

# Multi-fidelity surrogate modeling using long short-term memory networks

Paolo Conti<sup>a,\*</sup>, Mengwu Guo<sup>b</sup>, Andrea Manzoni<sup>c</sup>, Jan S. Hesthaven<sup>d</sup>

<sup>a</sup>*Department of Civil Engineering, Politecnico di Milano*

<sup>b</sup>*Department of Applied Mathematics, University of Twente*

<sup>c</sup>*MOX – Department of Mathematics, Politecnico di Milano*

<sup>d</sup>*Institute of Mathematics, École Polytechnique Fédérale de Lausanne*

---

## Abstract

When evaluating quantities of interest that depend on the solutions to differential equations, we inevitably face the trade-off between accuracy and efficiency. Especially for parametrized, time-dependent problems in engineering computations, it is often the case that acceptable computational budgets limit the availability of high-fidelity, accurate simulation data. Multi-fidelity surrogate modeling has emerged as an effective strategy to overcome this difficulty. Its key idea is to leverage many low-fidelity simulation data, less accurate but much faster to compute, to improve the approximations with limited high-fidelity data. In this work, we introduce a novel data-driven framework of multi-fidelity surrogate modeling for parametrized, time-dependent problems using long short-term memory (LSTM) networks, to enhance output predictions both for unseen parameter values and forward in time simultaneously – a task known to be particularly challenging for data-driven models. We demonstrate the wide applicability of the proposed approaches in a variety of engineering problems with high- and low-fidelity data generated through fine versus coarse meshes, small versus large time steps, or finite element full order versus deep learning reduced-order models. Numerical results show that the proposed multi-fidelity LSTM networks not only improve single-fidelity regression significantly, but also outperform the multi-fidelity models based on feed-forward neural networks.

*Keywords:* machine learning, multi-fidelity regression, LSTM network, parametrized PDE, time-dependent problem

---

## 1. Introduction

The construction of surrogate models is of paramount importance for multi-query, real-time simulations governed by parameterized, time-dependent partial differential equations (PDEs), as the surrogate models provide an efficient approximation of both parametric variation and time evolution in output quantities of interest. On the other hand, we often encounter the situation where many sources of data – large in number, easily accessible or fast to compute, but not perfectly accurate – are available. These low-fidelity (LF) data are ideally well correlated to the high-fidelity (HF) quantities that we aim to evaluate, and thus are expected to provide useful information in surrogate modeling. However, the LF data cannot guarantee a satisfactory credibility, because they are often generated by less reliable computations, such as coarse numerical discretizations, simplified physical assumptions, data fitting, and reduced-order modeling. To overcome this technical hurdle, multi-fidelity (MF) surrogate modeling methods have been developed for an effective data fusion from multiple sources. In particular, such methods are designed to detect trends from the ample LF data and enhance the predictive accuracy where the HF data are scarce. By exploiting the possibility to fast compute large amounts of LF data over the time-parameter domain of interest, the MF methods approximate the correlation between fidelity levels and hence infer the corresponding HF output quantities with limited HF data. A successful MF strategy should not only provide high-quality model predictions, but also achieve an overall efficiency improvement that is guaranteed by a substantially reduced cost of HF evaluations. As recognized, MF surrogate modeling is especially useful for real-time assessment

---

\*Corresponding author.

*Email addresses:* [paolo.conti@polimi.it](mailto:paolo.conti@polimi.it) (Paolo Conti), [m.guo@utwente.nl](mailto:m.guo@utwente.nl) (Mengwu Guo), [andrea1.manzoni@polimi.it](mailto:andrea1.manzoni@polimi.it) (Andrea Manzoni), [Jan.Hesthaven@epfl.ch](mailto:Jan.Hesthaven@epfl.ch) (Jan S. Hesthaven)

and monitoring, in which the running time of multi-query simulations over time and parameters should be much shorter than the operation time on the real-world asset to allow decision support.

Several MF strategies have been proposed to model the correlation among the data from different fidelity levels and found many applications in various areas of scientific computing [39]. A widely used MF surrogate modeling technique is co-kriging [35, 1], in which vector-valued Gaussian processes are used for the regression with multiple data sources. Such a non-parametric Bayesian method has become a popular tool of MF data fusion because of its good flexibility and intrinsic uncertainty quantification. However, Gaussian process regression suffers from the curse of dimensionality and thus significantly impacts the generalization performance of co-kriging in high-dimensional problems. As well known, artificial neural networks (NNs) are capable of handling high-dimensionality [21, 3]. Combined with their remarkable flexibility and multi-purpose nature, all these features made the NN-based models overwhelmingly successful in computational science and engineering, especially in the newly emerging area of *scientific machine learning* [2]. For instance, deep NNs are used to solve forward and inverse problems governed by PDEs [41, 25]; convolutional NNs have promoted massive advancements in image recognition [46]; deep auto-encoders provide a successful strategy for manifold learning and model reduction [7, 28, 12, 14]; and transformers [47] and recurrent NNs, especially long short-term memory (LSTM) networks [22, 17], have shown their effectiveness in time series analysis, e.g., for speech recognition [18, 43]. Importantly, NNs appear to be a promising candidate for MF surrogate modeling, because their strong expressive power for nonlinearity should be able to detect and represent the correlation among MF data sets, even with high-dimensional inputs. In particular, [31] tested a deep NN model with weighted linear and nonlinear components in the bi-fidelity correlation; a multi-step NN model was proposed and incorporated into Monte Carlo sampling in [32]; deep NNs were embedded into the kernel functions of co-kriging in [40]; Gaussian process latent variables were structured into a multi-layer network for MF modeling in [8]; a multi-step Bayesian NN model was developed for MF, physics-informed deep learning [30]; and LF information, such as initial guess of PDE solutions, was equipped to improve the training of physics-informed NNs [9]. More recently, MF data fusion has been explored with transfer learning [42], deep operator networks [23, 29], convolutional auto-encoders [38], and reinforcement learning [26]. Several multi-layer, multi-fidelity NN models were introduced in [20], inspired by the modeling assumptions of co-kriging combined with the interlink between multi-layer NNs and Gaussian processes [27, 19]. In addition, multi-fidelity NN surrogate models have been implemented in several engineering problems, such as structural health monitoring [44, 45] and aerodynamics [49].

Considering the demonstrated power of LSTM networks in time-series analysis, it is a natural and promising choice to embed LSTM units into the MF surrogate models for time-parameter-dependent problems. Though there have been a few LSTM-based MF techniques proposed towards specific applications, such as autonomous driving [24] and turbulent flow simulation [16], a general MF methodology that non-intrusively fuses a hierarchy of temporal, parametric data with LSTM networks is still lacking in the literature. To fill this gap, three data-driven MF surrogate models are proposed in this work using LSTM networks, generically for the approximation of both time evolution and parameter dependency in output quantities of interest. The proposed multi-fidelity NN models are intended as an extension of those in [20] to time-dependent problems. The choice of employing non-intrusive NN models ensures a remarkable flexibility when estimating quantities of interest for which we have no other information than the input-output data pairs. In general, however, due to a lack of incorporated physical model information, the predictive accuracy of data-driven, non-intrusive techniques may deteriorate dramatically outside the training data coverage. Poor generalization may especially be expected when extrapolating over time. In spite of this, our proposed models mitigate these challenges by combining MF strategies with LSTM networks, and present a good robustness in generalization performance. On one hand, the LSTM layers in our proposed models not only fit the data, but also recognize the underlying pattern of temporal evolution, thus enabling an effective prediction forward in time. On the other hand, the MF techniques leverage the dense sampling of cheap LF data to capture the time-parameter dependency towards accurate HF estimations. In this sense, the integration of LF data can be regarded as a ‘regularization’ that incorporates additional system information to the HF approximation.

A major goal of this work is to highlight the importance of embedding LSTM networks into multi-fidelity schemes. Thus, we compare the proposed multi-fidelity LSTM models with both the single-fidelity LSTM networks that are trained with either LF or HF data, and the MF models based on feed-forward NNs (rather than LSTM networks). To verify the effectiveness, applicability, and generality of the proposed models, we exemplify them in a diverse collection of benchmark problems, exploiting LF and HF levels of different nature.

We consider (i) a nonlinear diffusion-reaction system, namely, a parameterized one-dimensional FitzHugh-Nagumo membrane model that describes a simplified problem of action-potential propagation in an excitable cell; (ii) the temporal evolution of drag and lift coefficients in a fluid flow past a cylinder, governed by the unsteady Navier-Stokes equations; and (iii) a nonlinear dynamical system of the Lotka-Volterra type that describes the prey-predator interaction among three species. We consider output quantities with increased difficulty to approximate their time-parameter dependency. Throughout these numerical examples, we show that the proposed MF models are capable of predicting at unseen parameter locations and forward in time simultaneously. We address the following two cases for temporal extrapolation: (a) the LF data cover the whole time interval while the HF data are only available in a shorter time window, thus requiring the inference of HF outputs from the LF data outside the HF data coverage (see Sect. 6), and (b) the LF and HF data cover the same limited portion of the time interval, and the output predictions have to be carried out for new parameter values AND future states, at which we have no information of any fidelity at all (see Sect. 5). To the best of our knowledge, the proposed models are the first techniques that exploit both the multi-fidelity NNs’ capability in parametric generalization and the LSTM networks’ power in temporal pattern recognition, all in a general data-driven framework. Moreover, we utilize Bayesian optimization for hyperparameter tuning [5], especially for the identification of optimal NN architectures.

This paper is structured as follows. In Sect. 2, we introduce the proposed multi-fidelity LSTM network models and their main features. In Sect. 3, we present three numerical examples, including their governing physical models, quantities of interest, and bi-fidelity data sources. The results of these numerical tests are presented and discussed in Sects. 4, 5, and 6. Extension to more than two fidelity levels is discussed in Sect. 7 and, finally, conclusions are drawn in Sect. 8. We also provide a hyperparameter summary in the Appendix.

## 2. Multi-fidelity LSTM surrogate models

The central task of this work is the surrogate modeling of certain time-parameter-dependent quantities of interest, i.e., the approximation of the map from the time and parameter inputs  $\mathbf{x} = (t, \boldsymbol{\mu}) \in [t_0, T] \times \mathcal{P} = \mathcal{D}$  to the out quantities  $\mathbf{f}$ . Here  $t_0$  and  $T$  are the start and end instances of the time interval of interest,  $\mathcal{P}$  is the parameter domain, and  $\mathcal{D}$  denotes the full input domain. Towards this end, we present several types of NN architectures for the construction of MF surrogate models. These models are built upon the data of different fidelity levels:

- **High-fidelity (HF) data:** These data represent the best achievable accuracy and are obtained from detailed numerical simulations, e.g., the finite element approximations with suitably refined spatio-temporal discretizations. For parametrized, nonlinear, time-dependent differential problems, collecting these data can be so demanding that it may clash with computational budget restrictions. In particular, the multi-query evaluations at a large number of parameter locations may result in extremely expensive, or even impracticable, cost.
- **Low-fidelity (LF) data:** These data are less accurate but more accessible, inexpensive or faster to compute. For instance, LF data can be obtained from coarser discretizations, less strict convergence criteria, reduced-order/surrogate models, and/or simplified modeling assumptions.

Multiple levels of fidelity can be identified by adjusting the trade-off between accuracy and computational cost, or by considering different sources of data. In this work, we focus on bi-fidelity problems, i.e., we aim to estimate time-parameter-dependent HF quantities of interest with very limited HF observations, while leveraging abundant LF observations. Regarding this modeling task, it is useful to introduce the following notation:

- The LF parameter set is given by  $\mathcal{P}_{\text{LF}} = \{\boldsymbol{\mu}^{(j)} : 1 \leq j \leq N_{\text{LF}}^{\mu}\}$ , where  $\boldsymbol{\mu}^{(j)} \in \mathbb{R}^{p_{\mu}}$ ,  $j = 1, \dots, N_{\text{LF}}^{\mu}$ , are the parameter instances of the LF observations, and  $p_{\mu}$  is the dimension of the parameter space.
- We consider  $N_{\text{LF}}^t$  equally spaced instants  $\{t_n\}_{n=1}^{N_{\text{LF}}^t}$  over the time interval  $[t_0, T_{\text{LF}}]$ , where  $t_0$  is the initial time and  $T_{\text{LF}}$  is the time instant associated with the last LF measurement.
- The NN approximation of LF functions is denoted by  $\mathbf{f}_{\text{LF}}$ .

- The LF training set is given by  $\mathcal{T}_{\text{LF}} = \{(\mathbf{x}_{\text{LF}}^{(i)}, \mathbf{y}_{\text{LF}}^{(i)}) : 1 \leq i \leq N_{\text{LF}}\}$ , where  $N_{\text{LF}} = N_{\text{LF}}^t N_{\text{LF}}^\mu$  is the total number of time-parameter combinations. Here  $\mathbf{x}_{\text{LF}}^{(i)} \in \mathbb{R}^{p_{\text{in}}}$  collects each time-parameter location of the LF data inputs,  $p_{\text{in}} = p_\mu + 1$ , while  $\mathbf{y}_{\text{LF}}^{(i)} \in \mathbb{R}^{p_{\text{out}}}$  consists of the corresponding LF observations of the  $p_{\text{out}}$  quantities of interest that we aim to estimate.
- The notation for the HF counterparts is given by replacing the subscript LF with HF.

Typically, one has  $N_{\text{HF}}^\mu \ll N_{\text{LF}}^\mu$  based on the assumption that the availability of HF data is much limited due to their significantly higher computational cost compared to the LF data.

