# DPM: A deep learning PDE augmentation method
# (with application to large-eddy simulation)

Jonathan B. Freund[*], Jonathan F. MacArt[†], and Justin Sirignano[‡§]

November 22, 2019

## Abstract

Machine learning for scientific applications faces the challenge of limited data. We propose a framework that leverages *a priori* known physics to reduce overfitting when training on relatively small datasets. A deep neural network is embedded in a partial differential equation (PDE) that expresses the known physics and learns to describe the corresponding unknown or unrepresented physics from the data. Crafted as such, the neural network can also provide corrections for erroneously represented physics, such as discretization errors associated with the PDE's numerical solution. Once trained, the deep learning PDE model (DPM) can make out-of-sample predictions for new physical parameters, geometries, and boundary conditions.

Our approach optimizes over the functional form of the PDE. Estimating the embedded neural network requires optimizing over the entire PDE, which itself is a function of the neural network. Adjoint partial differential equations are used to efficiently calculate the high-dimensional gradient of the objective function with respect to the neural network parameters. A stochastic adjoint method (SAM), similar in spirit to stochastic gradient descent, further accelerates training.

The approach is demonstrated and evaluated for turbulence predictions using large-eddy simulation (LES), a filtered version of the Navier–Stokes equation containing unclosed sub-filter-scale terms. High-fidelity direct numerical simulations (DNS) of decaying isotropic turbulence provide the training and testing data. The DPM outperforms the widely-used constant-coefficient and dynamic Smagorinsky models, even for filter sizes so large that these established models become qualitatively incorrect. It also significantly outperforms *a priori* trained models, which do not account for the full PDE. For comparable accuracy, the overall cost is reduced. Simulations of the DPM are accelerated by efficient GPU implementations of network evaluations. Measures of discretization errors, which are well-known to be consequential in LES, suggest that the ability of the training formulation to correct for these errors

---

[*]Mechanical Science & Engineering and Aerospace Engineering, University of Illinois at Urbana–Champaign, jbfreund@illinois.edu

[†]The Center for Exascale Simulation of Plasma-coupled Combustion, Coordinated Science Laboratory, University of Illinois at Urbana–Champaign, jmacart@illinois.edu

[‡]Department of Industrial & Systems Engineering, University of Illinois at Urbana–Champaign, jasirign@illinois.edu

[§]The author list is alphabetical.

is crucial to its success. A relaxation of the discrete enforcement of the divergence-free constraint is also considered, instead allowing the DPM to approximately enforce incompressibility physics.

# 1 Introduction

It is well-understood that machine learning in science and engineering is challenged by limited availability of data. Experiments can be expensive and time-consuming, and some quantities are difficult or impossible to measure with current techniques. Similarly, high-fidelity numerical simulations are computationally costly or even infeasible for many applications. To reduce overfitting on relatively small datasets, our framework leverages the *a priori* known and representable physics by embedding a deep neural network $\mathbf{h}_\theta(u, v, w)$ in a PDE description of the known and represented physics. This yields the class of PDE models:

$$\frac{\partial \mathbf{u}}{\partial t}(t, \mathbf{x}) = \underbrace{\mathbf{f}_\nu\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big)}_{\text{Known and represented physics}} + \underbrace{\mathbf{h}_\theta\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big)}_{\text{Unknown/unrepresented physics}}, \qquad (1.1)$$

where $\mathbf{u}(t, \mathbf{x}) = \big(u_1(t, \mathbf{x}), \ldots, u_d(t, \mathbf{x})\big) \in \mathbb{R}^d$, $\mathbf{x} = (x_1, \ldots, x_q) \in \Omega \subset \mathbb{R}^q$, $\mathbf{u_x}(t, \mathbf{x}) = \big(u_{1,\mathbf{x}}(t, \mathbf{x}), \ldots, u_{d,\mathbf{x}}(t, \mathbf{x})\big) \in \mathbb{R}^{d \times q}$, and $\mathbf{u_{xx}}(t, \mathbf{x}) = \big(u_{1,x_1 x_1}(t, \mathbf{x}), \ldots, u_{1,x_q x_q}(t, \mathbf{x}), \ldots, u_{d,x_1 x_1}(t, \mathbf{x}), \ldots, u_{d,x_q x_q}(t, \mathbf{x})\big) \in \mathbb{R}^{d \times q}$. The term $\mathbf{f}_\nu : \mathbb{R}^{(2q+1)d} \to \mathbb{R}^d$ represents the known physics in the PDE, with $\nu$ the physical and scenario parameters that specify the application. Generalization to mixed-second-order and higher-order derivatives is straightforward. Appropriate initial and boundary conditions complete the PDE. The neural network $\mathbf{h}_\theta : \mathbb{R}^{(2q+1)d} \to \mathbb{R}^d$ will learn to describe the unknown or unrepresented physics (more precisely defined below) using the neural network parameters $\theta \in \mathbb{R}^{d_\theta}$, which are estimated from the data.

Two classes of unknown or unrepresented physics are anticipated. In the first, the governing equations are too computationally expensive to solve. In this case, deep learning can be used to develop a reduced model. For the Navier–Stokes turbulence example that we consider, the reduced model includes a sub-grid-scale stress closure that attempts to represent the physical effects of small unresolved turbulence scales on those represented in the large-eddy simulation (LES). The second case occurs when important physics is omitted, either because it is deemed too costly (for example, multi-component diffusion in reacting flows) or if it is truly unknown. In general, unknown or unrepresented physics will diminish the PDE model's accuracy. In such cases, the selected PDE model is incomplete although not incorrect, and thus it can be used as a building block for the development of machine learning models. The proposed deep learning formulation is designed to augment these governing equations based on additional though limited high-fidelity simulation or experimental data.

Equation (1.1) is designated a deep learning PDE model (DPM). The model parameters $\theta$ are estimated from trusted data $\mathbf{V}(t, \mathbf{x}; \nu)$, which are assumed to be available at certain times $t_1, \ldots, t_{N_t}$ and for certain physical or scenario parameters $\nu_1, \ldots, \nu_M$. These are used to numerically estimate a neural network $\mathbf{h}_\theta$ that generates the solution $\mathbf{u}$ to (1.1) that most closely matches the trusted data $\mathbf{V}$. This is done by minimizing

the objective function

$$L(\theta) = \sum_{m=1}^{M} \sum_{n=1}^{N_t} \int_{\Omega} \|\mathbf{u}(t_n, \mathbf{x}; \nu_m) - \mathbf{V}(t_n, \mathbf{x}; \nu_m)\| \, d\mathbf{x}, \tag{1.2}$$

where we explicitly denote the dependence of $\mathbf{u}$ on $\nu$ via the notation $\mathbf{u}(t, \mathbf{x}; \nu)$.

Minimizing (1.2) is challenging since it is a function of the PDE (1.1), which in turn is a nonlinear function of the neural network parameters $\theta$. A stochastic adjoint method (SAM) is proposed for computational efficiency (Section 2). It accelerates the optimization by calculating gradients with respect to the neural network parameters using adjoint PDEs. Although the dimension of $\theta$ can be large (easily $\gtrsim 10^5$ parameters), the number of adjoint PDEs matches the number of PDEs $d$ in (1.1).

Once trained on available datasets, the deep learning PDE model can be used for out-of-sample predictions in new scenarios with different $\nu$ or different boundary conditions. As a demonstration, we evaluate the accuracy of the DPM for out-of-sample initial conditions and physical parameters for the challenging case of LES of turbulence. In LES, the nonlinear Navier–Stokes equations are filtered and, to accelerate numerical solutions, are solved on coarse grids that purposefully do not resolve the full range of turbulence scales (see Section 3.1). This reduces the computational cost compared to full-resolution direct numerical simulation (DNS), making LES tractable for many problems of engineering interest. However, the filtering step introduces unclosed terms into the LES equations (Section 3.2) that correspond to unrepresented physics in (1.1). The accuracy of LES strongly depends upon accurately modeling this unrepresented physics.

LES truncates the turbulence energy spectrum at or near the mesh resolution, and it is therefore inherently linked to discretization errors, which are largest for the smallest scales closest to the nominal truncation point. Some estimates of discretization errors are made in Section 4. Coping with modeling errors concurrently with discretization error is a fundamental challenge in LES [1–3]. The DPM training (Section 5) is able to correct for this discretization error, just as if it were incorrect or missing physics. Overall, it recovers resolved spectra on significantly coarser LES meshes than established models (Section 5.4).

The application of deep learning is an exciting direction in scientific computing [4–16]. Especially relevant recent efforts include closure models for the Reynolds-averaged Navier–Stokes (RANS) equations using an *a priori* estimation method in which the optimization is de-coupled from the PDE [11, 12]. Although this approach works for RANS, which is simpler in the sense that only a time-averaged quantity is sought for the closure of time-averaged PDEs, it performs poorly for LES (Section 5.5), probably due to the more complex link between the resolved turbulence dynamics and the necessarily unsteady sub-grid-scale energy transfer. Wang *et al.* [14] developed closure models for LES using *a priori* training.

