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in uq. *Analysis and Applications*, 17(1):19–55, 2019.

[42] D. F. Specht. A general regression neural network. *IEEE transactions on neural networks*, 2(6):568–576, 1991.

[43] C. Szegedy, A. Toshev, and D. Erhan. Deep neural networks for object detection. In *Advances in neural information processing systems*, pages 2553–2561, 2013.

[44] R. Temam. *Navier-Stokes Equations*, volume 2. North-Holland Amsterdam, 1984.

[45] S. Volkwein. Proper orthogonal decomposition: Theory and reduced-order modelling. *Lecture Notes, University of Konstanz*, 4(4), 2013.

[46] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang. Residual attention network for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3156–3164, 2017.

[47] Q. Wang, J. S. Hesthaven, and D. Ray. Non-intrusive reduced order modeling of unsteady flows using artificial neural networks with application to a combustion problem. *Journal of computational physics*, 2019.

[48] T. Young, D. Hazarika, S. Poria, and E. Cambria. Recent trends in deep learning based natural language processing. *ieee Computational intelligenCe magazine*, 13(3):55–75, 2018.