When dealing with sequential data, it is typically a good practice to organize the data into subsequences. For this reason, we group training data in batch tensors of shapes  $n_{\text{batch}} \times K \times p_{\text{in}}$  and  $n_{\text{batch}} \times K \times p_{\text{out}}$  for input and output data, respectively, in which  $n_{\text{batch}}$  is the batch size and  $K$  is the length of batch subsequences. Thus, for each parameter instance  $\boldsymbol{\mu}$ , we have a time-parameter input subsequence  $\{\mathbf{x}_n\}_{n=1}^K = \{(t_n, \boldsymbol{\mu})\}_{n=1}^K$  and a corresponding output subsequence  $\{\mathbf{y}_n\}_{n=1}^K$  containing the quantities of interest. Among the commonly used techniques for sequential data processing, recurrent NNs – in particular LSTM networks – represent the state-of-the-art in a wide range of applications. In this work, we employ LSTM cells as the main block to construct multi-fidelity NN surrogate models. We briefly present the structure and main features of an LSTM unit in the following subsection.

### 2.1. LSTM unit

Commonly in all recurrent networks, there is a *recurrent state* in an LSTM cell (see Fig. 1), here denoted by  $\mathbf{h}$ . The main feature of an LSTM unit is the presence of a *cell state*  $\mathbf{c}$ , which accounts for long-term dependencies. At each time-step  $n$ , the quantities  $(\mathbf{c}_n, \mathbf{h}_n)$  are computed from the previous step  $(\mathbf{c}_{n-1}, \mathbf{h}_{n-1})$  together with the input values  $\mathbf{x}_n$ , namely the time instant  $t_n$  and parameter values  $\boldsymbol{\mu}$  in our case. These variables  $(\mathbf{x}_n, \mathbf{h}_{n-1}, \mathbf{c}_{n-1})$  are passed to a three-fold *gate mechanism* with a *forget gate*  $\mathbf{F}_n$ , an *update gate*  $\mathbf{U}_n$  and an *output gate*  $\mathbf{O}_n$ , to produce new values for the cell and recurrent states  $(\mathbf{c}_n, \mathbf{h}_n)$  and to compute the output quantity  $\hat{\mathbf{y}}_n$ . Finally,  $(\mathbf{c}_n, \mathbf{h}_n)$  are propagated forward through such a recurrent mechanism. The gate mechanism is the core of an LSTM unit, designed to allow a refined memory management and overcome major drawbacks of recurrent networks, such as exploding and vanishing gradients. The *gates* in such a mechanism regulate the contribution of current and previous time-step information to determine how the states  $\mathbf{c}$  and  $\mathbf{h}$  change over time. In particular, the cell and recurrent states are updated as follows:

$$\mathbf{c}_n = \mathbf{F}_n \circ \mathbf{c}_{n-1} + \mathbf{U}_n \circ \tilde{\mathbf{c}}_n, \quad \mathbf{h}_n = \mathbf{O}_n \circ \tanh \mathbf{c}_n, \quad (1)$$

in which the operator  $\circ$  denotes the element-wise product, and  $\tilde{\mathbf{c}}_n$  represents the new candidate cell state to replace  $\mathbf{c}_{n-1}$ , computed as

$$\tilde{\mathbf{c}}_n = \tanh(\mathbf{W}_c[\mathbf{h}_{n-1}, \mathbf{x}_n] + \mathbf{b}_c). \quad (2)$$

The gates  $\mathbf{F}_n, \mathbf{U}_n, \mathbf{O}_n$  are defined as

$$\mathbf{F}_n = \sigma(\mathbf{W}_f[\mathbf{h}_{n-1}, \mathbf{x}_n] + \mathbf{b}_f), \quad \mathbf{U}_n = \sigma(\mathbf{W}_u[\mathbf{h}_{n-1}, \mathbf{x}_n] + \mathbf{b}_u), \quad \mathbf{O}_n = \sigma(\mathbf{W}_o[\mathbf{h}_{n-1}, \mathbf{x}_n] + \mathbf{b}_o). \quad (3)$$

Here,  $\{\mathbf{W}_c, \mathbf{W}_f, \mathbf{W}_u, \mathbf{W}_o\}$  and  $\{\mathbf{b}_c, \mathbf{b}_f, \mathbf{b}_u, \mathbf{b}_o\}$  are the trainable parameters – weights and biases – of the LSTM unit, and  $\sigma$  denotes the sigmoid activation function. Each gate takes both the concatenation of the current input parameters  $\mathbf{x}_n$  and the previous recurrent state  $\mathbf{h}_{n-1}$  as input variables, and provides an output vector of values between 0 and 1 that represent how much information is preserved from the pre-gate variables. The extremes 1 and 0 represent a complete preservation of information and a total discard, respectively. We note from (1) that the current cell state  $\mathbf{c}_n$  is obtained by the sum of the previous cell state and the new candidate, weighted respectively by the forget and update gates, through which the contribution of past and present information is effectively managed. Once we have updated the cell state  $\mathbf{c}_n$ , the recurrent state  $\mathbf{h}_n$  is updated accordingly. Different from the cell state  $\mathbf{c}$ , which is an internal variable of the LSTM recursive mechanism, the recurrent state  $\mathbf{h}$  also serves as the output of an LSTM unit and is passed to the next unit as input. The updated  $\mathbf{h}$  of the output layer should represent the estimation of the quantities of our interest  $\hat{\mathbf{y}}$ . The discrepancy between the network output  $\hat{\mathbf{y}}$  and the available data  $\mathbf{y}$  is measured by a mean squared error (MSE) loss function, and it is minimized to determine all the trainable parameters of the network model. For further details, we refer to [36, 22, 17].

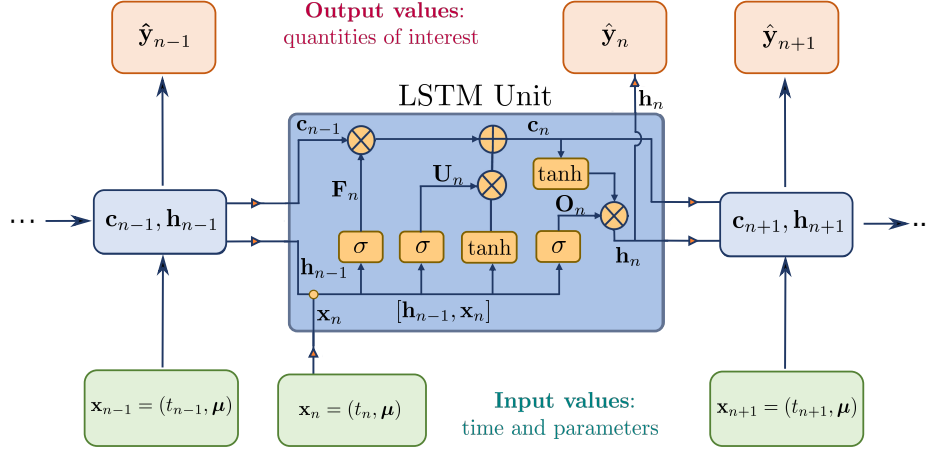


Figure 1: The visualization of an LSTM unit with its input-output setting [36, 48].

## 2.2. Multi-fidelity models

Now we introduce a variety of multi-fidelity NN architectures for the approximation of time-parameter-dependent output quantities  $\mathbf{f} = \mathbf{f}(\mathbf{x}) = \mathbf{f}(t, \boldsymbol{\mu})$ . Targeting an effective treatment of time-dependency, the models described below extend the strategies proposed in [20] through an incorporation of LSTM networks. The proposed MF architectures are shown in Fig. 2 and induce three different paradigms:

1. **“2-step” LSTM model.** The 2-step LSTM model is composed of two distinct NNs,  $NN_{\text{LF}}$  and  $NN_{\text{HF}}$ , each consisting of a sequence of LSTM layers followed by dense layers. The first network  $NN_{\text{LF}}$  (shown in blue in Fig.2(a)) is trained on the LF data  $\mathcal{T}_{\text{LF}}$  to learn the LF function  $\mathbf{f}_{\text{LF}}$ . Once  $NN_{\text{LF}}$  is trained, we compute the LF outputs  $\mathbf{f}_{\text{LF}}(\mathbf{x}_{\text{HF}})$  for each HF training input time sequence  $\mathbf{x}_{\text{HF}}$  through  $NN_{\text{LF}}$ . Next, we train the second network  $NN_{\text{HF}}$  (shown in red in 2(a)) to approximate  $\mathbf{f}_{\text{HF}}$ , using  $[\mathbf{x}_{\text{HF}}, \mathbf{f}_{\text{LF}}(\mathbf{x}_{\text{HF}})]^T$  as input data and the available HF evaluations  $\mathbf{y}_{\text{HF}}$  as output data.
2. **“3-step” LSTM model.** This architecture is an upgrade of the 2-step LSTM model by considering an extra concatenation level, generated by a third network  $NN_{\text{Lin}}$  (shown in green in Fig.2(b)).  $NN_{\text{Lin}}$  is trained with the same input-output data as  $NN_{\text{HF}}$  in the 2-step LSTM model. It does not use any nonlinear activation and thus reads the outputs as linear combinations of the inputs. Therefore,  $NN_{\text{Lin}}$  captures linear correlations between the HF and LF data sets. Then, the third and last network  $NN_{\text{HF}}$  exploits both the HF training set  $\mathcal{T}_{\text{HF}}$  and the outputs of the previous NNs to approximate the HF function  $\mathbf{f}_{\text{HF}}$ , i.e.,  $NN_{\text{HF}}$  is trained with inputs  $[\mathbf{x}_{\text{HF}}, \mathbf{f}_{\text{LF}}(\mathbf{x}_{\text{HF}}), \mathbf{f}_{\text{Lin}}(\mathbf{x}_{\text{HF}})]^T$  and output  $\mathbf{y}_{\text{HF}}$ .
3. **“Intermediate” LSTM model.** Different from the previous two models, the intermediate LSTM model only has a single network that simultaneously performs vector-valued learning of  $[\mathbf{f}_{\text{HF}}(\mathbf{x}), \mathbf{f}_{\text{LF}}(\mathbf{x})]^T$ . As shown in Fig. 2(c), a single input layer processes the time-parameter instances at both LF and HF levels, i.e.,  $\mathbf{x}_{\text{LF}}$  and  $\mathbf{x}_{\text{HF}}$ . Moreover, the corresponding outputs are placed at different locations: while  $\mathbf{y}_{\text{HF}}$  is the final output layer of the network,  $\mathbf{y}_{\text{LF}}$  is located at an intermediate LSTM layer, for which this model is named. To learn  $\mathbf{f}_{\text{LF}}$  and  $\mathbf{f}_{\text{HF}}$  simultaneously, the corresponding MSEs,  $\mathcal{L}_{\text{LF}}$  and  $\mathcal{L}_{\text{HF}}$ , respectively, are weighted by a hyperparameter  $\alpha \in [0, 1]$ . Hence one can perform a single optimization minimizing the following loss function:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{HF}} + (1 - \alpha) \mathcal{L}_{\text{LF}}. \quad (4)$$

in which the hyperparameter  $\alpha$  regulates the contributions from the two fidelity levels.

The mapping  $(t, \boldsymbol{\mu}) \mapsto \mathbf{f}(t, \boldsymbol{\mu})$  from time-parameter inputs to output quantities of interest can feature complex behaviors and thus be highly nonlinear. In this case, a large HF dataset is typically required to ensure predictive accuracy. By incorporating LF information into HF networks, the proposed MF models aim at simplifying the estimation of the HF input-output mapping so that it can be learnt from a small HF

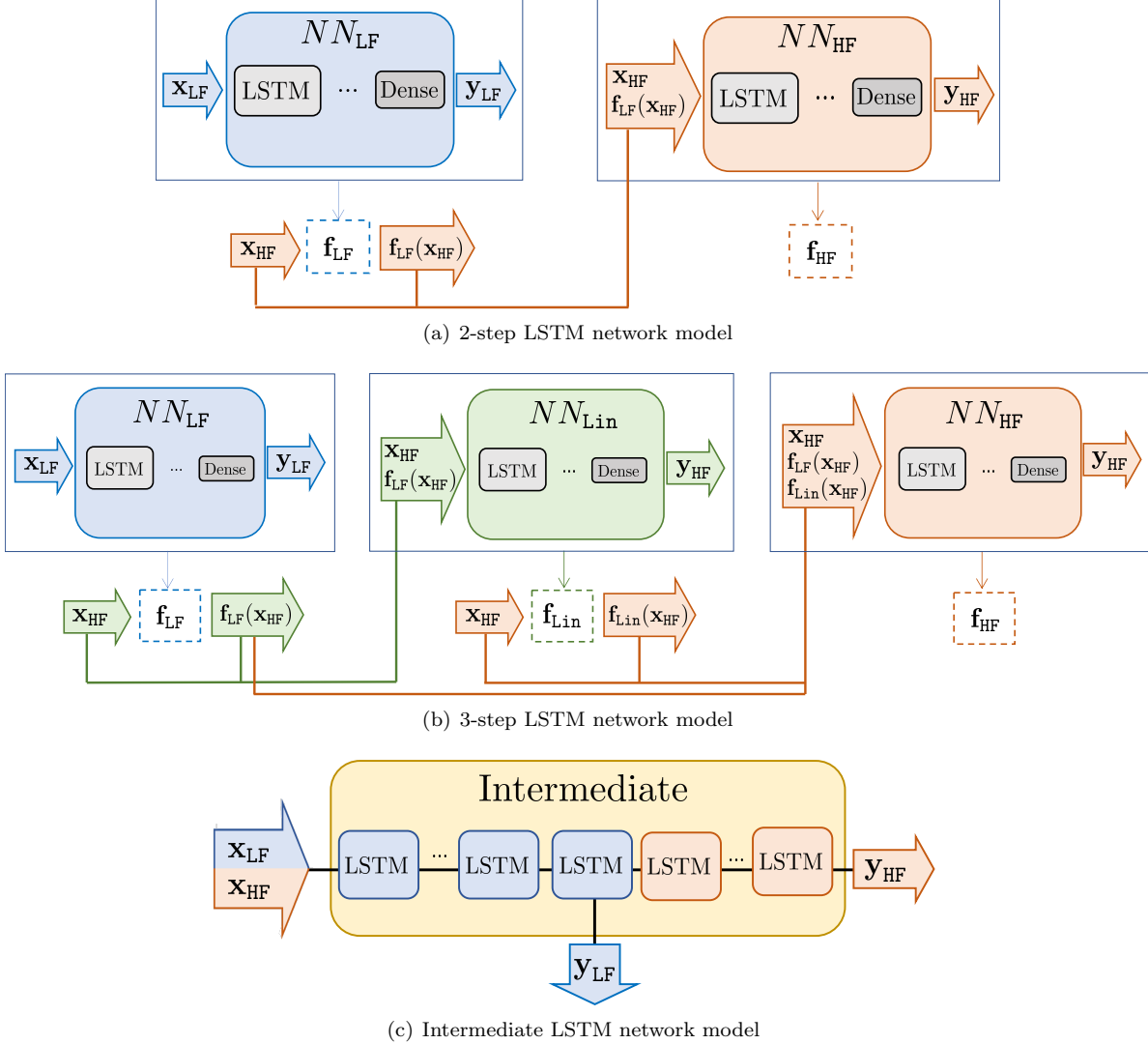


Figure 2: Multi-fidelity LSTM network models.