Nonlinear physics with a strong coupling between the unknown terms (to be estimated) and the output variables motivates the proposed optimization over the entire PDE (1.1). Berg and Nystrom [13] also developed an adjoint-based method for estimating a neural network coefficient function in a PDE from data, demonstrating it for a Poisson equation. In that case, the function does not depend upon the PDE solution. In our case, adjoint equations are used to estimate the functional form of the nonlinear PDEs. In addition, Holland *et al.* [17] recently used discrete adjoint methods to train a neural network for closure of RANS.

The stochastic adjoint method to accelerate training is introduced in Section 2. The DPM is then formulated for large-eddy simulation in Section 3. Numerical discretization errors relevant to any LES are quantified and discussed in Section 4. A specific numerical application to large-eddy simulation of isotropic turbulence is introduced and analyzed in Section 5, which includes discussion of the relatively poor performance of the corresponding *a priori* approach (Section 5.5).

## 2 Stochastic Adjoint Method

### 2.1 Adjoint-based PDE gradient

The objective function (1.2) can be iteratively minimized using gradient descent,

$$\theta_{k+1} = \theta_k - \alpha_k \nabla_\theta L(\theta_k), \tag{2.1}$$

where $\alpha_k$ is the so-called learning rate. This requires calculating the gradient $\nabla_\theta L(\theta)$, which is computationally challenging since $L(\theta)$ depends upon the solution of the PDE (1.1). A naïve approach would directly apply the gradient $\nabla_\theta$ to (1.1) and derive a PDE for $\nabla_\theta \mathbf{u}$. However, the dimension of this PDE would match the dimension $d_\theta$ of the neural network parameters $\theta$, which is large (often over $10^5$ parameters), so such an approach is computationally intractable. Calculating $\nabla_\theta L$ using numerical differentiation with finite difference approximations would similarly require an intractable number of evaluations.

The adjoint PDE provides a computationally-efficient method for calculating $\nabla_\theta L(\theta)$. It requires solving only $p$ PDEs per parameter update, which matches the dimension of $\mathbf{u}$ in (1.1). For any deep learning model, $p$ will be many times smaller than $d_\theta$. The adjoint PDE can be viewed as a continuous-time PDE version of the usual neural network backpropagation algorithm. For notational convenience we start by considering $M = 1$ and $d = q = 1$. The generalization to $M, d, q \geq 1$ is straightforward (see Section 2.3). For $x \in \Omega$ and $0 \leq t < t_N$,

$$
\begin{aligned}
-\frac{\partial \widehat{u}}{\partial t} = {}& \widehat{u} \nabla_u f_{\nu_m}(u, u_x, u_{xx}) \\
& - \frac{\partial}{\partial x}\left[\widehat{u} \nabla_v f_{\nu_m}(u, u_x, u_{xx})\right] \\
& + \frac{\partial^2}{\partial x^2}\left[\widehat{u} \nabla_w f_{\nu_m}(u, u_x, u_{xx})\right] \\
& + \widehat{u} \nabla_u h_\theta(u, u_x, u_{xx}) \\
& - \frac{\partial}{\partial x}\left[\widehat{u} \nabla_v h_\theta(u, u_x, u_{xx})\right] \\
& + \frac{\partial^2}{\partial x^2}\left[\widehat{u} \nabla_w h_\theta(u, u_x, u_{xx})\right] \\
& + \sum_{n=1}^{N-1} \delta(t - t_n) \nabla_u \left\| u(t_n, x; \nu_m) - V(t_n, x; \nu_m) \right\|,
\end{aligned}
$$

4

with the final condition

$$\widehat{u}(t_N, x) \quad = \quad \nabla_u \left\| u(t_N, x; \nu_m) - V(t_N, x; \nu_m) \right\|. \tag{2.2}$$

The gradient of the objective function (1.2) is then

$$\nabla_\theta L(\theta) \quad = \quad \int_0^{t_N} \int_\Omega \widehat{u}(t, x) \nabla_\theta h_\theta \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \, dx dt. \tag{2.3}$$

For $M > 1$, $M$ such adjoint PDEs must be solved. The adjoint equations are generalized to $d, q \geq 1$ in Section 2.3.

## 2.2   Stochastic optimization

When $M$ and $N_t$ are large, the adjoint PDE method still requires significant computation time per gradient descent iteration (2.1). To accelerate training, we introduce a stochastic adjoint method (SAM), which randomly samples time intervals, allowing for a larger number of training iterations per computational time. Its steps are:

- Select uniformly at random a scenario $m \in \{1, \ldots, M\}$ and time $n \in \{0, \ldots, N-1\}$,

- Solve (1.1) on $[t_n, t_{n+1}]$ with initial condition $V_{(t_n, x; \nu_m)}$ taken from available data,

- Evaluate the objective function

$$J(\theta) = \int_\Omega \left\| u(t_{n+1}, x; \nu_m) - V(t_{n+1}, x; \nu_m) \right\| \, dx, \tag{2.4}$$

- Calculate $\nabla_\theta J(\theta)$ via its (time-reversed) adjoint PDE, which satisfies, for $x \in \Omega$ and $t_n \leq t < t_{n+1}$,

$$
\begin{aligned}
-\frac{\partial \widehat{u}}{\partial t}(t, x) \quad = \quad & \widehat{u}(t, x) \nabla_u f_{\nu_m} \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) - \frac{\partial}{\partial x} \left[ \widehat{u}(t, x) \nabla_v f_{\nu_m} \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \right] \\
+ \quad & \frac{\partial^2}{\partial x^2} \left[ \widehat{u}(t, x) \nabla_w f_{\nu_m} \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \right] + \widehat{u}(t, x) \nabla_u h_\theta \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \\
- \quad & \frac{\partial}{\partial x} \left[ \widehat{u}(t, x) \nabla_v h_\theta \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \right] + \frac{\partial^2}{\partial x^2} \left[ \widehat{u}(t, x) \nabla_w h_\theta \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \right] \\
\widehat{u}(t_{n+1}, x) \quad = \quad & \nabla_u \left\| u(t_{n+1}, x; \nu_m) - V(t_{n+1}, x; \nu_m) \right\|, \\
\nabla_\theta J(\theta) \quad = \quad & \int_{t_n}^{t_{n+1}} \int_\Omega \widehat{u}(t, x) \nabla_\theta h_\theta \big( u(t, x), u_x(t, x), u_{xx}(t, x) \big) \, dx dt,
\end{aligned} \tag{2.5}
$$

- Update the neural network parameters

$$\theta_{k+1} = \theta_k - \alpha_k \nabla_\theta J(\theta_k), \tag{2.6}$$

  and

- Repeat until a convergence criterion satisfied.

## 2.3 Adjoint equations in the multi-dimensional case

The adjoint equations for $d, q \geq 1$ follow from (2.2). Let $v_j$ and $w_j$ respectively be the arguments $\mathbf{u}_{x_j}$ and $\mathbf{u}_{x_j x_j}$ for the $\mathbf{f}$ and $\mathbf{h}_\theta$ functions. Then,

$$
\begin{aligned}
-\frac{\partial \widehat{\mathbf{u}}}{\partial t} = {} & \nabla_u \mathbf{f}_{\nu_m}(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}} \\
& - \sum_{j=1}^{q} \frac{\partial}{\partial x_j}\left[\nabla_{v_j} \mathbf{f}_{\nu_m}(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}}\right] \\
& + \sum_{j=1}^{q} \frac{\partial^2}{\partial x_j^2}\left[\nabla_{w_j} \mathbf{f}_{\nu_m}(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}}\right] \\
& + \nabla_u \mathbf{h}_\theta(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}} \\
& - \sum_{j=1}^{q} \frac{\partial}{\partial x_j}\left[\nabla_{v_j} \mathbf{h}_\theta(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}}\right] \\
& + \sum_{j=1}^{q} \frac{\partial^2}{\partial x_j^2}\left[\nabla_{w_j} \mathbf{h}_\theta(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}})^\top \widehat{\mathbf{u}}\right] \\
& + \sum_{n=1}^{N_t-1} \delta(t - t_n)\nabla_u \left\| \mathbf{u}(t_n, \mathbf{x}; \nu_m) - \mathbf{V}(t_n, \mathbf{x}; \nu_m)\right\|,
\end{aligned}
\tag{2.7}
$$

with the final condition

$$
\widehat{\mathbf{u}}(t_N, \mathbf{x}) = \nabla_u \left\| \mathbf{u}(t_N, \mathbf{x}; \nu_m) - \mathbf{V}(t_N, \mathbf{x}; \nu_m)\right\|,
\tag{2.8}
$$

and the gradient of the objective function (1.2) is then

$$
\nabla_\theta L(\theta) = \int_0^{t_N} \int_\Omega \widehat{\mathbf{u}}(t, \mathbf{x}) \cdot \nabla_\theta \mathbf{h}_\theta\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big) \, d\mathbf{x} dt.
\tag{2.9}
$$

## 2.4 Adjoint equations with a divergence-free constraint

When subject to a divergence-free constraint $\nabla_\mathbf{x} \cdot \mathbf{u} = 0$, such as in the case of the incompressible Navier–Stokes equations augmented with a neural-network model $\mathbf{h}_\theta$, the DPM framework (1.1) becomes

$$
\frac{\partial \mathbf{u}}{\partial t} = -\nabla_\mathbf{x} p + \mathbf{f}_\nu(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}) + \mathbf{h}_\theta(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}),
\tag{2.10}
$$

$$
\nabla_\mathbf{x} \cdot \mathbf{u} = 0.
\tag{2.11}
$$

The divergence-free condition is enforced via the pressure-like variable $p$. Thus, for $x \in \Omega$, the pressure satisfies

$$
\nabla^2 p = \nabla_\mathbf{x} \cdot \mathbf{f}_\nu(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}) + \nabla_\mathbf{x} \cdot \mathbf{h}_\theta(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}).
\tag{2.12}
$$

For computation, it is advantageous to time discretize (2.11) via an operator splitting method and then derive the adjoint equations for the time-discretized equations. This ensures that the adjoint equations are

fully compatible with the discretized PDE. A standard projection method is used for this [18]:

$$\mathbf{u}^*(t + \Delta t, \mathbf{x}) = \mathbf{u}(t, \mathbf{x}) + \Delta t \left[ f_\nu\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big) + h_\theta\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big) \right], \quad (2.13)$$

$$\nabla_\mathbf{x}^2 p(t, \mathbf{x}) = \frac{\nabla_\mathbf{x} \cdot \mathbf{u}^*(t + \Delta t, \mathbf{x})}{\Delta t}, \quad (2.14)$$

$$\mathbf{u}(t + \Delta t, \mathbf{x}) = \mathbf{u}^*(t + \Delta t, \mathbf{x}) - \Delta t \nabla_\mathbf{x} p(t, \mathbf{x}), \quad (2.15)$$

which yields the adjoint

$$\nabla_\mathbf{x}^2 \widehat{p}(t, \mathbf{x}) = -\frac{\nabla_\mathbf{x} \cdot \widehat{\mathbf{u}}(t + \Delta t, \mathbf{x})}{\Delta t},$$

$$\widehat{\mathbf{u}}^*(t + \Delta t, \mathbf{x}) = \widehat{\mathbf{u}}(t + \Delta t, \mathbf{x}) + \nabla_\mathbf{x} \widehat{p}(t, \mathbf{x}) \Delta t,$$

$$\begin{aligned}
\widehat{\mathbf{u}}(t, \mathbf{x}) = \widehat{\mathbf{u}}^*(t + \Delta t, \mathbf{x}) + \Bigg\{ & \nabla_u \mathbf{f}_{\nu_m}\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \\
& - \sum_{j=1}^q \frac{\partial}{\partial x_j}\left[ \nabla_{v_j} \mathbf{f}_{\nu_m}\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \right] \\
& + \sum_{j=1}^q \frac{\partial^2}{\partial x_j^2}\left[ \nabla_{w_j} \mathbf{f}_{\nu_m}\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \right] \\
& + \nabla_u \mathbf{h}_\theta\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \\
& - \sum_{j=1}^q \frac{\partial}{\partial x_j}\left[ \nabla_{v_j} \mathbf{h}_\theta\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \right] \\
& + \sum_{j=1}^q \frac{\partial^2}{\partial x_j^2}\left[ \nabla_{w_j} \mathbf{h}_\theta\big(\mathbf{u}, \mathbf{u_x}, \mathbf{u_{xx}}\big)^\top \widehat{\mathbf{u}}^* \right] \Bigg\} \Delta t \\
& + \sum_{n=1}^{N_t - 1} \mathbf{1}(t = t_n) \nabla_u \left\| \mathbf{u}(t_n, \mathbf{x}; \nu_m) - \mathbf{V}(t_n, \mathbf{x}; \nu_m) \right\|,
\end{aligned} \quad (2.16)$$

where $\mathbf{1}(t = t_n)$ is unity when $t = t_n$ and zero otherwise. The starting condition for the backward-in-time adjoint solve is

$$\widehat{\mathbf{u}}(t_N, \mathbf{x}) = \nabla_u \left\| \mathbf{u}(t_N, \mathbf{x}; \nu_m) - \mathbf{V}(t_N, \mathbf{x}; \nu_m) \right\|. \quad (2.17)$$

The gradient of the objective function (1.2) is

$$\nabla_\theta L(\theta) = \Delta t \sum_{t = 0, \Delta t, \ldots, t_{N_t - 1}} \int_\Omega \widehat{\mathbf{u}}^*(t + \Delta t, \mathbf{x}) \cdot \nabla_\theta \mathbf{h}_\theta\big(\mathbf{u}(t, \mathbf{x}), \mathbf{u_x}(t, \mathbf{x}), \mathbf{u_{xx}}(t, \mathbf{x})\big) \, d\mathbf{x}. \quad (2.18)$$

This adjoint equation is used in our demonstrations.

# 3    Sub-grid-scale closure for incompressible turbulence

The momentum equation, when filtered for LES, contains an unclosed sub-grid-scale (SGS) term that has been the subject of extensive modeling efforts. Accounting for it, along with any numerical discretization errors, is the goal of the deep learning model in this application. The Navier–Stokes equations are introduced in Section 3.1, followed by a discussion of filtering and the LES governing equations in Section 3.2.

## 3.1 Navier–Stokes governing equations

Incompressible fluid flow is governed by the Navier–Stokes equations, which comprise a momentum balance

$$\frac{\partial u_i}{\partial t} + \frac{\partial u_i u_j}{\partial x_j} = -\frac{1}{\rho}\frac{\partial p}{\partial x_i} + \frac{\mu}{\rho}\frac{\partial}{\partial x_j}\left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} - \frac{2}{3}\frac{\partial u_k}{\partial x_k}\delta_{ij}\right),\tag{3.1}$$

and an incompressibility constraint

$$\frac{\partial u_k}{\partial x_k} = 0,\tag{3.2}$$

where $\rho$ is the fluid density, $\mu$ is the dynamic viscosity, $u_i$ is the $i^{\text{th}}$ velocity component, $p$ is the pressure, and $\delta_{ij}$ is the Kronecker delta. Repeated indices imply summation. We take $\rho$ and $\mu$ to be constant.

For DNS, (3.1) and (3.2) are discretized with sufficient resolution that all turbulence scales are accurately resolved. This typically requires a mesh spacing comparable to the Kolmogorov scale to achieve mesh independence for typical statistical observables [19]. Evidence suggests that the DNS time step must be smaller than the Kolmogorov time scale in order to avoid spurious dissipation [20], though in general the time step size requirements depend on the application and numerical methods [19]. For many applications, which have Reynolds numbers that lead to excessively small Kolmogorov scales, DNS is prohibitively expensive even for relatively simple flows [19].

## 3.2 Large-eddy simulation

Using LES, computational expense is reduced by resolving only the largest turbulence scales. This involves, implicitly or explicitly, a spatial filtering operation applied to the Navier–Stokes equations, which leads to terms that depend on unrepresented scales. Thus, the filtered equations are said to be unclosed. Sub-grid-scale models are introduced to represent the effect of unrepresented scales in terms of the resolved scales, thereby closing the equations. In practice, the accuracy of LES calculations can be improved by better sub-grid-scale closures.

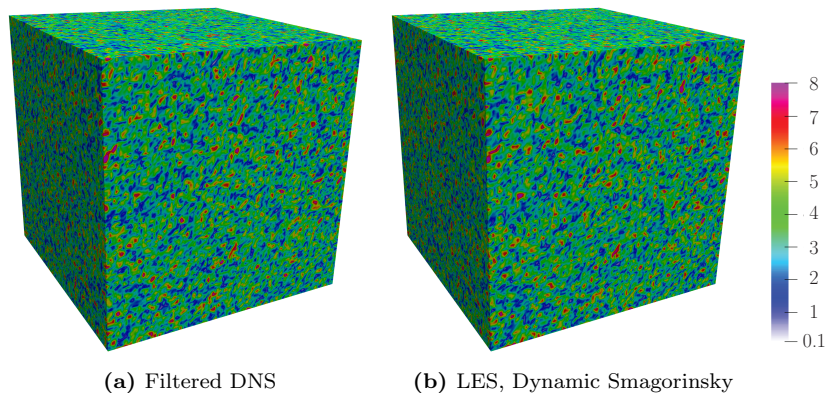A filtered quantity is denoted $\overline{\phi}$, which indicates

$$\overline{\phi}(\mathbf{x},t) = \int_{\Omega} G(\mathbf{r},\mathbf{x})\phi(\mathbf{x}-\mathbf{r},t)\,\mathrm{d}\mathbf{r},\tag{3.3}$$

where $\mathbf{x},\mathbf{r}\in\Omega\subset\mathbb{R}^p$ and the filter kernel is $G(\mathbf{r},\mathbf{x})$. Common choices are box, Gaussian, and spectral cutoff filters. We choose the common box filter for simplicity and because it replicates the common practice of implicit filtering, in which the filtering operation (3.3) is not explicitly performed so the LES mesh spacing itself truncates the small-scale features. A box filter on a uniform grid with spacing $\overline{\Delta}$ has $G=1$ within a $\overline{\Delta}^3$ cube centered at $\mathbf{x}$ [21]. Elsewhere, $G=0$.