dataset. The HF part of the network models mainly learn the correlations between the two fidelity levels, which simplifies and accelerates the HF approximation, and prevents convergence to sub-optimal solutions. Instead of directly computing the expensive HF evaluations at new time-parameter instances, the proposed MF models seek to efficiently infer the HF outputs through the approximated fidelity correlations with cheaply obtained LF information. This not only avoids the high computational cost of HF evaluations, but also extends accurate HF predictions to the entire time-parameter domain of interest (LF data coverage), which can be considerably larger than the HF data coverage.

For each model, the numbers of LSTM and dense layers and the number of neurons per layer can be determined through hyperparameters optimization (HPO), which automatically identifies the optimal network setting along with other hyperparameters, such as the learning rate, optimization algorithm and batch size [6, 20].

We note that the intermediate LSTM model carries out simultaneous regressions on the LF and HF data, while the 2-step and 3-step LSTM models, referred to as *multi-level* models, approximate the LF and HF functions sequentially through a hierarchy of separately trained networks. This affects how the contribution of each fidelity level is weighted by the model setting, and how much trust is put in the LF data for the estimation of HF quantities. While the role of the LF data is adjusted explicitly by the hyperparameter  $\alpha$  in

the intermediate LSTM model, this is left implicit in the multi-level models and determined in the training process.

The *multi-level* models fuse multiple datasets by including additional networks. However, the computational burden of training the second- and/or third-step NNs is not heavy, because the correlation between the two fidelity levels should be uncomplicated enough to be learned with small amounts of HF data and rather simple network architectures. The 3-step model is expected to present advantages over the 2-step model in the case where considerable linear components exist in the correlation between HF and LF data. In addition to the  $NN_{LF}$  outputs, the  $NN_{Lin}$  outputs provide additional, important features for simplifying the HF approximation by  $NN_{HF}$ . Conversely, if the correlation between the two fidelity levels is predominantly nonlinear, the 2-step model is sufficient.

All the presented models can be extended to more than two fidelity levels, either by increasing the number of “steps” in the multi-level approaches or by locating all the fidelity levels except for the HF in different hidden layers of the intermediate LSTM model. More details are provided in Sect. 7. To assess the approximation capabilities of the proposed models, however, we limit ourselves to the bi-fidelity case and compare the proposed models with state-of-the-art regression techniques on a diversified collection of numerical benchmarks, which are introduced in the following section.

### 3. Numerical tests

In Sections 4, 5, and 6, we apply the proposed multi-fidelity LSTM models to the following numerical examples:

- (I) In Sect. 4 we analyze the propagation of an electrical signal in excitable cells described by a one-dimensional nonlinear PDE-ODE (coupled) system. The problem is parameterized by a coefficient that determines the amplitude and the steepness of the action potential front. The action potential at a given spatial location is then considered as the quantity of interest. HF data are generated from highly accurate finite element approximations, while LF data come from a reduced-order model based on deep learning.
- (II) The estimation of drag and lift coefficients for a fluid flow around a cylindrical obstacle as functions of the Reynolds number is considered in Sect. 5. Data are generated by numerical approximation of the unsteady Navier-Stokes equations, and the fidelity levels are defined by the quality of spatial and temporal discretizations.
- (III) The last example, as discussed in Sect. 6, considers a Lotka-Volterra system that represents a three-population prey-predator nonlinear interaction. The system is characterized by a parameter that regulates the amplitude and the frequency of the oscillatory pattern for each population size. Data are generated by numerical time integration, and the distinction between the fidelity levels is given by the size of time steps.

To assess the effectiveness and generality of the proposed models, the numerical examples are chosen to involve different types of governing equations and MF data generations.

In example (I), the quantity of interest – a point-wise action potential – presents a single wave front. In this case, we carry out MF regression with respect to both the input parameter and time using the proposed LSTM models. Their performance is compared with those of the MF approximation based only on feed-forward NNs and the LSTM regression trained solely with single-fidelity data.

Example (II) considers the drag and lift coefficients as output quantities. In the full developed state of the fluid flow, these quantities exhibit a periodic oscillatory pattern that varies with respect to the Reynolds number. From the regression point of view, this task is much more challenging than capturing the single peak of the solution in example (I). We repeat the comparison with other regression techniques and, in addition, test the robustness of the proposed LSTM models by varying the size of HF data set. Moreover, we evaluate the simultaneous parametric interpolation (over the parameter domain) and predictive extrapolation (beyond the training time interval), assessing the generalization accuracy while narrowing the training time window.

Finally, in example (III), we extend the proposed methods to the evaluation of vector-valued output quantities. We again test our models’ extrapolation performance over time, yet with a more challenging task

than example (II), as the outputs exhibit aperiodic oscillations. To track the aperiodic pattern, we let the LF solution evolve up to the final time of interest, while the HF data only cover a shorter time window.

Such diversity of the numerical examples aims to emphasize the wide applicability of the proposed MF strategies guaranteed by their non-intrusive nature, i.e., the surrogate models are directly learned from data without requiring access to the governing physical systems or numerical solvers. Instead of being sorted by increasing complexity of the governing physical systems, the numerical examples are listed in an order with increased difficulty of time-parameter-dependent surrogate modeling, e.g., from interpolation to extrapolation, from periodic to aperiodic behavior capture, and from scalar-valued to vector-valued function approximation.

For all the numerical examples, the comparison with the single-fidelity NNs trained on either LF or HF data set aims at highlighting the benefits of the MF modeling. Moreover, we compare the proposed LSTM architectures with the MF models solely based on feed-forward networks, such as those in [20], to stress the critical role played by LSTM layers for enabling good generalization properties with respect to both time and parameters. All the different network models are evaluated on a test set  $\mathcal{T}_{\text{test}} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N_{\text{test}}}$  that covers the whole time-parameter domain of interest. The goodness of fit is measured by the following MSE:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2, \quad (5)$$

where  $\{\hat{\mathbf{y}}_i\}_{i=1}^{N_{\text{test}}}$  are the predicted output values by the network models.

#### 4. Numerical example (I): Propagation of electrical signal

The first benchmark problem arises from computational biology and deals with the propagation of electrical potential in excitable cells, described by the following FitzHugh-Nagumo membrane model [11, 33]:

$$\begin{aligned} \mu \frac{\partial \nu}{\partial t} - \mu^2 \frac{\partial^2 \nu}{\partial x^2} + I_{\text{ion}}(\nu) + \omega &= 0, & x \in (0, L), t \in (0, T), \\ \frac{\partial \omega}{\partial t} + (\gamma \omega - b \nu) &= 0, & x \in (0, L), t \in (0, T), \\ \frac{\partial \nu}{\partial x}(0, t) = -i_0(t), \quad \frac{\partial \nu}{\partial x}(L, t) &= 0, & t \in (0, T), \\ \nu(x, 0) = 0, \quad \omega(x, 0) &= 0, & x \in (0, L). \end{aligned} \quad (6)$$

where  $\nu$  is the cardiac transmembrane electrical potential (excitation variable),  $\omega$  is the recovery variable which accounts for the refractoriness of heart cells, and  $t$  denotes a rescaled time. As given in [37], we choose  $I_{\text{ion}}(\nu) = \nu(\nu - 0.1)(\nu - 1)$ ,  $T = 2$ ,  $L = 1$ ,  $\gamma = 2$ , and  $b = 0.5$ . This problem is parametrized by  $\mu \in \mathcal{P} = [0.005, 0.05]$ , and its solution is characterized by the traveling wave fronts whose shapes are determined by the parameter  $\mu$ . The model can be extended to two- or three-dimensional spatial domains to represent the propagation of electrical stimuli on thin portions of tissue, or even on realistic geometries at the organ scale. For such a PDE-ODE coupled system – even more substantially, for two- or three-dimensional cases – full order solvers may be computationally demanding, as they usually involve very small time steps to properly detect signal propagation, as well as fine spatial meshes because of the steep fronts of the action potential. Therefore, we take advantage of MF strategies to achieve a reasonable approximation with improved efficiency and controlled accuracy. For the case at hand, the quantity of interest is the solution value of the electrical potential  $\nu$  at  $\bar{x} = 0.5$ , regarded as a function of the parameter  $\mu$  and time  $t$ .

##### 4.1. Multi-fidelity setting

The HF model is obtained by a linear finite element discretization over the spatial domain  $\Omega = (0, L)$  with the first-order semi-implicit scheme for time integration, and the total number of degrees of freedom is 1024. The LF model is constructed through a deep-learning-based reduced order modeling technique, POD-DL-ROM [14], which efficiently constructs non-intrusive reduced-order models for nonlinear parametrized time-dependent problems starting with a prior dimensionality reduction through the proper orthogonal



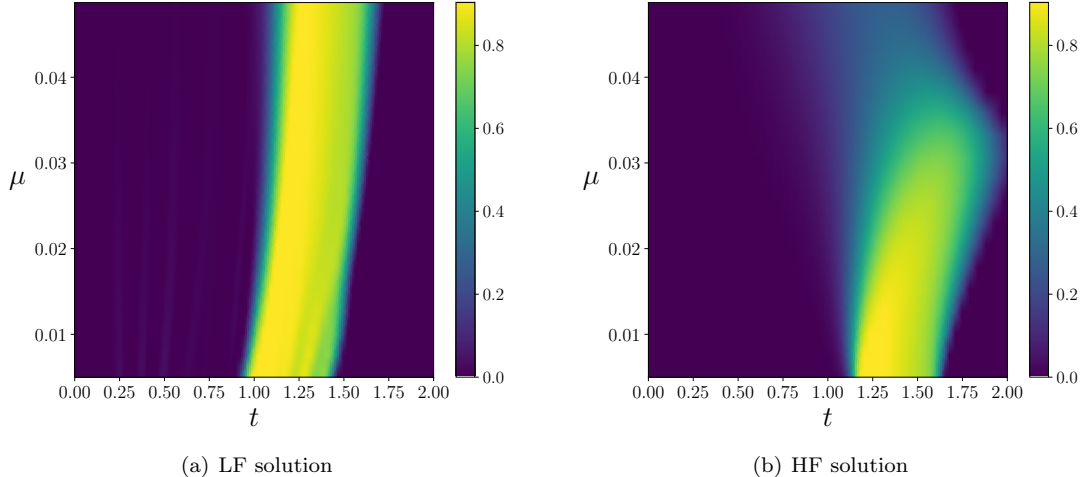


Figure 3: LF and HF solutions in example (I).

decomposition. 2 degrees of freedom (latent variables) are adopted in the LF reduced-order model throughout this example. We refer to [37] and [15, 14] for further details about the construction of the HF and the LF models, respectively. Solutions for the quantity of interest at the two fidelity levels are shown in Fig. 3.

Training data are obtained by solving either the LF or the HF model for parameter-time locations sampled over  $\mathcal{P} \times [0, T]$ . We consider  $N_{\text{HF}}^{\mu} = 4$  (resp.  $N_{\text{LF}}^{\mu} = 25$ ) parameter values uniformly spaced over  $\mathcal{P}$  for the sampling of HF (resp. LF) training data. For both fidelity levels, we choose  $\Delta t = 0.1$  as time step size, such that  $N_{\text{HF}}^t = N_{\text{LF}}^t = 20$ . We consider full-length time sequences by taking  $K = 20$ . In the testing stage, we employ a test set consisting of the HF evaluations on a uniform tensor grid over  $\mathcal{P} \times [0, T]$  with  $N_{\text{test}} = 18$  parameter locations and 20 time instants.

#### 4.2. Results and discussions

Our goal in this subsection is to demonstrate the advantages of the proposed MF LSTM surrogates over single-fidelity regressions. As shown in Table 1 and Fig. 4, we train eight different NN models and compare their predictions on the test set. LF and HF feed-forward networks are the single-fidelity models with only feed-forward dense layers, trained on  $\mathcal{T}_{\text{LF}}$  and  $\mathcal{T}_{\text{HF}}$ , respectively. LF LSTM and HF LSTM are their respective extensions incorporating LSTM layers. Intermediate, 2-step and 3-step LSTM networks are the proposed MF models with LSTM layers, while the 3-step feed-forward network is the counterpart of the latter without LSTM layers. The corresponding test MSE values are collected in Table 1, and the absolute discrepancy values between the model predictions and the HF test values are illustrated in Fig. 4. This allows us to better understand how the generalization quality varies over the parameter-time domain  $\mathcal{P} \times [0, T]$ .

We first note that, in the single-fidelity modeling with LF data  $\mathcal{T}_{\text{LF}}$ , i.e., the LF feed-forward and the LF LSTM cases, prediction errors are very large. Although a quite large training set is available, predictions are poor due to the low accuracy of the LF data, hence the possible advantage of changing architecture is limited by the data quality. In fact, using LSTMs does not yield any remarkable improvement, and both the LF feed-forward and the LF LSTM models show comparable predictive capability. In the single-fidelity modeling with HF data  $\mathcal{T}_{\text{HF}}$ , the HF LSTM model results in a clear improvement compared to the HF feed-forward model. As seen in Fig 4(d), however, there are regions where the discrepancy between HF and predicted solutions is large. Roughly equispaced over the parameter range, these regions correspond to the parameter configurations for which we have no HF information. This implies that, although incorporating LSTM layers has resulted in an improvement when approximating the time dependency, the single-fidelity HF LSTM model has poor generalization with respect to  $\mu$  due to the limited availability of training data.

On the other hand, we note that the MF LSTM models (intermediate, 2-step, and 3-step) achieve much better accuracy and allow to successfully generalize with respect to both parameter and time. These models, especially the 3-step LSTM (Fig. 4(h)), show uniformly low prediction errors over the time-parameter domain. We stress that these advantages result from the coupling of MF strategies with LSTM-based

Table 1: Mean squared errors (MSE) calculated on the test set in example (I).

Model	LF feed-forward	HF feed-forward	LF LSTM	HF LSTM
<b>Test MSE</b>	$1.06 \times 10^{-1}$	$2.99 \times 10^{-3}$	$1.11 \times 10^{-1}$	$7.87 \times 10^{-4}$
Model	MF 3-step feed-forward	MF Intermediate	MF 2-step LSTM	MF 3-step LSTM
<b>Test MSE</b>	$1.08 \times 10^{-2}$	$5.87 \times 10^{-4}$	$3.61 \times 10^{-4}$	$6.73 \times 10^{-5}$

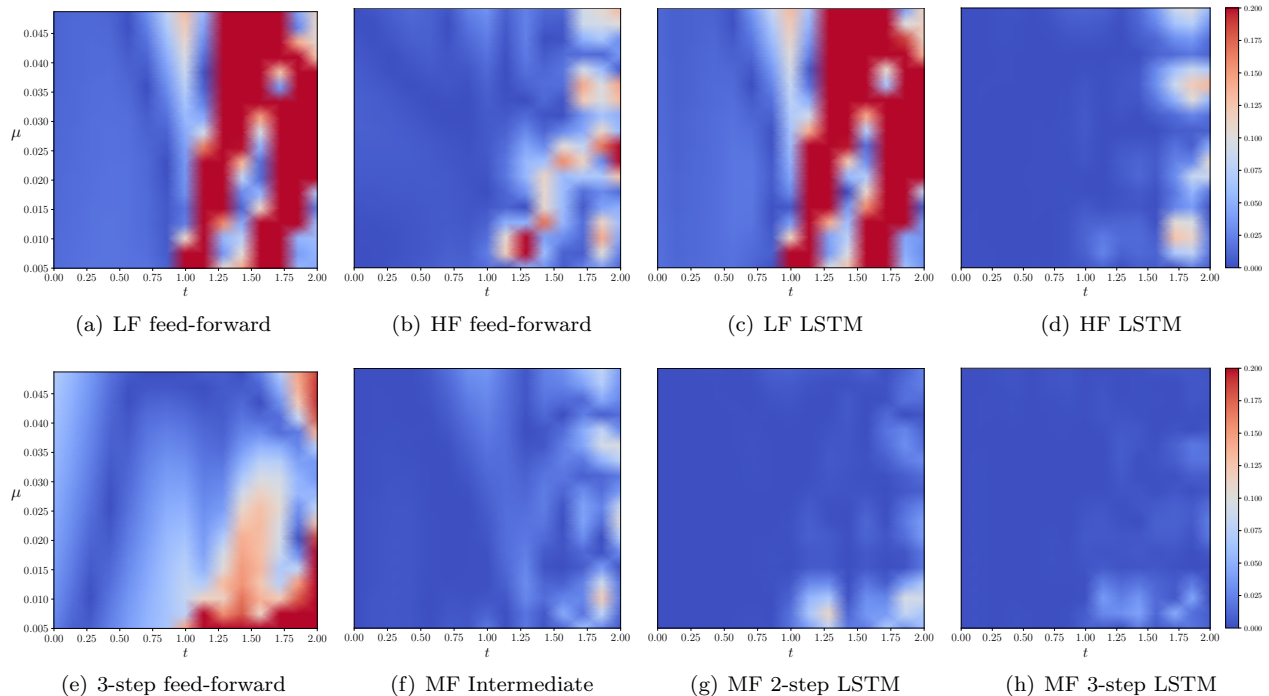


Figure 4: Absolute values of the errors in different NN model predictions of the electrical potential  $\nu$  at  $\bar{x} = 0.5$ .

architectures. In fact, when replacing LSTM layers with dense ones in the best performing MF LSTM model – the 3-step LSTM, we notice a dramatic deterioration of the predictive capability, as can be clearly seen in both Table 1 and Fig. 4(e). In conclusion, by leveraging the high quality of limited HF data and the intensive domain exploration of the LF data, the MF LSTM models are shown to be efficient tools for accurately predicting time-parameter-dependent quantities of interest that exhibit complex behaviors, e.g., the steep fronts of action-potential in this example.

## 5. Numerical example (II): Fluid flow around a cylinder

We now consider the estimation of both drag and lift coefficients associated with a viscous, incompressible fluid flow around a cylinder. The problem is described by the following unsteady Navier-Stokes equations:

$$\begin{aligned} \rho \frac{\partial \mathbf{v}}{\partial t} - \rho \mathbf{v} \cdot \nabla \mathbf{v} - \nabla \cdot \boldsymbol{\sigma}(\mathbf{v}, p) &= \mathbf{0} & (\mathbf{x}, t) \in \Omega \times (0, T), \\ \nabla \cdot \mathbf{v} &= 0 & (\mathbf{x}, t) \in \Omega \times (0, T), \end{aligned} \quad (7)$$

in which  $\mathbf{v}(\mathbf{x}, t)$  and  $p(\mathbf{x}, t)$  are respectively the velocity and pressure fields,  $\rho = 1.0 \text{ kg/m}^3$  is the fluid density,  $\boldsymbol{\sigma}(\mathbf{v}, p) = -p\mathbf{I} + 2\nu\boldsymbol{\epsilon}(\mathbf{v})$  is the stress tensor,  $\boldsymbol{\epsilon}(\mathbf{v}) = \frac{1}{2}(\nabla \mathbf{v} + \nabla^T \mathbf{v})$  is the strain tensor, and  $\nu$  is the kinematic viscosity. The domain,  $\Omega = (0, 2.2) \times (0, 0.41) \setminus B_r(0.2, 0.2)$  with  $r = 0.05$ , represents a 2D channel, while the

omitted disc,  $B_r$ , is the obstacle. Moreover, we prescribe the following boundary and initial conditions:

$$\begin{aligned} \mathbf{v} &= \mathbf{0} & (\mathbf{x}, t) \in \Gamma_{D_1} \times (0, T), \\ \mathbf{v} &= \mathbf{h} & (\mathbf{x}, t) \in \Gamma_{D_2} \times (0, T), \\ \boldsymbol{\sigma}(\mathbf{v}, p)\mathbf{n} &= \mathbf{0} & (\mathbf{x}, t) \in \Gamma_N \times (0, T), \\ \mathbf{v}(\mathbf{x}, 0) &= \mathbf{0} & \mathbf{x} \in \Omega, \end{aligned} \quad (8)$$

i.e., a no-slip condition on  $\Gamma_{D_1}$ , a parabolic inflow

$$\mathbf{h}(\mathbf{x}, t) = \left( \frac{4U(t)x_2(0.41 - x_2)}{0.41^2}, 0 \right), \quad \text{with } U(t) = \begin{cases} 0.75(1 - \cos(\pi t)), & t < 1 \\ 1.5, & t \geq 1 \end{cases} \quad (9)$$

on the inlet  $\Gamma_{D_2}$ , an open boundary condition on the outlet  $\Gamma_N$ , and a homogeneous initial condition.

Varying the parameter  $\nu$  results in a changing Reynolds number. In this specific problem, the Reynolds number can be defined as  $Re = LU_{\text{mean}}/\nu$ , in which the average free stream velocity for the parabolic inflow is  $U_{\text{mean}} = 2 \cdot U_{\text{max}}/3 = 2 \cdot 1.5/3 = 1$ , and  $L = 2r = 2 \cdot 0.05 = 0.1$  represents the characteristic length, and thus  $Re = 0.1/\nu$ . See also [14, 13] for further details on this setting.

We are interested in estimating the drag and lift coefficients, denoted by  $C_D$  and  $C_L$ , respectively. They are two important dimensionless quantities defined as follows:

$$C_D = \frac{2|F_x|}{\rho U_{\text{mean}}^2}, \quad C_L = \frac{2|F_y|}{\rho U_{\text{mean}}^2}. \quad (10)$$

Here  $F_x$  and  $F_y$  denote the two components of the total force  $\mathbf{F}$  acting on the cylinder, written as

$$\mathbf{F} = \oint_{\partial B_r} [-p\mathbf{I} + \nu(\nabla\mathbf{v} + \nabla^T\mathbf{v})] \cdot \mathbf{n} \, ds, \quad (11)$$

in which  $\mathbf{n}$  denotes the outer normal vector along the boundary  $\partial B_r$ .

Our goal is to predict the evolution of drag and lift coefficients over time as the Reynolds number  $Re$  changes. We consider  $\mu = Re \in [70, 160]$ , a range in which the flow is unsteady, and we confine our surrogate modeling within the time interval  $t \in [14.5, 15.0]$ , during which the flow becomes fully developed and presents a periodic behavior.

### 5.1. Multi-fidelity setting

Both the HF and LF models are constructed through the numerical approximation of (7) using the MATLAB library redbKIT [34], which exploits finite elements and the backward differentiation formula for the spatial and temporal discretizations, respectively. The two fidelity levels are distinguished by the mesh and time step sizes. In the HF model we adopt a time step  $\Delta t_{\text{HF}} = 0.01$  and a fine computational mesh, while we use  $\Delta t_{\text{LF}} = 0.02$  and a coarse mesh in the LF model. The two meshes are displayed in Fig. 6.

As for the data sets, we first consider uniform grids of  $N_{\text{LF}}^\mu = 19$  and  $N_{\text{HF}}^\mu = 10$  Reynolds number values over the parameter interval  $\mathcal{P} = [70, 160]$ . For each parameter value, we include  $N_{\text{HF}}^\tau = N_{\text{LF}}^\tau = 26$  uniform steps over the time interval  $[t_0, T] = [14.5, 15.0]$ , and compute the corresponding values of drag and lift

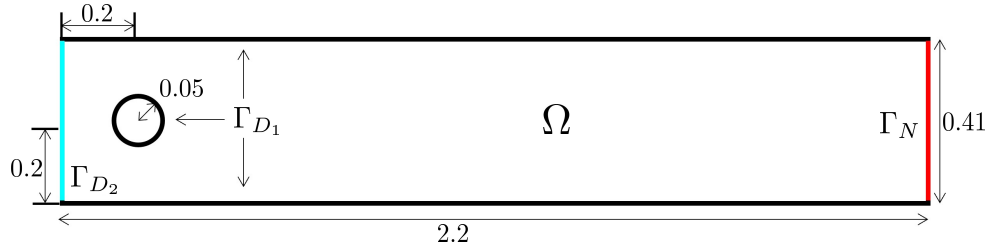


Figure 5: Geometry for the 2-D channel flow around a cylinder. All lengths are measured in meters.

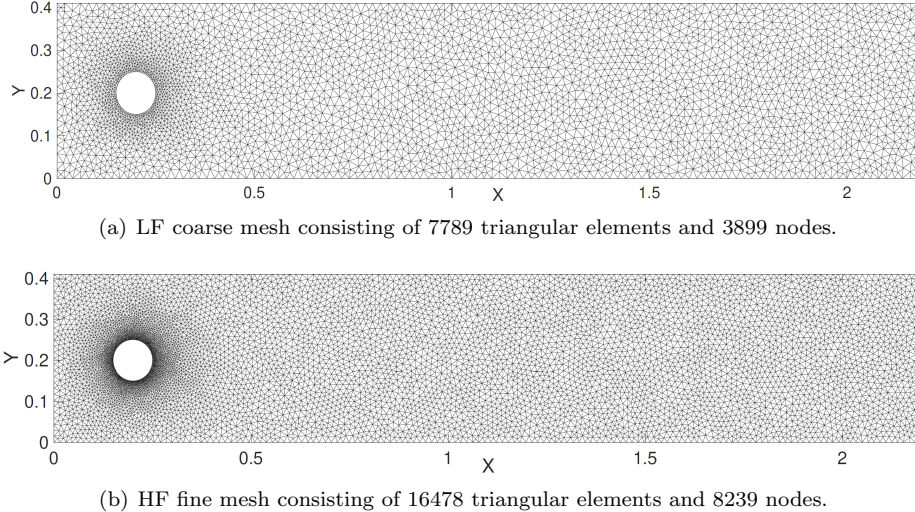


Figure 6: Computational meshes used for the finite element discretization of the Navier-Stokes equations (7).

coefficients. In the LSTM models, we consider full-length time sequence by setting the batch subsequence length  $K = 26$ . NNs are trained to approximate the drag and lift coefficients as functions of time  $t$  and Reynolds number  $Re$ . Moreover, the test set comprises the HF evaluations of drag and lift coefficients for 19 Reynolds number values and 26 time instants that are uniformly spaced in  $\mathcal{P}$  and  $[t_0, T]$ , respectively.