Governing equations are obtained by applying (3.3) to (3.1) and (3.2), which yields

$$\frac{\partial \overline{u}_i}{\partial t} = -\frac{\partial \overline{u}_i \overline{u}_j}{\partial x_j} - \frac{1}{\rho}\frac{\partial \overline{p}}{\partial x_i} + \frac{\mu}{\rho}\frac{\partial}{\partial x_j}\left(\frac{\partial \overline{u}_i}{\partial x_j} + \frac{\partial \overline{u}_j}{\partial x_i} - \frac{2}{3}\frac{\partial \overline{u}_k}{\partial x_k}\delta_{ij}\right) + \frac{\partial \tau_{ij}^r}{\partial x_j}\tag{3.4}$$

$$\frac{\partial \overline{u}_k}{\partial x_k} = 0,\tag{3.5}$$

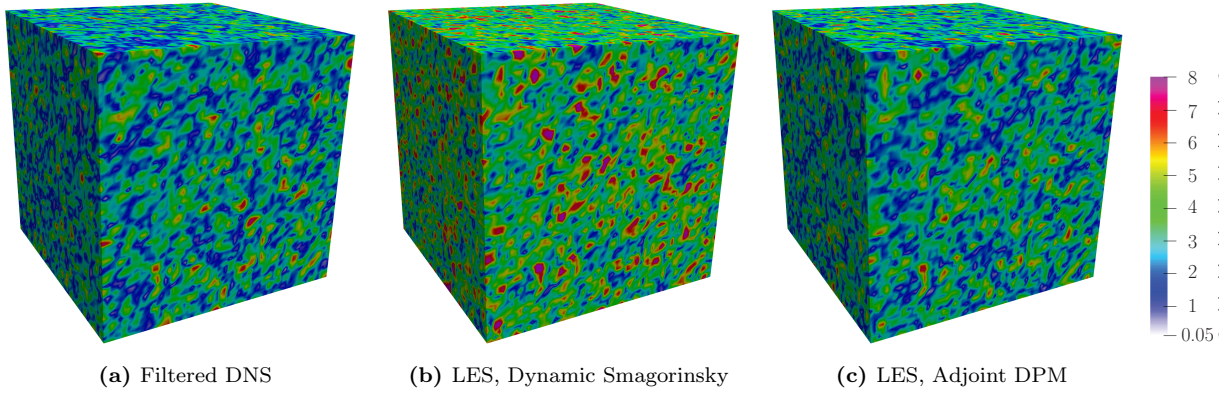**(a)** Filtered DNS        **(b)** LES, Dynamic Smagorinsky

**Figure 1:** Vorticity magnitude ($\times 10$) in decaying isotropic turbulence, normalized by the eddy-turnover time scale $t_{\ell,0} = k/\varepsilon$ (see Section 5.1), shown when energy is $36\,\%$ of the initial condition: (a) filtered $N = 2048^3$ DNS using $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$ for an effective grid resolution $N = 128^3$, and (b) $N = 128^3$ LES solution using the dynamic Smagorinsky [24] sub-grid-scale model.

where the SGS stress tensor $\tau_{ij}^r = \overline{u_i u_j} - \overline{u}_i \overline{u}_j$ requires closure.

The most common closures are based on gradient-diffusion with the eddy-viscosity coefficient determined as suggested by Smagorinsky [22, 23] or dynamically estimated based on resolved scales [24, 25]. Other options include scale-similarity models [26] and filter-scale Taylor expansions [21, 27, 28]. These and related models are extensively reviewed elsewhere [23, 29, 30]. Figure 1 compares a $\overline{\Delta}/\Delta x = 16$ filtered $N = 2048^3$ DNS field on a $N = 128^3$ grid with a corresponding $N = 128^3$ LES using the dynamic Smagorinsky model. For this Reynolds number and relatively small LES filter width (*i.e.*, a wide range of resolved scales), the dynamic Smagorinsky model yields a qualitatively-correct LES solution in Figure 1. This is also supported quantitatively based on turbulence statistics such as we consider subsequently.

While these common sub-grid-scale models are accurate for sufficiently small $\overline{\Delta}$ (and, if different, $\Delta x_{\mathrm{LES}}$), resolving still fewer scales would increase computational efficiency, if it can be done accurately. Even a modest improvement can significantly reduce the overall cost of a three-dimensional, time-dependent simulation. Our objective is to develop turbulence closures using limited data available from such costly high-fidelity DNS data. Figure 2(a) shows a coarser $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$ filtered DNS solution on a $N = 64^3$ grid, obtained from a $N = 1024^3$ DNS. Figure 2(b) shows the analogous $N = 64^3$ LES using the same dynamic Smagorinsky model. For this coarser filter (*i.e.*, a narrower range of resolved scales), obvious visual differences suggest that the model produces a qualitatively-incorrect solution. Conversely, the adjoint-trained DPM model produces the qualitatively-correct solution shown in Figure 2(c), even for this coarse mesh. Section 5 includes quantitative comparisons between the DPM and Smagorinsky and dynamic Smagorinsky models.

**(a)** Filtered DNS        **(b)** LES, Dynamic Smagorinsky        **(c)** LES, Adjoint DPM

**Figure 2:** Vorticity magnitude ($\times 10$) in decaying isotropic turbulence, normalized by the eddy-turnover time scale $t_{\ell,0} = k/\varepsilon$ (see Section 5.1), shown when energy is 9.2 % of the initial condition: (a) filtered $N = 1024^3$ DNS using $\overline{\Delta}/\Delta x = 16$ for an effective grid resolution $N = 64^3$, (b) $N = 64^3$ LES solution using the dynamic Smagorinsky [24] sub-grid-scale model, and (c) $N = 64^3$ LES solution using the adjoint-trained DPM model developed in this work.

# 4    Discretization errors and learning corrections

In addition to closure error, LES calculations suffer from discretization error. These discretization errors can themselves be viewed as unclosed terms, which arise from discretization of the governing equations on a coarse grid. Sub-grid-scale closures formulated as PDE terms, as well as any *a priori*-trained machine learning turbulence models (see Section 5.5), do not account for them. Yet in standard LES practice they can be comparable to $\boldsymbol{\tau}_r$ modeling errors, depending on the filter width and spatial discretization scheme [1–3]. This section discusses and quantifies discretization errors for LES and shows how the DPM can learn to compensate for discretization errors. This trained model will of course be linked to both the mesh density and the numerical schemes, and training dataset will need to be sufficiently large to allow for generalizations. However, resolutions will not vary widely for efficient LES applications since it is most efficient to run with the coarsest mesh possible, which will aid generalization.

The filtered Navier–Stokes solution $\overline{u}_i$ satisfies

$$\frac{\partial \overline{u}_i}{\partial t} = \widehat{\mathcal{A}}_i^{\Delta}(\overline{\mathbf{u}}, \overline{p}) + \underbrace{\left( \mathcal{A}_i(\overline{\mathbf{u}}, \overline{p}) - \widehat{\mathcal{A}}_i^{\Delta}(\overline{\mathbf{u}}, \overline{p}) \right)}_{\text{Discretization error}} + \underbrace{\nabla \cdot \boldsymbol{\tau}^r}_{\text{Closure mismatch}}$$

$$0 = \mathcal{D}^{\Delta}(\overline{\mathbf{u}}) + \underbrace{\left( \nabla \cdot \overline{\mathbf{u}} - \mathcal{D}^{\Delta}(\overline{\mathbf{u}}) \right)}_{\text{Discretization error}}, \tag{4.1}$$

where $\mathcal{A}_i(\overline{\mathbf{u}}, \overline{p}) = -\frac{\partial \overline{u}_i \overline{u}_j}{\partial x_j} - \frac{1}{\rho}\frac{\partial \overline{p}}{\partial x_i} + \frac{\mu}{\rho}\frac{\partial}{\partial x_j}\left( \frac{\partial \overline{u}_i}{\partial x_j} + \frac{\partial \overline{u}_j}{\partial x_i} - \frac{2}{3}\frac{\partial \overline{u}_k}{\partial x_k}\delta_{ij} \right)$ are terms in the continuous flow equation, $\widehat{\mathcal{A}}_i^{\Delta}(\overline{\mathbf{u}}, \overline{p})$ is the discrete approximation to $\mathcal{A}_i(\overline{\mathbf{u}}, \overline{p})$ for the selected mesh and numerical schemes, and $\mathcal{D}^{\Delta}(\overline{\mathbf{u}})$ is the corresponding discrete approximation to $\nabla \cdot \overline{\mathbf{u}}$. When trained on the coarse LES mesh with the filtered (downsampled) DNS data, the DPM $\mathbf{h}_\theta$ will learn the closure mismatch and discretization errors in (4.1).

To assess the discretization errors in the LES, several quantities calculated with the finite-difference

| Filtering | $\dfrac{\overline{\Delta}}{\Delta x_{\mathrm{DNS}}}$ | $\dfrac{\Delta}{\Delta x_{\mathrm{DNS}}}$ | $\dfrac{\langle|\delta\overline{u}_1|\rangle}{\langle|\nabla\overline{\mathbf{u}}|\rangle_{\mathrm{DNS}}}$ | $\max(\mathcal{D}^\Delta(\overline{\mathbf{u}}))_{\mathrm{DNS}}$ | $\max(\mathcal{D}^\Delta(\overline{\mathbf{u}}))_{\mathrm{LES}}$ | $\dfrac{\max(\mathcal{D}^\Delta(\overline{\mathbf{u}}))_{\mathrm{LES}}}{\langle|\nabla\overline{\mathbf{u}}|\rangle_{\mathrm{DNS}}}$ | $\dfrac{\langle|\mathcal{D}^\Delta(\overline{\mathbf{u}})|\rangle_{\mathrm{LES}}}{\langle|\nabla\overline{\mathbf{u}}|\rangle_{\mathrm{DNS}}}$ |
|---|---|---|---|---|---|---|---|
|  | 8 | 8 | 0.601 | $2.545 \times 10^{-8}$ | $2.257 \times 10^{5}$ | 7.893 | 0.827 |
| Implicit | 16 | 16 | 0.854 | $1.566 \times 10^{-8}$ | $1.216 \times 10^{5}$ | 6.288 | 1.081 |
|  | 32 | 32 | 1.076 | $0.596 \times 10^{-8}$ | $0.579 \times 10^{5}$ | 6.651 | 1.208 |
|  | 32 | 16 | 0.654 | $6.399 \times 10^{-9}$ | $4.734 \times 10^{4}$ | 5.435 | 0.895 |
| Explicit | 32 | 8 | 0.325 | $7.195 \times 10^{-9}$ | $2.855 \times 10^{4}$ | 3.277 | 0.476 |
|  | 32 | 4 | 0.141 | $6.930 \times 10^{-9}$ | $1.325 \times 10^{4}$ | 1.522 | 0.211 |

**Table 1:** Error in the finite-difference approximation to the filtered velocity gradient $\delta\overline{u}_1 = \overline{u}_{1,x,\mathrm{LES}} - \overline{u}_{1,x,\mathrm{DNS}}$, evaluated on LES grids of varying resolution $n_{x,\mathrm{DNS}}/n_{x,\mathrm{LES}}$, relative to the average magnitude of the filtered velocity gradient evaluated on the DNS grid. The maximum of the discrete filtered-velocity divergence is also shown evaluated on the DNS and LES grids. Averages $\langle\cdot\rangle$ are over the full simulation domain at fixed time.

schemes used in the DNS and LES are listed in Table 1. The starting-point DNS velocity field $u_i$ is computed on a $N = 2048^3$ mesh. Filtered velocity components $\overline{u}_i$ are obtained with a box filter of width $\overline{\Delta}$. With the DNS providing the trusted solution, we take as reference the $\overline{u}_i$ derivative in the $x$-direction at $(t, x + \Delta x/2)$ to be a dense-mesh finite difference

$$\overline{u}_{i,x,\mathrm{DNS}}(t, x + \Delta_{\mathrm{LES}}/2) = \frac{\overline{u}_i(t, x + \Delta x_{\mathrm{DNS}}) - \overline{u}_i(t,x)}{\Delta x_{\mathrm{DNS}}}, \tag{4.2}$$
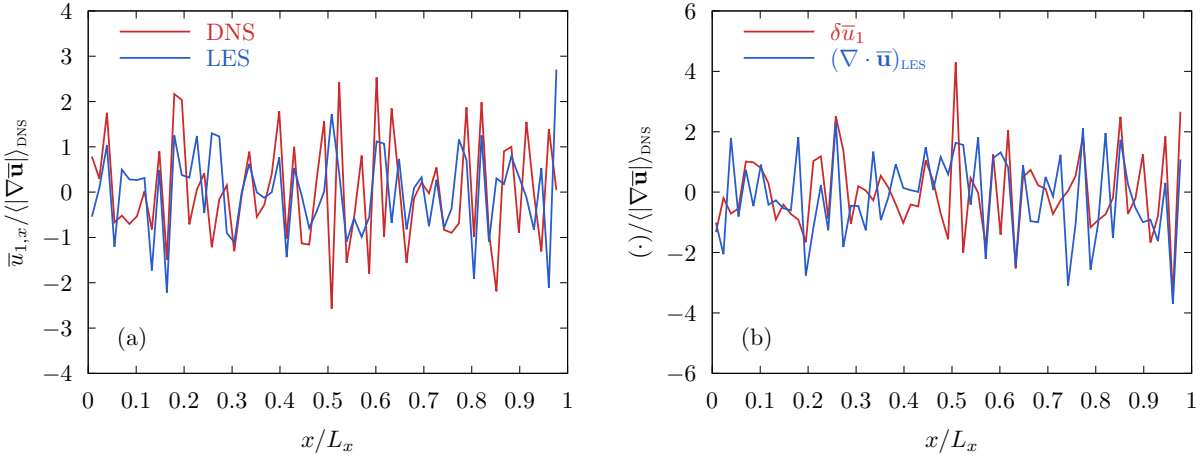
which is compatible with how the derivatives were calculated in the DNS. On the coarse LES mesh, the corresponding filtered velocity gradient is evaluated as

$$\overline{u}_{i,x,\mathrm{LES}}(t, x + \Delta_{\mathrm{LES}}/2) = \frac{\overline{u}_i(t, x + \Delta_{\mathrm{LES}}) - \overline{u}_i(t,x)}{\Delta_{\mathrm{LES}}},$$

where $\Delta_{\mathrm{LES}}$ is the coarse-mesh resolution. Table 1 shows that the difference between the DNS-mesh derivative and the LES-mesh derivative $\delta\overline{u}_i = \overline{u}_{i,x,\mathrm{LES}} - \overline{u}_{i,x,\mathrm{DNS}}$ is comparable to the average velocity gradient magnitude evaluated on the DNS mesh $\langle|\nabla\overline{\mathbf{u}}|\rangle_{\mathrm{DNS}}$. Figure 3(a) compares them for a representative segment of the domain. In Table 1, the error $\delta\overline{u}_1$ increases with the mesh size, becoming greater than the average velocity gradient magnitude for $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 32$. With such errors, even if the unclosed term $\nabla \cdot \boldsymbol{\tau}^r$ is modeled exactly, the LES calculation would be inaccurate.

It is noteworthy that the discrete divergence-free constraint $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$ is not satisfied by the filtered DNS solution when evaluated on the coarse mesh. Table 1 shows that the filtered velocity satisfies the discrete divergence-free condition on the DNS grid within a small factor of the Poisson equation solver tolerance $(10^{-9})$, while the filtered velocity on the coarse LES grid does not. $\mathcal{D}^\Delta(\overline{\mathbf{u}})$ is comparable to the average velocity gradient, as is shown in Figure 3(b). Similarly, the average divergence residual is comparable to the average DNS velocity gradient magnitude.

Numerical errors such as these are widely recognized. Nonetheless, LES solvers typically enforce the condition $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$. This accurately reflects the commutative property of the linear divergence opera-

**Figure 3:** Typical errors versus position: (a) finite-difference approximations to the filtered velocity gradient $\overline{u}_{1,x}$ evaluated on DNS ($N = 2048^3$) and LES ($N = 64^3$, $\overline{\Delta}/\Delta x = 32$) grids, normalized by the average DNS velocity gradient magnitude; and (b) $\delta \overline{u}_1 = \overline{u}_{1,x,\text{LES}} - \overline{u}_{1,x,\text{DNS}}$ and the normalized velocity divergence on the LES grid.

tor (3.2) and filter (3.3), though it does not reflect the discretization error of the discrete divergence operator evaluated on the coarse LES grid. Enforcing $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$ to hold exactly for the coarse LES grid is a choice. (Perot [31] discusses this more generally for incompressible flow simulation: the common practice of enforcing a discrete-exact version of incompressibility rather than a discrete-exact version of momentum is also a choice.) Another perspective is that the coarse-mesh (LES) divergence should not be zero in the LES. If there were a hypothetical sub-grid-scale field (and not necessarily a unique one [32]) that exactly reflected the realistic SGS dynamics, the coarse-mesh sampling of it would not yield $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$, and would deviate further from this condition the coarser the mesh. In an extension in Section 5.6, the DPM is shown to also be able to exploit this flexibility to provide an accurate model that reproduces key features of the turbulence by learning, rather than explicitly representing, the physics of the incompressibility constraint.