Next, we repeat the analysis with a smaller HF training set consisting of  $N_{\text{HF}}^{\mu} = 6$  equally spaced Reynolds number values, while keeping  $N_{\text{HF}}^T = 26$ , to assess the robustness of the NN models. In this work, we assume that LF data are cheap/easy to obtain in a sufficiently large amount, so that the mapping from the time-parameter inputs to LF outputs can be approximated accurately. A general guidance for choosing the number of LF data is to verify that the LF training and testing errors are controlled within the same order of magnitude, i.e., a good generalization on the LF level is achieved. Hence, what primarily limits the model performance is either the *quality* of LF approximation or the *quantity* of HF data. In this example, we analyze the latter case while fixing the LF model, because our focus is on exploring the models' prediction and extrapolation capabilities when the HF data are limited. However, we refer to [20] for a detailed discussion about the impact of varied LF model quality on the feed-forward versions of intermediate, 2-step and 3-step architectures in a parametric differential problem solved through a reduced basis method.

The LF and HF solutions are shown in Fig. 7. We note that the drag and lift coefficients exhibit oscillatory patterns over time, varying with respect to the Reynolds number. This may be problematic for the feed-forward NNs that treat time as a generic input entry.

### 5.2. Results and discussions 1: Interpolation over the time-parameter domain

Similar to example (I), we train eight different NN models on the aforementioned data sets, and compare the proposed MF LSTM models with both the single-fidelity regressions and the MF feed-forward networks. In Tables 2 and 3, we collect the prediction errors on the test set from all eight NN models. Our MF LSTM models achieve the best prediction accuracy. In particular, the 2-step and 3-step LSTM models outperform the others, as highlighted from the errors reported in the Table. The discrepancy between the HF solutions and network model predictions is displayed in Figs. 8 and 9, the former for the lift with  $N_{\text{HF}}^{\mu} = 10$  and the latter for the drag with  $N_{\text{HF}}^{\mu} = 6$ .

We notice that the feed-forward networks, which treat the time and parameter as equal input entries, fail to approximate the oscillatory temporal patterns in the output quantities (see (a), (b), and (e) of Figs. 8 and 9). On the other hand, the single-fidelity LSTM regressions (subfigures (c) and (d)) manage to capture the oscillations; however they are still unable to guarantee high-quality predictions either due to low data quality or limited data availability. The proposed MF LSTM models (subfigures (f), (g), and (h)), especially the 2-step and 3-step LSTM, achieve good predictive accuracy all over the time-parameter domain  $\mathcal{P} \times [t_0, T]$ .

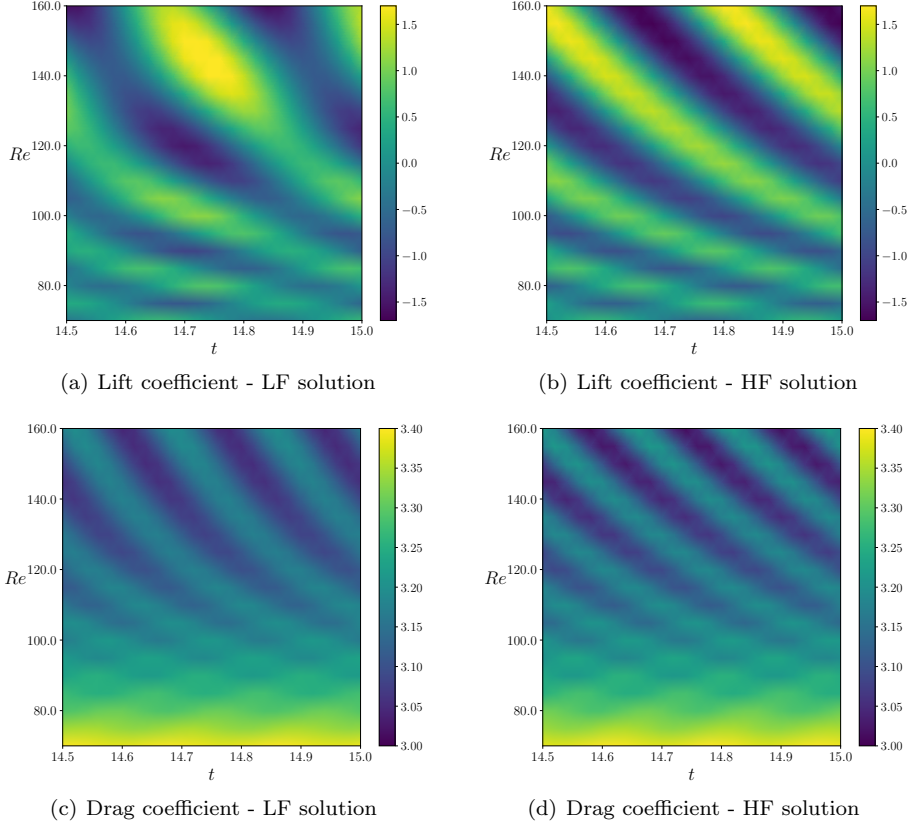


Figure 7: LF and HF solutions of the lift (above) and drag (below) coefficients.

Table 2: Test mean square errors (MSE) of the single-fidelity models in example (II).

Output quantity	#HF data	LF feed-forward	HF feed-forward	LF LSTM	HF LSTM
Lift coefficient	6	$8.64 \times 10^{-2}$	$1.12 \times 10^{-1}$	$2.45 \times 10^{-1}$	$1.12 \times 10^{-1}$
Lift coefficient	10	"	$8.63 \times 10^{-2}$	"	$1.92 \times 10^{-2}$
Drag coefficient	6	$9.45 \times 10^{-3}$	$9.43 \times 10^{-3}$	$8.77 \times 10^{-3}$	$1.33 \times 10^{-2}$
Drag coefficient	10	"	$9.41 \times 10^{-3}$	"	$1.63 \times 10^{-3}$

Table 3: Test mean square errors (MSE) of the multi-fidelity models in example (II).

Output quantity	#HF data	3-step feed-forward	Intermediate	2-step LSTM	3-step LSTM
Lift coefficient	6	$8.67 \times 10^{-2}$	$2.83 \times 10^{-2}$	<b><math>1.14 \times 10^{-2}</math></b>	$2.24 \times 10^{-2}$
Lift coefficient	10	$8.63 \times 10^{-2}$	$6.64 \times 10^{-3}$	$1.25 \times 10^{-3}$	<b><math>1.22 \times 10^{-3}</math></b>
Drag coefficient	6	$9.44 \times 10^{-3}$	$7.59 \times 10^{-3}$	<b><math>1.58 \times 10^{-3}</math></b>	$2.28 \times 10^{-3}$
Drag coefficient	10	$9.41 \times 10^{-3}$	$7.99 \times 10^{-3}$	$2.91 \times 10^{-4}$	<b><math>2.86 \times 10^{-4}</math></b>

In addition, we show in Fig. 10 the predictions from the single-fidelity (LF and HF) LSTM models and the best performing MF models, i.e., the 3-step (resp. 2-step) LSTM for the lift (resp. drag) coefficient with  $N_{\text{HF}}^{\mu} = 10$  (resp.  $N_{\text{HF}}^{\mu} = 6$ ). Comparing these figures to the LF and HF solutions that generate the data (see Fig. 7), it is reasonable that the LF LSTM reconstruction based on a large data set is in good agreement with the LF ground truth, yet still differs significantly from the HF solution because of the coarse spatio-temporal discretization, and that the HF LSTM reconstruction exhibits a better fit to the HF ground truth yet is still not ideal due to the limited data coverage. By taking advantage of both the large amount

of LF data and the limited HF measurements, the MF results, however, achieve a close-to-HF accuracy in learning the temporal patterns with the LSTM layers.

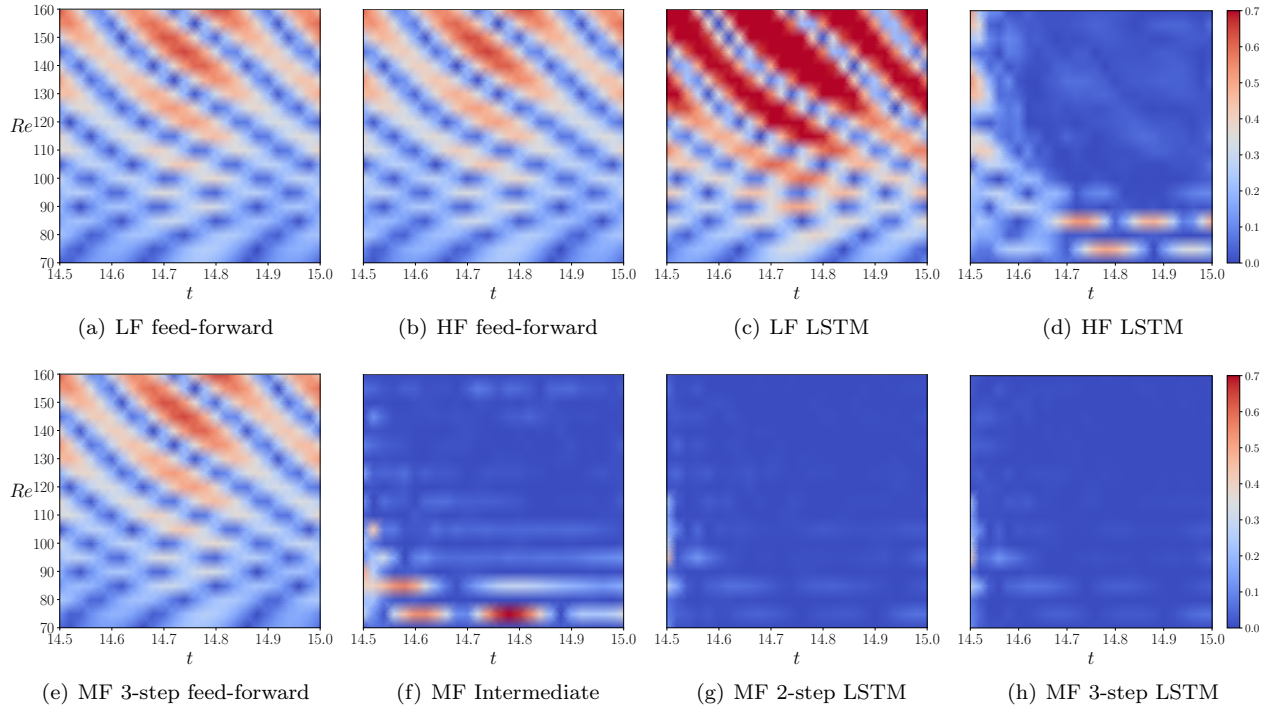


Figure 8: Absolute values of the errors in different NN model predictions of the lift coefficient with  $N_{\text{HF}}^{\mu} = 10$  HF training time series.

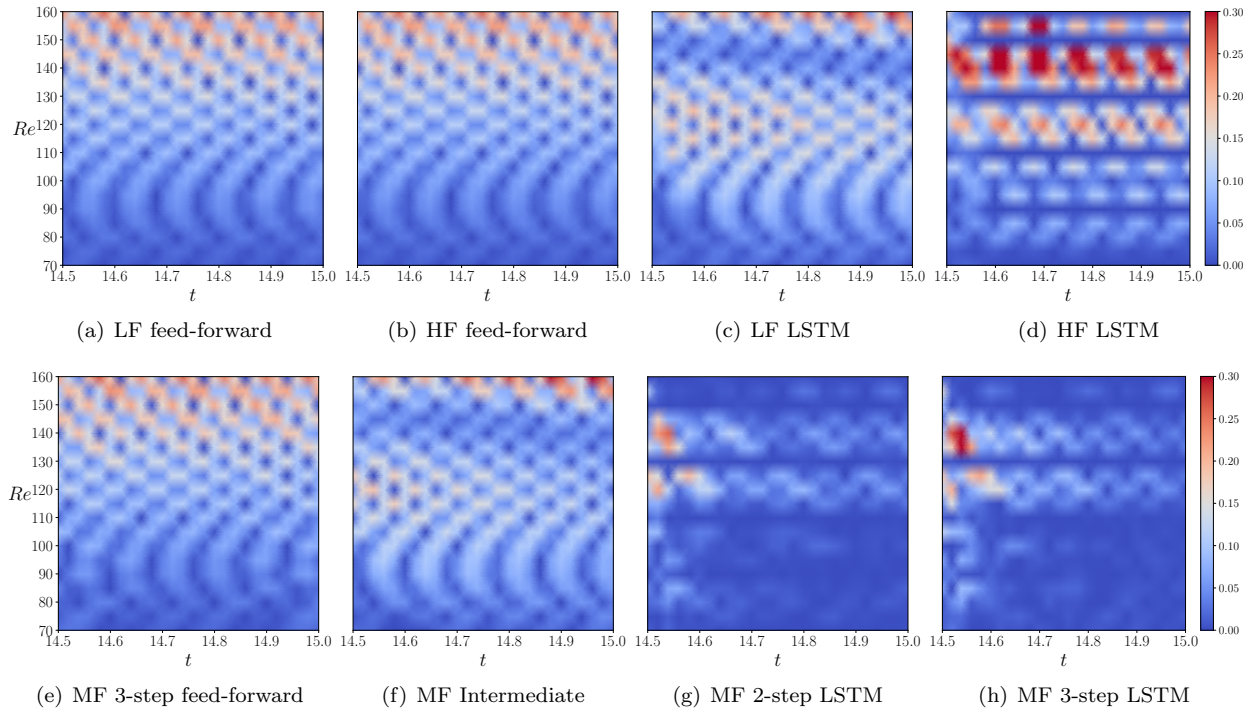


Figure 9: Absolute values of the errors in different NN model predictions of the drag coefficient with  $N_{\text{HF}}^{\mu} = 6$  HF training time series.