# 5 Numerical results

The DNS data used for training and out-of-sample testing are introduced in Sections 5.1 and 5.2, the neural network architecture is described in Section 5.3, and model comparisons are presented in Section 5.4. Their behavior is compared to a more straightforward *a priori* training regimen in Section 5.5. The extension mentioned at the end of the previous section, which learns rather than enforces strict adherence to the divergence-free $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$ constraint, is discussed in Section 5.6. Finally, computational cost is quantified and discussed in Section 5.7.

| $N$ | Case | $\mu/\mu_0$ | $u_{\mathrm{rms},0}$ | $t_{\ell,0} = k/\varepsilon$ | $t_{\eta,0} = (\mu\varepsilon/\rho)^{1/2}$ | $\varepsilon_0$ | Train | Test | $Re_{t,0}$ |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 0.50 | 101 | $1.27 \times 10^{-1}$ | $2.03 \times 10^{-3}$ | $1.21 \times 10^{5}$ | • | | |
| | 2 | 0.75 | 152 | $8.49 \times 10^{-2}$ | $1.35 \times 10^{-3}$ | $4.09 \times 10^{5}$ | | • | |
| $1024^3$ | 3 | 1.00 | 203 | $6.38 \times 10^{-2}$ | $1.01 \times 10^{-3}$ | $9.69 \times 10^{5}$ | • | | 1749 |
| | 4 | 1.25 | 255 | $5.09 \times 10^{-2}$ | $8.12 \times 10^{-4}$ | $1.89 \times 10^{6}$ | | • | |
| | 5 | 1.50 | 304 | $4.24 \times 10^{-2}$ | $6.79 \times 10^{-4}$ | $3.27 \times 10^{6}$ | | • | |
| | 6 | 2.00 | 406 | $3.18 \times 10^{-2}$ | $5.09 \times 10^{-4}$ | $7.74 \times 10^{6}$ | • | | |

**Table 2:** The DNS datasets that are used for training and testing the models are described. All cases have an initial integral scale $L_{ii} = 1.00$. The domain for all cases is a periodic cube with sides of length $L = 66.5$, the initial Kolmogorov length scale for all cases is $\eta = 3.18 \times 10^{-2}$, and the initial Reynolds number based on the Taylor microscale is $Re_\lambda \approx 162$. Cases used for training and testing ensembles are indicated.

## 5.1 DNS data

The decaying isotropic turbulence DNS were initialized with a standard spectrum [33] at an initial Reynolds number $Re_{t,0} = \rho u_{\mathrm{rms}}\ell/\mu = 1749$, where $u_{\mathrm{rms}} = \langle u_i u_i \rangle^{1/2}$ is the domain-averaged root-mean-squared velocity, $\ell = k^{3/2}/\varepsilon$ is the pseudo-integral scale of the initial spectrum, $k = \langle u_i u_i \rangle/2$ is the turbulence kinetic energy (TKE), and $\varepsilon$ is the TKE dissipation rate. The reference density and viscosity are $\rho = 1$ and $\mu_0 = 1$.

Though isotropic turbulence is nominally Reynolds-number isoparametric, we chose to train the DPM in a manner more generally applicable. We require that the model, if possible, learn the Reynolds number scaling. Though obvious for isotropic turbulence, this is emblematic of requiring the model to learn corresponding parameters in more complex scenarios. For training and testing cases, we adjust the dimensional time scales of the turbulence such that the model inputs vary by factors of 10 for nominally the same instantaneous Reynolds number in the decay. However, this is a side concern: the main point is that the DPM learns the instantaneous structure of the turbulence sufficiently to close the LES governing equations (including numerical error terms) better than established models.

Table 2 lists initial parameters for the six $N = 1024^3$ DNS datasets that are used to train and test models. Three each are used for training and out-of-sample testing. The decay time scales are adjusted via $\mu$ and $u_{\mathrm{rms},0}$ for fixed initial Reynolds number $Re_{t,0}$. Data were produced for $\mu/\mu_0 = \{0.5, 0.75, 1.0, 1.25, 1.5, 2.0\}$, each for two random phasings of the initial conditions, for a total of 12 simulations. As shown in Table 2, the initial TKE dissipation rate varies by a factor of 64.

The DNS were first allowed to decay for $t \approx 0.05 t_{\ell,0}$. After this period, the velocity fields were stored every $0.01 t_{\ell,0}$. These were box filtered with $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$. Initial conditions on the coarse LES mesh were obtained by restricting data from the fine DNS mesh, resulting in the irrecoverable information loss inherent to LES [27, 32]. For simplicity, we only consider cases in which the filter width and LES grid resolution are equivalent, which is common practice though it lacks significant theoretical basis [34, 35].

| Field | Mesh $\Delta x$ | Description |
|---|---|---|
| $\mathbf{u}$ | DNS | trusted DNS data |
| $\overline{\mathbf{u}}$ | DNS | filtered DNS data (3.3) |
| $\overline{\mathbf{U}}$ | LES | sub-sampled filtered DNS data (every $k$-th point for $k = \frac{\Delta x_{\mathrm{DNS}}}{\Delta x_{\mathrm{LES}}}$) |
| $\overline{\mathbf{w}}$ | LES | divergence-free projected filtered DNS training target |
| $\widetilde{\mathbf{u}}$ | LES | numerical solution of LES equations |

**Table 3:** Notation for data variables used in simulations and model training.

These DNS have a large number of integral scales per domain length, $L/L_{ii} \approx 66$, which is helpful for statistical sampling. For comparison, recent simulations aimed at maximizing the Reynolds number have $L/L_{ii} \approx 5$ [36], which is nearly correlated across the periodic domain.

The DNS datasets were produced using the *NGA* code [37,38], which uses a fractional-step method [39]. For the DNS cases listed in Table 2, the initial CFL number was 0.4, and the time step $\Delta t$ was fixed throughout the simulations. Space is discretized using second-order central differences on a staggered mesh [40].

Adjoint-based *a posteriori* model training was performed using a new implementation that matches the *NGA* finite-difference discretization but is Python-native and interfaces with the *PyTorch* [41] machine learning and automatic differentiation library. It was co-verified with *NGA* and validated against standard accepted turbulence results. While the adjoint could in theory be entirely computed using existing automatic differentiation software, this would be prohibitively expensive in practice. Instead, we solve the adjoint equations (2.16) on the staggered mesh using the Python-native solver. Only gradients of the neural network ($\nabla_\theta \mathbf{h}_\theta$) are evaluated using automatic differentiation.

## 5.2  Pre-processing of training data

It is important to clearly define the training problem, including how the training data is processed, which specifies the unclosed terms targeted by $\mathbf{h}_\theta$. The following develops the specifics for our demonstration, building upon the general discussion of discretization errors in Section 4.

Variable labels associated with different stages are summarized in Table 3. We assume that the DNS fields $\mathbf{u}$ are sufficiently trusted that we disregard any errors they might entail. The filtered DNS solution $\overline{\mathbf{u}}$ notionally satisfies the continuous filtered flow equations,

$$\frac{\partial \overline{u}_i}{\partial t} = \mathcal{A}_i(\overline{\mathbf{u}}, \overline{p}) + \nabla \cdot \boldsymbol{\tau}^r$$
$$0 = \nabla \cdot \overline{\mathbf{u}}, \tag{5.1}$$

where we use the compact $\mathcal{A}$ notation of (4.1) for the terms in the flow equations. This stage only includes the explicit sub-grid-scale stress closure $\nabla \cdot \boldsymbol{\tau}^r$, and it remains represented numerically on the fine $\Delta x_{\mathrm{DNS}}$ mesh. To be used in training on the $\Delta x_{\mathrm{LES}}$ coarse mesh, it is downsampled: $\overline{\mathbf{U}}_{ijk} = \overline{\mathbf{u}}(t, i\Delta_{\mathrm{LES}}, j\Delta_{\mathrm{LES}}, k\Delta_{\mathrm{LES}})$,

14

where for simplicity $\overline{\Delta} = \Delta x_{\text{LES}}$. The downsampled filtered DNS satisfies

$$
\begin{aligned}
\frac{\partial \overline{U}_{ijk}}{\partial t} &= \widehat{\mathcal{A}}^{\Delta}(\overline{\mathbf{U}}, \overline{p})_{ijk} + \left[\mathcal{A}(\overline{\mathbf{u}}, \overline{p})(i\overline{\Delta}, j\overline{\Delta}, k\overline{\Delta}) - \widehat{\mathcal{A}}^{\Delta}(\overline{\mathbf{U}}, \overline{p})_{ijk}\right] + \nabla \cdot \boldsymbol{\tau}^{r}(i\overline{\Delta}, j\overline{\Delta}, k\overline{\Delta}) \\
0 &= \mathcal{D}^{\Delta}(\overline{\mathbf{U}})_{ijk} + \left[\nabla \cdot \overline{\mathbf{u}}(i\Delta, j\Delta, k\Delta) - \mathcal{D}^{\Delta}(\overline{\mathbf{U}})_{ijk}\right],
\end{aligned}
\tag{5.2}
$$

where on this coarse mesh $\overline{\mathbf{U}}$ is not divergence free due to the residual $\nabla \cdot \overline{\mathbf{u}}(i\Delta, j\Delta, k\Delta) - \mathcal{D}^{\Delta}(\overline{\mathbf{U}})_{i,j,k}$. Appendix A includes details on the discrete forms underlying the projection used to ensure that the training data $\overline{\mathbf{w}}$ satisfies $\mathcal{D}^{\Delta}(\overline{\mathbf{w}}) = 0$ in order to be compatible with standard LES practice (though this is revisited in Section 5.6). The divergence-free result of this projection can be expressed $\mathbf{w} = \overline{\mathbf{U}} + \mathcal{H}^{\Delta}\overline{\mathbf{U}}$, where $\mathcal{H}^{\Delta}$ is a (constant) linear operator with the property that $\|\mathcal{H}^{\Delta}\overline{\mathbf{U}}_t\| \to 0$ as $\Delta \to 0$. This final projection stage adds still further unclosed terms to the momentum equation, now governing $\overline{\mathbf{w}}$ on the coarse LES mesh:

$$
\begin{aligned}
\frac{\partial \mathbf{w}_{ijk}}{\partial t} = \widehat{\mathcal{A}}^{\Delta}(\mathbf{w}, \overline{p})_{ijk} + \underbrace{\nabla \cdot \boldsymbol{\tau}^{r}}_{\text{SGS closure}} + \underbrace{\left[\widehat{\mathcal{A}}^{\Delta}(\overline{\mathbf{U}}, \overline{p})_{ijk} - \widehat{\mathcal{A}}^{\Delta}(\mathbf{w}, \overline{p})_{ijk}\right]}_{\text{Projection effect on } \mathcal{A}^{\Delta} \text{ evaluation}} \\
+ \underbrace{\frac{\partial \mathcal{H}^{\Delta}\overline{\mathbf{U}}}{\partial t}}_{\text{Change in the projection}} + \underbrace{\left[\mathcal{A}(\overline{\mathbf{u}}, \overline{p})(i\Delta, j\Delta, k\Delta) - \widehat{\mathcal{A}}^{\Delta}(\overline{\mathbf{U}}, \overline{p})_{ijk}\right]}_{\text{Finite-difference error}}
\end{aligned}
\tag{5.3}
$$

with the discrete divergence-free constraint

$$
\mathcal{D}^{\Delta}(\mathbf{w})_{ijk} = 0.
\tag{5.4}
$$

The last four terms of (5.3) are, in a sense, unclosed, representing the mismatch targeted by the neural network $\mathbf{h}_{\theta}$ by optimizing $\theta$ in

$$
\frac{\partial \widetilde{\mathbf{u}}}{\partial t} = \widehat{\mathcal{A}}^{\Delta}(\widetilde{\mathbf{u}}, \overline{p}) + \mathbf{h}_{\theta},
\tag{5.5}
$$

where $\widetilde{\mathbf{u}}$ is a numerical LES solution, and optimization seeks to minimize

$$
L(\theta) = \sum_{n=1}^{N_t} \sum_{i,j,k=1}^{N} \|\widetilde{\mathbf{u}} - \overline{\mathbf{w}}\|.
\tag{5.6}
$$

## 5.3   Model architecture and hyperparameters

The deep neural network $F_{\theta}(z)$ has the following architecture:

$$
\begin{aligned}
H^1 &= \sigma(W^1 z + b^1) \\
H^2 &= \sigma(W^2 H^1 + b^2) \\
H^3 &= G^1 \odot H^2 \qquad \text{with} \quad G^1 = \sigma(W^5 z + b^5) \\
H^4 &= \sigma(W^3 \odot H^3 + b^3) \\
H^5 &= G^2 \odot H^4 \qquad \text{with} \quad G^2 = \sigma(W^6 z + b^6) \\
F_{\theta}(z) &= W^4 H^5 + b^4
\end{aligned}
\tag{5.7}
$$

where $\sigma$ is a tanh() element-wise nonlinearity, $\odot$ denotes element-wise multiplication, and the parameters are $\theta = \{W^1, W^2, W^3, W^4, W^5, W^6, b^1, b^2, b^3, b^4, b^5, b^6\}$. The final machine learning model $\mathbf{h}_\theta$, which is used in the DPM (1.1), applies a series of derivative operations on $F_\theta(z)$. The input $z$ to $F_\theta(z)$ at an LES grid point $(i, j, k)$ includes the velocity components and their first and unmixed second derivatives at $(i, j, k)$ as well as the 6 closest neighboring grid points. This selection does not strictly enforce Galilean invariance, but this was not found to be a challenge. The issue of Galilean invariance should be considered further, especially in regard to further extrapolation from the training data that we consider here. Enforcing this (or many other) invariances would be straightforward. Each layer includes $N_{\mathrm{H}} = 200$ hidden units, for a total of approximately $245,000$ parameters, all initialized using a standard Xavier initialization [42]. Inputs to the neural network are normalized by the same set of constants for all cases. Although large, this network proves both effective (Section 5.4) and efficient (Section 5.7).
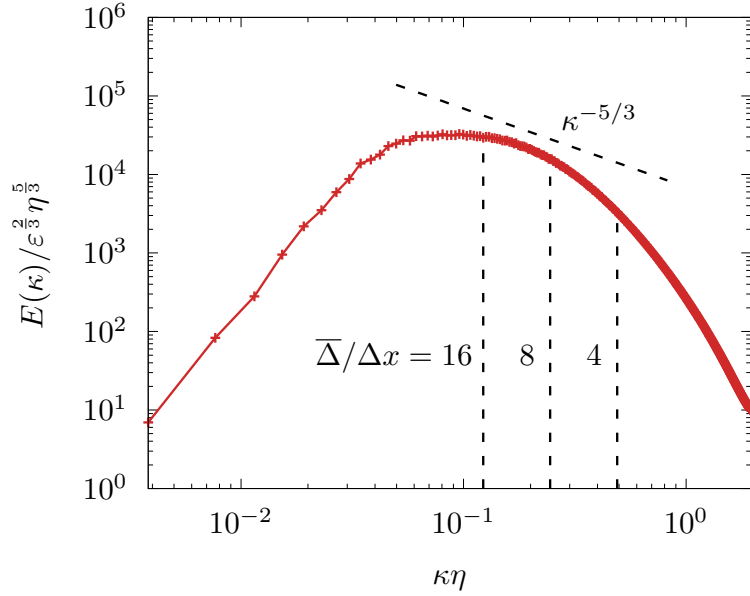
Training is distributed across multiple compute nodes, with each working with a randomly-selected velocity field from the training data in Table 2. We advance the LES solutions over five LES time steps, which is the equivalent of 50 DNS time steps, then solve the adjoint over the equivalent reverse-time interval. Parameter updates use the RMSprop algorithm with a standard decaying learning rate magnitude schedule. Training is accelerated by distributing computations across multiple GPU nodes.

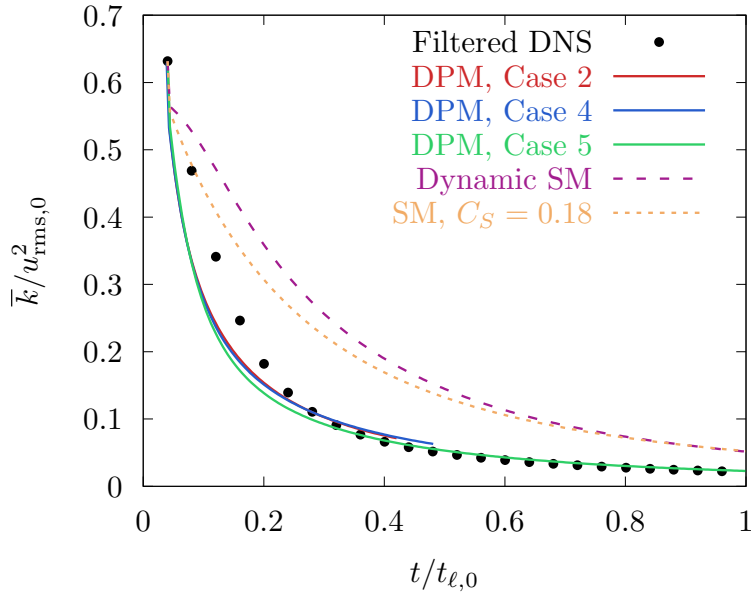## 5.4    Out-of-sample model comparison

To quantify model needs, and ultimately performance, in terms of representing turbulence scales, velocity fields are Fourier transformed in all three directions to provide standard wavenumber-magnitude $\kappa \equiv (\kappa_i \kappa_i)^{1/2}$ spectra $E(\kappa)$. Figure 4 shows this for the full range of $\kappa$ for the unfiltered $N = 1024^3$ DNS data with respect to the nominal filter cutoffs for different cases. The most challenging $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$ filter nominal cutoff scale is close to the most energetic turbulence scales. The $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 4$ and $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 8$ cutoffs, more typical for LES, include more represented scales and should therefore be more accurate. However, the higher resolution they afford is costly: halving the filter increases the operation count by approximately a factor of 16. Models that perform well for large $\overline{\Delta}$ are therefore attractive.