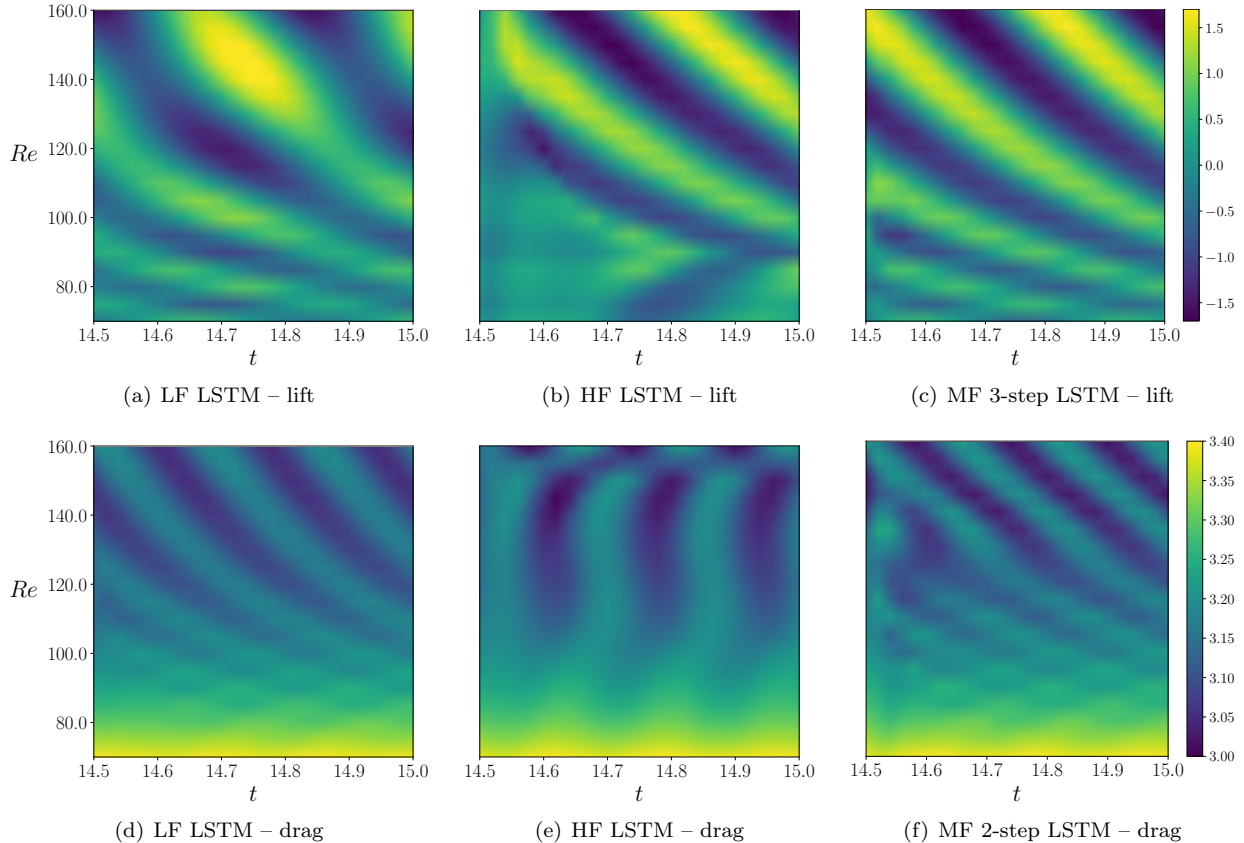


Figure 10: Single-fidelity (LF and HF) and the best performing MF LSTM regressions for the lift coefficient with  $N_{\text{HF}}^{\mu} = 10$  (top) and the drag coefficient with  $N_{\text{HF}}^{\mu} = 6$  (bottom).

### 5.3. Results and discussions 2: Prediction forward in time

So far, we have shown that the MF LSTM networks allow to accurately generalize/interpolate within the training region  $\mathcal{P} \times [t_0, T]$ . In this subsection we investigate the generalization capability of predicting forward in time, i.e., for the future states beyond the training range, which is known to be challenging for data-driven models.

Starting with the same training data set as in the previous subsection, we progressively move backward the final training time  $t^*$ , from  $t^* = T$  to  $t^* = t_0$ , while using the same instances of the Reynolds number. Though the training time interval is shortened, we still predict over the entire  $[t_0, T]$ , which includes both reconstruction and extrapolation in time. In Fig. 11, we report the prediction errors with respect to  $t^*$  through different LSTM network models. Given sufficient time steps to effectively merge LF and HF data, the MF LSTM model shows a consistent advantage over the single-fidelity models, both for the lift and drag coefficients (see Fig. 11). For example, in the drag case (Fig. 11(b)), a smaller test error is achieved with the 2-step LSTM trained on half of the whole time interval ( $t^* = 14.75$ ) than the single-fidelity LSTM, either HF or LF, with the full training set over time ( $t^* = T = 15.0$ ). This suggests that we can improve the predictive accuracy by using MF LSTM with only more than half of the training data, which may enable a significant reduction in the cost of HF data generation. Moreover, a description of uncertainty in the model predictions can be obtained with an ensemble-based technique [10]. Here, we train the last-step ( $NN_{\text{HF}}$ ) of the 2-step LSTM model multiple times with the same data and randomly generated network initialization, and then compute the statistical moments of the corresponding samples of predictions. This allows to construct uncertainty bounds for MF approximation, see Fig. 11. The low predictive uncertainty further highlights the consistent robustness of the proposed 2-step MF model as the extrapolation time window varies. In Fig. 12 we depict the approximation of both the lift and drag coefficients by the 2-step LSTM model with  $t^* = 14.80$ . The solution shows a very good agreement with the HF ground truth, even over the time interval

$[t^*, T]$  where we have no training data at all – neither HF nor LF.

Results have highlighted the effectiveness of the proposed MF LSTM networks in time-parameter-dependent surrogate modeling. For output quantities that present an oscillatory behaviour, the proposed models allow to learn the temporal pattern and predict the evolution of future states. Good generalization has been seen not only in the parameter ranges where few HF data are available, but also at future time instances outside the coverage of training data.

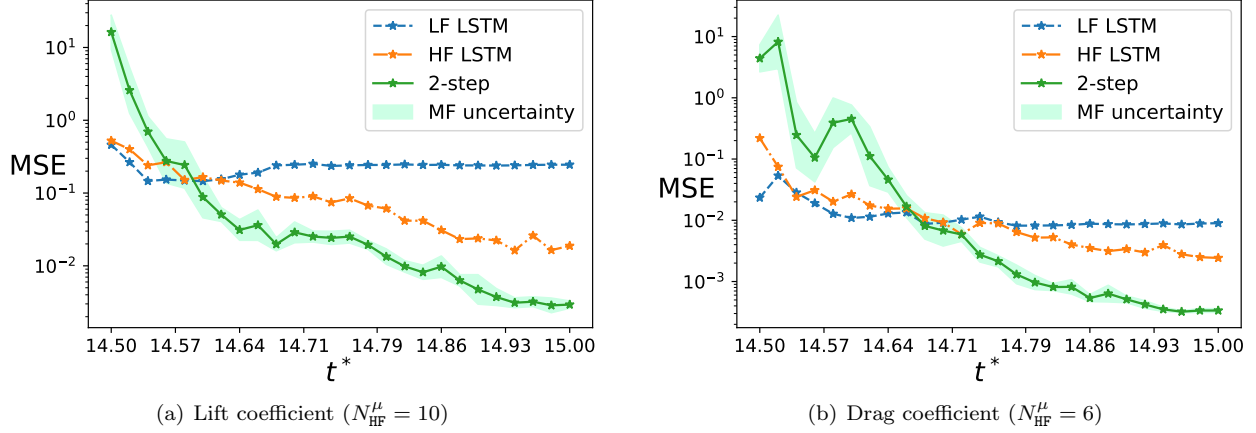


Figure 11: Test error on the whole time interval  $[t_0, T] = [14.5, 15.0]$  versus the final training time  $t^*$ . The testing time window  $[t_0, T]$  remains fixed, while training time length ( $t^* - t_0$ ) is increased as  $t^*$  goes from  $t_0$  to  $T$ . For each value of  $t^*$ , the second network of 2-step model is trained 20 times with random initialization, and the corresponding model predictions are used to compute uncertainty bounds on the test set through an ensemble-based technique. The test error mean (green) and  $\pm$  one standard deviation (shaded region) over the samples in the ensemble are shown, providing an uncertainty quantification for the MF predictions.

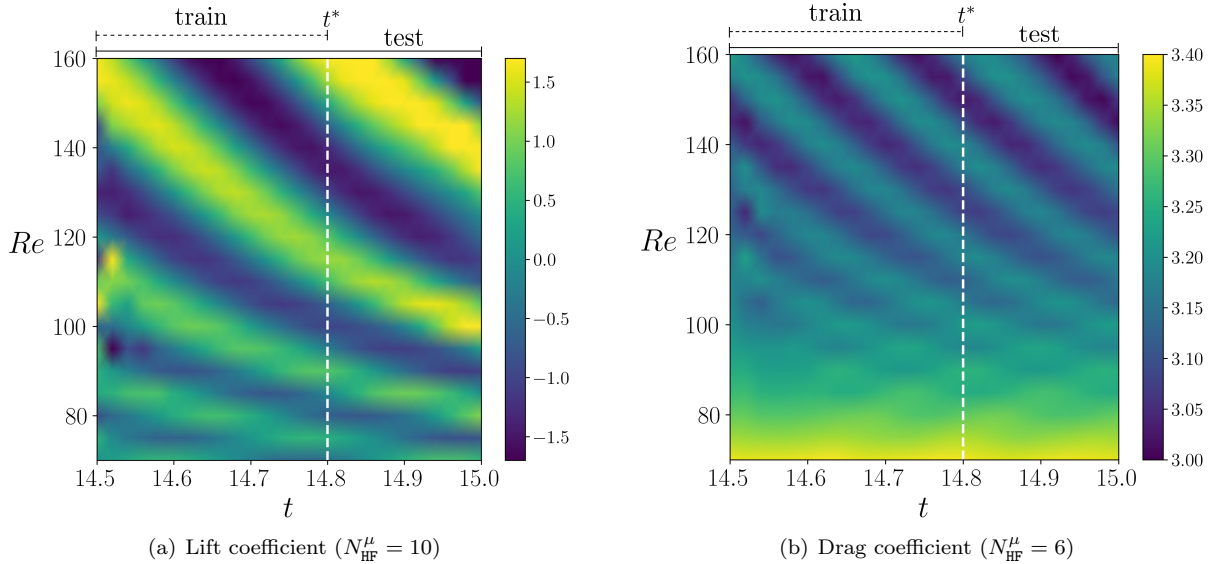


Figure 12: 2-step LSTM model prediction for the lift (left) and drag (right) coefficients. The model is trained with data up to  $t^* = 14.80$  (white dashed line) and tested over the whole time interval  $[14.5, 15.0]$ .



## 6. Numerical example (III): Lotka-Volterra system

In this section we consider a Lotka–Volterra system that describes a nonlinear, three-species prey-predator interaction:

$$\begin{cases} \frac{dy_1}{dt}(t) = y_1(t)(\mu - 0.1y_1(t) - 0.5y_2(t) - 0.5y_3(t)), & t \in (0, T), \\ \frac{dy_2}{dt}(t) = y_2(t)(-\mu + 0.5y_1(t) - 0.3y_3(t)), & t \in (0, T), \\ \frac{dy_3}{dt}(t) = y_3(t)(-\mu + 0.2y_1(t) + 0.5y_2(t)), & t \in (0, T), \\ y_i(0) = 0.5, & \forall i = 1, 2, 3. \end{cases} \quad (12)$$

Here we employ the proposed MF LSTM models to approximate vector-valued output quantities  $\mathbf{y}(t; \mu) = [y_1(t; \mu), y_2(t; \mu), y_3(t; \mu)]^T \in \mathbb{R}^3$ , i.e., the number of individuals in each population/species. In this MF regression task, we aim at estimating the system solution  $\mathbf{y}$  up to time  $T = 15.0$ , as the parameter  $\mu$  changes in  $\mathcal{P} = [1.0, 3.0]$ . Such a regression task is more challenging than those in the previous examples because the vector-valued quantity of interest  $\mathbf{y}$  exhibits aperiodic oscillations. Moreover, the parametric variation in  $\mu$  induces different aperiodic patterns in the amplitude and frequency of the oscillations. In this case, we only let the HF data cover a limited time window until  $T_{\text{HF}} < T$ , and intend to infer the aperiodic evolution of  $\mathbf{y}$  from the LF data that cover the entire domain  $\mathcal{D} = \mathcal{P} \times [0, T]$ .

### 6.1. Multi-fidelity setting

The training data are generated by the second-order Runge-Kutta (RK2) scheme with time steps  $\Delta t_{\text{LF}} = 0.25$  and  $\Delta t_{\text{HF}} = 0.0025$  for the LF and HF models, respectively. The LF data are taken at  $N_{\text{LF}}^\mu = 20$  uniform  $\mu$ -values over  $\mathcal{P} = [1.0, 3.0]$ . For each  $\mu$ -instance, we collect the LF training data along with the time integration with  $\Delta t_{\text{LF}} = 0.25$  up to the final time  $T_{\text{LF}} = T = 15$ . Similarly, the HF data are taken at  $N_{\text{HF}}^\mu = 10$  uniform  $\mu$ -values. Though the HF time integration is evaluated with  $\Delta t_{\text{HF}} = 0.0025$ , we only collect the HF training data every  $\Delta t_{\text{LF}} = 0.25$  to be consistent with the LF data, and stop the data collection at  $t = T_{\text{HF}} = 10$ . Thus, no HF data are available in the time interval  $[10.0, 15.0]$ . For the LSTM training, we group the data into batch subsequences of length  $K = 25$ . Fig. 13 shows the training data for  $\mu = 1.0$  and  $3.0$ , the minimum and maximum  $\mu$ -values. We note that the frequency of output oscillations increases with  $\mu$ . This implies that the LF coarse time discretization impacts the data quality more significantly for a larger value of  $\mu$ . In fact, at  $\mu = 3.0$  we spot a much larger discrepancy between the HF and LF data in comparison to the case of  $\mu = 1.0$ . We generate the test set similarly to the HF training data, but include  $N_{\text{test}} = 30$   $\mu$ -values in  $\mathcal{P}$  and reach the final time  $T = 15$ .

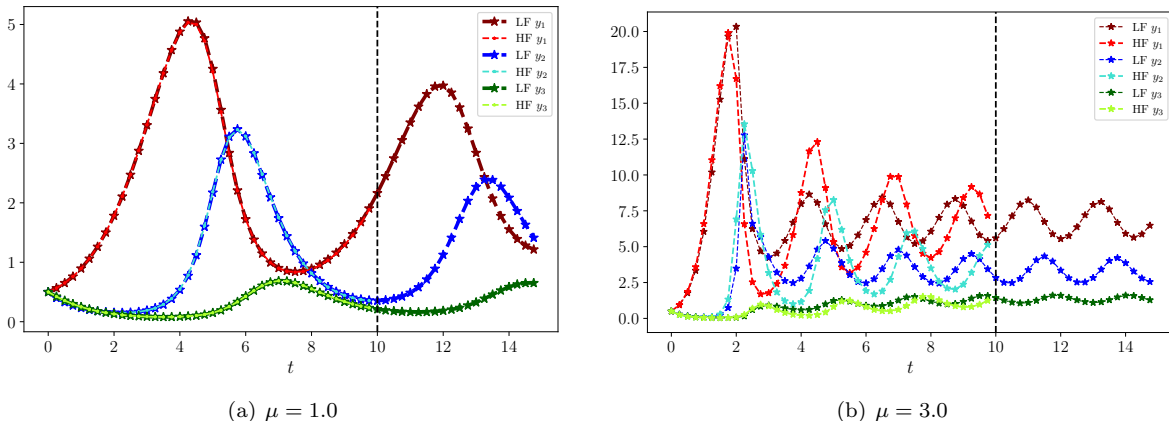


Figure 13: HF and LF training data for two parameter instances. The parameter  $\mu$  regulates both the amplitude and frequency of the oscillations of  $\mathbf{y}$ . As  $\mu$  grows, the amplitude and frequency increase, so does the discrepancy between the two data levels.

### 6.2. Results and discussions

Here we compare the single-fidelity and MF LSTM models on the given data sets. For the sake of brevity, we only present the results by the 2-step architecture among the MF LSTM networks. Table 4 collects the

prediction MSE on the test set that covers the whole domain  $\mathcal{D}$ , and we can clearly observe that the MF model outperforms the single-fidelity ones. Furthermore, Fig. 14 shows the comparison between the model predictions and the HF reference solution at a testing parameter instance  $\mu = 2.93$ .

The LF LSTM model fails to guarantee a low test error due to the poor data quality stemming from the coarse approximation. The HF LSTM model achieves a good accuracy in the range  $[0, T_{\text{HF}}] = [0, 10]$ , but cannot generalize to  $t > T_{\text{HF}}$  where no HF data are available for any  $\mu$ -value. This implies that even though HF LSTM can well interpolate over the parameter domain, it fails to extrapolate over time. This is understandable as the limited HF information is not sufficient for capturing the aperiodic oscillations or supporting the predictions of future states. On the other hand, our 2-step LSTM model performs well: it manages to both interpolate for the parameter  $\mu$  and predict forward in time. Here the full-domain LF data coverage over  $\mathcal{D}$  plays a critical role in enabling the MF time extrapolation, as the MF LSTM captures the aperiodic oscillations by well approximating the correlation between the two fidelity levels.

Table 4: Test errors for the single- and multi-fidelity models in example (III).

Model	LF LSTM	HF LSTM	MF 2-step LSTM
Test MSE	0.354	0.853	<b>0.043</b>

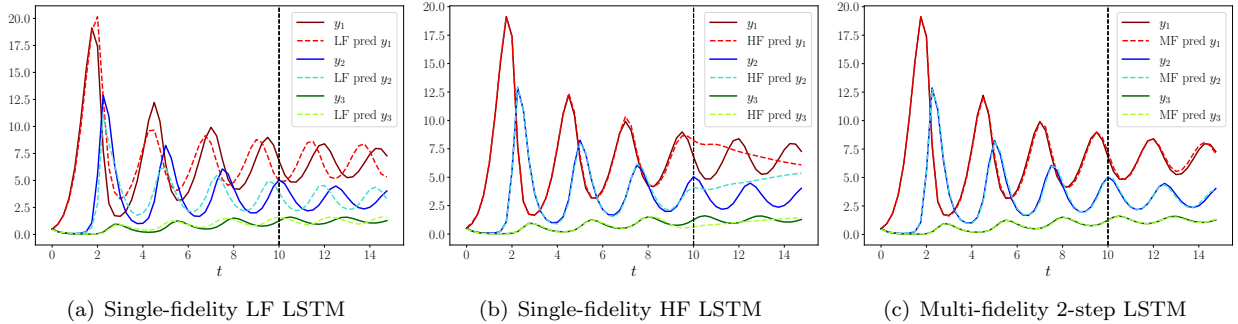


Figure 14: Model predictions by the single- and multi-fidelity LSTM networks compared to the HF solution ( $\mu = 2.93$ ). The black dashed line at  $T_{\text{HF}} = 10$  indicates the maximum time instant of the HF training data, i.e., no HF information is available over  $t \in [T_{\text{HF}}, T] = [10, 15]$ .

## 7. Extension to more fidelity levels

Although this work has been presented with bi-fidelity data, the proposed models in Sect. 2.2 can be extended to more than two fidelity levels. Let us consider  $M$  datasets  $\{\mathcal{T}_m = \{\mathbf{x}_m, \mathbf{y}_m\}\}_{m=1}^M$  given by a hierarchy of fidelity levels, sorted by increased accuracy of the data.

One can upgrade the proposed *multi-level* models in the following ways:

- *Series multi-level.* The number of “steps” is increased by concatenating multiple NNs – one for each fidelity level. This leads to a sequence of networks  $\{NN_m\}_{m=1}^M$ . The first network is fed with  $\mathcal{T}_1$  to learn  $\mathbf{f}_1$ , and then, sequentially, the  $m$ -th network  $NN_m$  is trained to approximate  $\mathbf{f}_m$  between the inputs  $[\mathbf{x}_m, \mathbf{f}_{m-1}(\mathbf{x}_m), \dots, \mathbf{f}_1(\mathbf{x}_m)]^T$  (all of them or a subset) and the outputs  $\mathbf{y}_m$ .
- *Parallel multi-level.* Each network of  $\{NN_m\}_{m=1}^{M-1}$  is trained independently on  $\mathcal{T}_m$  to approximate  $\mathbf{f}_m$ . The last network  $NN_M$  is fed with the inputs  $[\mathbf{x}_M, \mathbf{f}_{M-1}(\mathbf{x}_M), \dots, \mathbf{f}_1(\mathbf{x}_M)]^T$  (including all the other networks’ outputs) and the outputs  $\mathbf{y}_M$  to learn the HF function  $\mathbf{f}_M$ .
- A combination that mixes the *series* and *parallel* settings of networks.

The *series* approach should be advantageous over the *parallel* approach in the case where data come from sources of similar nature (e.g., refinement of numerical discretization) and each fidelity level adds further

information to the previous. Thus, it is meaningful to sequentially incorporate the features from previous fidelity levels into the prediction of the current level. On the other hand, the *parallel* strategy is more suitable for the case with multiple LF data sources that are not strongly correlated but individually provide useful information for the final HF prediction, so they are all together included in the inputs of  $NN_M$ . It is worth noting that the *parallel* strategy is expected to present better computational flexibility, because the first  $M - 1$  NNs can be trained independently and/or in parallel.

The “intermediate” model can be upgraded beyond bi-fidelity by locating the multiple levels of LF outputs in sequential hidden layers. Despite a similarity with the *series multi-level* extension, the upgraded “intermediate” model is determined by a single NN training whose loss function is a convex combination of all the losses of individual fidelity levels with coefficients  $\{\alpha_m \in \mathbb{R}^+ : 1 \leq m \leq M, \sum_{m=1}^M \alpha_m = 1\}$ , similar to (4) in Sect. 2.2. To regulate the contribution of each loss term, these coefficients can either be set manually or determined via the hyperparameter optimization.

It is sometimes the case that we may update the data for one of the fidelity levels or insert an additional level to an already trained MF model. The *parallel multi-level* approach hence only requires the network training for the updated/added data levels and the retraining of the last network, while freezing the other networks. With the “intermediate” setting or the *series multi-level* extension, however, the whole model must be retrained.

## 8. Concluding remarks

In this work we present several novel methods for the estimation of time-parameter-dependent quantities of interest using multi-fidelity techniques with LSTM neural networks. The proposed techniques are shown to be advantageous over both the single-fidelity neural networks and the multi-fidelity ones without LSTM units, especially when approximating time-dependency. The multi-fidelity LSTM strategy enables accurate estimation of the target quantities’ time evolution at a reasonable computational cost, as it leverages many low-fidelity, easily attainable data and only requires a limited number of high-fidelity, expensive data. The non-intrusive nature of the proposed models guarantees a wide applicability, which has been exemplified by a diverse collection of engineering applications. On the other hand, the multi-fidelity LSTM models show excellent predictive capabilities and generalization performance both in time and over parametric variation, which is typically deemed challenging for data-driven, non-intrusive surrogate models.

The proposed multi-fidelity LSTM models can be straightforwardly extended to more than two fidelity levels, as well as to the involvement of multiple physical systems, for which we can collect multi-fidelity data from several interacting systems governing the quantities of interest. Another promising, but more challenging, future application is the full-field approximation of PDE solutions. In this regard, a viable strategy is to consider the time-parameter-dependent coefficients of a low-dimensional reduced basis (e.g., through proper orthogonal decomposition on a set of HF snapshot solutions) as output quantities for the multi-fidelity LSTM models.

## Acknowledgment

The second author is financially supported by Sectorplan Bèta (the Netherlands) under the focus area *Mathematics of Computational Science*. The third author acknowledges the support from Fondazione Cariplo under the Grant n. 2019-4608. The authors would like to express their appreciation to Dr. Stefania Fresca for the fruitful discussions and for her help with numerical implementations.

## Appendix: hyperparameter summary

Features of an NN model that cannot be optimized during the training process are called *hyperparameters*. They are either related to the network structure (e.g., the numbers of layers and nodes) or associated with the training (e.g., the learning rate). The performance of an NN highly relies on the hyperparameter choices, and thus it is important to find an optimal set of hyperparameter values. In this work, we employ a Bayesian hyperparameter optimization technique that minimizes an objective function  $\mathcal{O}$  in a multi-dimensional domain  $\Lambda = \Lambda_1 \times \Lambda_2 \times \dots \times \Lambda_N$  of all  $N$  hyperparameters. Here  $\Lambda_i$  represents the range of the  $i$ -th hyperparameter value.  $\mathcal{O}$  is defined as the cross-validation error on the training set, because generating a HF validation set would be expensive and data reuse thus becomes necessary. Specifically, we use the Python package Hyperopt [4] (see [5, 6] for more details).

To ensure a fair comparison among different models, the same choice of domain  $\Lambda$  and objective function  $\mathcal{O}$  is considered in the hyperparameter optimization for each model, unless otherwise specified. The hyperparameters are estimated by a common Bayesian optimization method [4, 5], which balances between the minimization of error  $\mathcal{O}$  and the exploration of domain  $\Lambda$ . In the presented numerical examples, the performance of different models is assessed with optimized hyperparameter values.

Tables A.1, A.2, A.3 and A.4 collect the optimized HPs in all the numerical examples. For notation, we let  $\eta$  denote the learning rate, and let ‘depth’ and ‘width’ represent the numbers of layers and nodes, respectively. For the multi-level models we enumerate the level steps with the variable ‘step’. Although the intermediate model consists of a single NN, we use the variable ‘step’ to denote the ‘LF’ (resp. ‘HF’) portion of the network, i.e., the first (resp. second) part of the network dedicated to the estimation of the LF (resp. HF) output (see Fig. 2(c) in Section 2). We recall that  $\alpha$  is the coefficient regulating the contributions of the two fidelity levels to the loss function in the intermediate LSTM model. Moreover, all the NNs use the hyperbolic tangent activation function.

Table A.1: Optimized HP values in example (I) (Section 4).

Model	Step	Depth $\times$ width LSTM	Depth $\times$ width Dense	Optimizer	$\eta$	Batch size	$\alpha$
LF feed-forward	-	-	4 $\times$ 61	Adam	$2.54 \times 10^{-2}$	13	-
HF feed-forward	-	-	3 $\times$ 52	Adam	$1.98 \times 10^{-2}$	3	-
LF LSTM	-	1 $\times$ 98	0 $\times$ 0	Adam	$7.14 \times 10^{-3}$	1	-
HF LSTM	-	2 $\times$ 94	0 $\times$ 0	Adamax	$2.94 \times 10^{-2}$	4	-
3-step feed-forward	1	-	4 $\times$ 61	Adam	$2.54 \times 10^{-2}$	13	-
	2	-	3 $\times$ 41	Adamax	$3.89 \times 10^{-3}$	67	-
	3	-	1 $\times$ 32	Adam	$1.78 \times 10^{-4}$	151	-
2-step LSTM	1	1 $\times$ 98	0 $\times$ 0	Adam	$7.14 \times 10^{-3}$	1	-
	2	2 $\times$ 98	2 $\times$ 16	Adamax	$1.78 \times 10^{-3}$	3	-
3-step LSTM	1	1 $\times$ 98	0 $\times$ 0	Adam	$7.14 \times 10^{-3}$	1	-
	2	4 $\times$ 20	0 $\times$ 0	Adam	$1.80 \times 10^{-2}$	3	-
	3	3 $\times$ 82	3 $\times$ 20	Adamax	$6.98 \times 10^{-3}$	4	-
Intermediate	LF	3 $\times$ 62	0 $\times$ 0	Adam	$3.81 \times 10^{-4}$	27	0.51
	HF	2 $\times$ 72	1 $\times$ 118	”	”	”	”

Table A.2: Optimized HP values for the lift coefficient in example (II) (Section 5).

Model	Step	Depth $\times$ width LSTM	Depth $\times$ width Dense	Optimizer	$\eta$	Batch size	$\alpha$
LF feed-forward	-	-	$2 \times 20$	Adamax	$3.52 \times 10^{-4}$	177	-
HF feed-forward	-	-	$2 \times 34$	Adamax	$1.49 \times 10^{-4}$	149	-
LF LSTM	-	$4 \times 128$	$1 \times 46$	Adam	$1.96 \times 10^{-3}$	1	-
HF LSTM	-	$4 \times 24$	$0 \times 0$	Adamax	$5.46 \times 10^{-2}$	6	-
3-step feed-forward	1	-	$2 \times 20$	Adamax	$3.52 \times 10^{-4}$	177	-
	2	-	$4 \times 62$	Adamax	$5.52 \times 10^{-4}$	266	-
	3	-	$1 \times 44$	Adam	$8.73 \times 10^{-3}$	128	-
2-step LSTM	1	$4 \times 128$	$1 \times 46$	Adam	$1.96 \times 10^{-3}$	1	-
	2	$3 \times 64$	$2 \times 26$	Adam	$2.66 \times 10^{-3}$	5	-
3-step LSTM	1	$4 \times 128$	$1 \times 46$	Adam	$1.96 \times 10^{-3}$	1	-
	2	$3 \times 18$	$0 \times 0$	Adamax	$1.17 \times 10^{-2}$	6	-
	3	$4 \times 24$	$0 \times 0$	Adamax	$6.20 \times 10^{-4}$	8	-
Intermediate	LF	$1 \times 106$	$0 \times 0$	Adam	$2.65 \times 10^{-2}$	16	0.64
	HF	$3 \times 20$	$0 \times 0$	"	"	"	"

Table A.3: Optimized HP values for the drag coefficient in example (II) (Section 5).

Model	Step	Depth $\times$ width LSTM	Depth $\times$ width Dense	Optimizer	$\eta$	Batch size	$\alpha$
LF feed-forward	-	-	$4 \times 52$	Adam	$8.91 \times 10^{-2}$	162	-
HF feed-forward	-	-	$1 \times 38$	Adam	$1.71 \times 10^{-4}$	74	-
LF LSTM	-	$3 \times 36$	$1 \times 112$	Adam	$1.25 \times 10^{-3}$	1	-
HF LSTM	-	$3 \times 36$	$0 \times 0$	Adamax	$2.60 \times 10^{-2}$	6	-
3-step feed-forward	1	-	$4 \times 52$	Adam	$8.91 \times 10^{-2}$	162	-
	2	-	$2 \times 22$	Adamax	$5.74 \times 10^{-2}$	158	-
	3	-	$1 \times 39$	Adamax	$3.32 \times 10^{-3}$	107	-
2-step LSTM	1	$3 \times 36$	$1 \times 112$	Adam	$1.25 \times 10^{-3}$	1	-
	2	$3 \times 78$	$1 \times 50$	Adam	$2.87 \times 10^{-3}$	10	-
3-step LSTM	1	$3 \times 36$	$1 \times 112$	Adam	$1.25 \times 10^{-3}$	1	-
	2	$3 \times 34$	$0 \times 0$	Adam	$2.06 \times 10^{-3}$	10	-
	3	$1 \times 24$	$0 \times 0$	Adamax	$5.05 \times 10^{-2}$	8	-
Intermediate	LF	$2 \times 46$	$0 \times 0$	Adam	$2.87 \times 10^{-2}$	11	0.1
	HF	$3 \times 110$	$2 \times 22$	"	"	"	"

Table A.4: HP values in example (III) (Section 6). The same structure is used in all the networks, showing that the MF approach is robustly advantageous over the single-fidelity ones, regardless of the hyperparameter tuning.

Model	Step	Depth $\times$ width LSTM	Depth $\times$ width Dense	Optimizer	$\eta$	Batch size
LF LSTM	-	$3 \times 64$	$1 \times 32$	Adamax	$1.00 \times 10^{-3}$	50
HF LSTM	-	$3 \times 64$	$1 \times 32$	Adamax	$1.00 \times 10^{-3}$	50
2-step LSTM	1	$3 \times 64$	$1 \times 32$	Adamax	$1.00 \times 10^{-3}$	50
	2	$3 \times 64$	$1 \times 32$	Adamax	$5.00 \times 10^{-3}$	100

## References

- [1] M. A. Álvarez, L. Rosasco, and N. D. Lawrence. Kernels for vector-valued functions: A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266, 2012.
- [2] N. Baker, F. Alexander, T. Bremer, A. Hagberg, Y. Kevrekidis, H. Najm, M. Parashar, A. Patra, J. Sethian, S. Wild, et al. Workshop report on basic research needs for scientific machine learning: Core technologies for artificial intelligence. Technical report, USDOE Office of Science (SC), Washington, DC (United States), 2019.
- [3] C. Beck, A. Jentzen, et al. Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations. *Journal of Nonlinear Science*, 29(4):1563–1619, 2019.
- [4] J. Bergstra. Hyperopt: Distributed asynchronous hyperparameter optimization in python, 2013.
- [5] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. *Advances in Neural Information Processing Systems*, 24, 2011.
- [6] J. Bergstra, D. Yamins, and D. D. Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML’13, pages 115–123. JMLR.org, 2013.
- [7] K. Champion, B. Lusch, J. N. Kutz, and S. L. Brunton. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019.
- [8] K. Cutajar, M. Pullin, A. Damianou, N. Lawrence, and J. González. Deep Gaussian processes for multi-fidelity modeling. *arXiv:1903.07320*, 2019.
- [9] N. Demo, M. Strazzullo, and G. Rozza. An extended physics informed neural network for preliminary analysis of parametric optimal control problems. *arXiv:2110.13530*, 2021.
- [10] T. G. Dietterich. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, pages 1–15. Springer, 2000.
- [11] R. FitzHugh. Impulses and physiological states in theoretical models of nerve membrane. *Biophysical Journal*, 1(6):445–466, 1961.
- [12] S. Fresca, L. Dede, and A. Manzoni. A comprehensive deep learning-based approach to reduced order modeling of nonlinear time-dependent parametrized pdes. *Journal of Scientific Computing*, 87(2):1–36, 2021.
- [13] S. Fresca and A. Manzoni. Real-time simulation of parameter-dependent fluid flows through deep learning-based reduced order models. *Fluids*, 6(7), 2021.
- [14] S. Fresca and A. Manzoni. POD-DL-ROM: enhancing deep learning-based reduced order models for nonlinear parametrized pdes by proper orthogonal decomposition. *Computer Methods in Applied Mechanics and Engineering*, 388:114181, 2022.
- [15] S. Fresca, A. Manzoni, L. Dedè, and A. Quarteroni. Deep learning-based reduced order models in cardiac electrophysiology. *PLOS One*, 15(10):e0239416, 2020.
- [16] N. Geneva and N. Zabaras. Multi-fidelity generative deep learning turbulent flows. *Foundations of Data Science*, 2(4):391–428, 2020.
- [17] F. A. Gers, J. Schmidhuber, and F. Cummins. Learning to forget: Continual prediction with lstm. *Neural Computation*, 12(10):2451–2471, 2000.

- [18] A. Graves, S. Fernández, and J. Schmidhuber. Bidirectional lstm networks for improved phoneme classification and recognition. In *International conference on artificial neural networks*, pages 799–804. Springer, 2005.
- [19] M. Guo. A brief note on understanding neural networks as Gaussian processes. *arXiv:2107.11892*, 2021.
- [20] M. Guo, A. Manzoni, M. Amendt, P. Conti, and J. S. Hesthaven. Multi-fidelity regression using artificial neural networks: efficient approximation of parameter-dependent output quantities. *Computer methods in Applied Mechanics and Engineering*, 389:114378, 2022.
- [21] J. Han, A. Jentzen, and W. E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [22] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [23] A. A. Howard, M. Perego, G. E. Karniadakis, and P. Stinis. Multifidelity deep operator networks. *arXiv:2204.09157*, 2022.
- [24] M. Jain, K. Brown, and A. K. Sadek. Multi-fidelity recursive behavior prediction. *arXiv:1901.01831*, 2018.
- [25] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
- [26] S. Khairy and P. Balaprakash. Multifidelity reinforcement learning with control variates. *arXiv:2206.05165*, 2022.
- [27] J. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. Deep neural networks as Gaussian processes. *arXiv:1711.00165*, 2017.
- [28] K. Lee and K. T. Carlberg. Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders. *Journal of Computational Physics*, 404:108973, 2020.
- [29] L. Lu, R. Pestourie, S. G. Johnson, and G. Romano. Multifidelity deep neural operators for efficient learning of partial differential equations with application to fast inverse design of nanoscale heat transport. *arXiv:2204.06684*, 2022.
- [30] X. Meng, H. Babae, and G. E. Karniadakis. Multi-fidelity Bayesian neural networks: Algorithms and applications. *Journal of Computational Physics*, 438:110361, 2021.
- [31] X. Meng and G. E. Karniadakis. A composite neural network that learns from multi-fidelity data: Application to function approximation and inverse pde problems. *Journal of Computational Physics*, 401, 2019.
- [32] M. Motamed. A multi-fidelity neural network surrogate sampling method for uncertainty quantification. *International Journal for Uncertainty Quantification*, 10(4), 2020.
- [33] J. Nagumo, S. Arimoto, and S. Yoshizawa. An active pulse transmission line simulating nerve axon. *Proceedings of the IRE*, 50(10):2061–2070, 1962.
- [34] F. Negri. redbKIT Version 2.2. <http://redbkit.github.io/redbKIT/>, 2016.
- [35] A. O’Hagan and M. C. Kennedy. Predicting the output from a complex computer code when fast approximations are available. *Biometrika*, 87(1):1–13, 2000.
- [36] C. Olah. Understanding LSTM networks. 2015.
- [37] S. Pagani, A. Manzoni, and A. Quarteroni. Numerical approximation of parametrized problems in cardiac electrophysiology by a local reduced basis method. *Computer Methods in Applied Mechanics and Engineering*, 340:530–558, 2018.

- [38] L. Partin, G. Geraci, A. Rushdi, M. S. Eldred, and D. E. Schiavazzi. Multifidelity data fusion in convolutional encoder/decoder networks. *arXiv:2205.05187*, 2022.
- [39] B. Peherstorfer, K. Willcox, and M. Gunzburger. Survey of multifidelity methods in uncertainty propagation, inference, and optimization. *SIAM Review*, 60(3):550–591, 2018.
- [40] M. Raissi and G. Karniadakis. Deep multi-fidelity Gaussian processes. *arXiv:1604.07484*, 2016.
- [41] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [42] D. H. Song and D. M. Tartakovsky. Transfer learning on multi-fidelity data. *arXiv:2105.00856*, 2021.
- [43] M. Sundermeyer, R. Schlüter, and H. Ney. LSTM neural networks for language modeling. In *the 13th Annual Conference of the International Speech Communication Association*, 2012.
- [44] M. Torzoni, A. Manzoni, and S. Mariani. Health monitoring of civil structures: A MCMC approach based on a multi-fidelity deep neural network surrogate. In *the 1st Online Conference on Algorithms*, 2021.
- [45] M. Torzoni, A. Manzoni, and S. Mariani. A deep neural network, multi-fidelity surrogate model approach for Bayesian model updating in shm. In *European Workshop on Structural Health Monitoring*, pages 1076–1086. Springer, 2023.
- [46] B. B. Traore, B. Kamsu-Foguem, and F. Tangara. Deep convolution neural network for image recognition. *Ecological Informatics*, 48:257–268, 2018.
- [47] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- [48] Wikimedia Commons. File:long short-term memory.svg — wikimedia commons, the free media repository, 2020. [Online; accessed 01-October-2022].
- [49] L. Zheng, T. L. Hedrick, and R. Mittal. A multi-fidelity modelling approach for evaluation and optimization of wing stroke aerodynamics in flapping flight. *Journal of Fluid Mechanics*, 721:118–154, 2013.