The simplest quantity of interest for isotropic turbulence is the decaying kinetic energy of the resolved scales, $\overline{k}(t) \equiv \frac{1}{2}\langle \overline{u_i u_i} - \overline{u}_i \overline{u}_i \rangle$. Figure 5 compares $\overline{k}(t)$ with the filtered DNS. For all the three out-of-sample test cases listed in Table 2, the DPM outperforms the widely-used Smagorinsky model [22, 23], whether the coefficient is determined dynamically [24, 25] or fixed at $C_S = 0.18$ [21]. These all have the same Reynolds number and so follow the same decay-rate profile, though this was learned only based on the instantaneous fields at different initial dissipation rates, so the DPM indeed learns the correct rescaling. The eddy viscosity-type LES models perform better (not shown) with smaller $\overline{\Delta}$ (and thus significantly higher cost) but perform poorly for $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$.
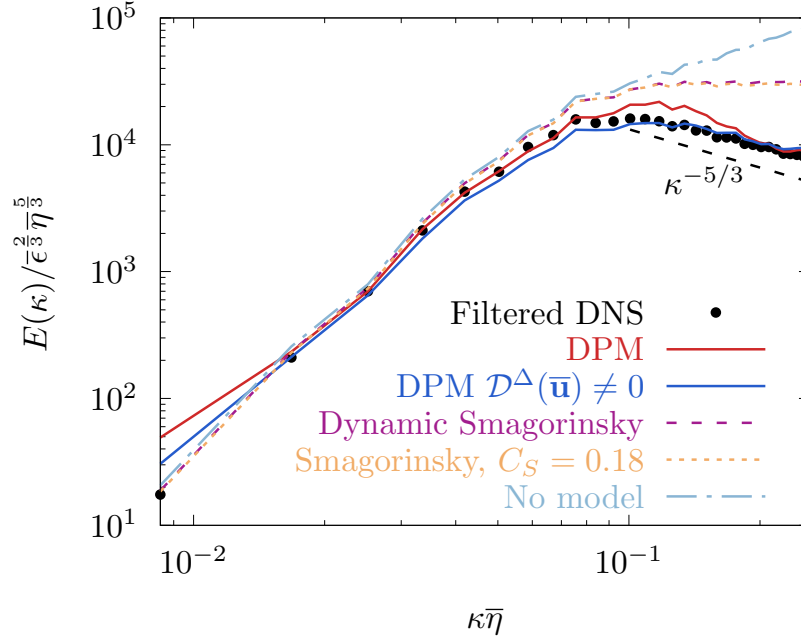
The more detailed comparison in Figure 6 shows that the adjoint-trained DPM reproduces the spectrum

**Figure 4:** Energy spectrum of the unfiltered $N = 1024^3$ isotropic turbulence DNS. The nominal cutoff wavenumbers for spectrally-sharp filters of size $\overline{\Delta}/\Delta x = 4$, 8, and 16 are indicated by vertical dashed lines. A reference Kolmogorov $\kappa^{-5/3}$ inertial range is also indicated.



**Figure 5:** Decay of $\overline{k}(t)$ for out-of-sample test case 2 ($\mu/\mu_0 = 0.75$), case 4 ($\mu/\mu_0 = 1.25$), and case 5 ($\mu/\mu_0 = 1.5$) from Table 2, for LES filter size $\overline{\Delta}/\Delta x = 16$. The adjoint PDE-trained DPM correctly learns the normalized evolution despite not having been trained on these dissipation rates.

17

**Figure 6:** Resolved energy spectra for filter size $\overline{\Delta}/\Delta x = 16$ for an out-of-sample ($\mu/\mu_0 = 1.5$) test case at time instant $t/t_{\ell,0} = 0.36$. DPM ($\mathcal{D}^\Delta(\overline{\mathbf{u}}) \neq 0$) shows the non-divergence-free model result of Section 5.6.

remarkably well, qualitatively better than the eddy-viscosity-based models. For this large $\overline{\Delta}/\Delta x_{\mathrm{DNS}} = 16$, the constant-coefficient and dynamic Smagorinsky energy spectra are nearly identical, though it is confirmed that, as expected, the dynamic model outperforms the constant-coefficient model for smaller $\overline{\Delta}$. The overpredicted high-wavenumber energy is consistent with the underpredicted energy decay rates in Figure 5. Also shown is a no-model LES, which is still further from the filtered DNS spectrum.
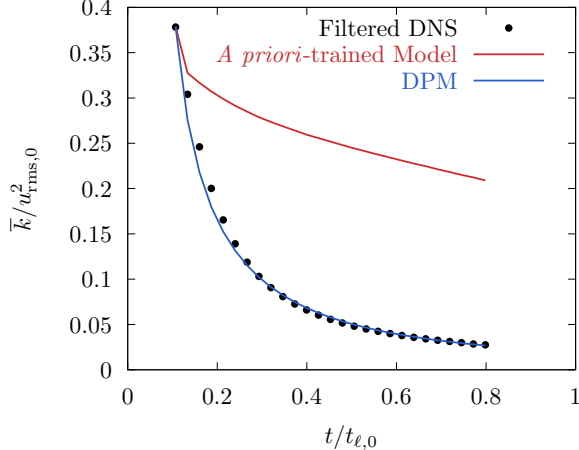
## 5.5  Comparison with *a priori* training

For comparison, we analyze the common practice of *a priori* training, which seeks sub-grid-scale corrections without regard to how they would couple with the LES PDEs. To do this, $\boldsymbol{\tau}^r$ from the filtered DNS data is used to directly train $\mathbf{h}_\theta$ by minimizing

$$J(\theta) = \|\mathbf{h}_\theta(\overline{\mathbf{u}}) - \nabla \cdot \boldsymbol{\tau}^r\|, \tag{5.8}$$

where $\overline{\mathbf{u}}$ and $\boldsymbol{\tau}^r$ are from the DNS without additional processing and do not change during the optimization. This de-coupled form also does not include errors due to the LES discretization. Once trained, using the same input and neural-network architecture as the DPM models introduced in Section 5.3, $\mathbf{h}_\theta$ closes the LES equation.

An advantage of *a priori* training is that it is easy to implement and only requires standard machine learning tools to minimize $J$. However, the DPM is expected to outperform *a priori* training because of

18

**Figure 7:** Performance of an *a priori*-trained deep learning model compared to an adjoint-trained DPM and filtered DNS data for filter size $\overline{\Delta}/\Delta x = 16$ and case $\mu/\mu_0 = 1.0$ from Table 2.

the mathematical inconsistency of the latter: the *a priori* approach interchanges optimization (training) with a nonlinear operation (the LES equations), which is expected to degrade subsequent predictions. For our demonstration of *a priori* training, $\mathbf{h}_\theta$ is trained on the trusted filtered DNS data $\overline{\mathbf{u}}$, whereas during *a posteriori* LES it only receives input from the solution of the discretized LES equations. This provides robustness to nonlinear accumulation of finite-difference error during the simulation. In summary, the neural network parameters $\theta$ that minimize (5.8) for the DNS data do not necessarily minimize the objective function (1.2), and so do not necessarily yield a consistent and accurate LES model.

Figure 7 shows the poor predictive performance of the *a priori*-trained deep learning closure model for LES. The DPM, trained with the adjoint method, performs substantially better than the *a priori*-trained model.

In summary, there are three main advantages sought with our adjoint-based approach. Foremost, it avoids the inconsistency of *a priori* training: optimization does not commute with a nonlinear function, which introduces error. Second, it accounts for numerical errors on the coarse LES grid. Including the full numerical discretization of the PDE in the adjoint-based optimization allows $\mathbf{h}_\theta$ to respond to the unavoidable LES discretization error (discussed in Section 4). This is linked to the particular resolution, but is better than neglecting this important challenge inherent in any LES. Finally, *a priori* training as designed in this example requires a full training-target description of the to-be-modeled $\nabla \cdot \boldsymbol{\tau}_r$. This is readily available in DNS, though not generally. For example, *a priori* training cannot so readily be used to estimate a deep learning PDE model from limited experimental data, which is almost always relatively sparse. It is mostly readily applied when the unclosed terms to be modeled are directly available as in (5.8). If the DPM is estimated from $M_s$ high-fidelity numerical simulation datasets and $M_e$ experimental datasets

19

with $N_e$ measurement points, our original objective function (1.2) becomes

$$L(\theta) = \underbrace{\sum_{m=1,\ldots,M_s} \sum_{n=1}^{N_t} \int_\Omega \|\mathbf{u}(t_n, \mathbf{x}; \nu_m) - \mathbf{V}(t_n, \mathbf{x}; \nu_m)\| \ d\mathbf{x}}_{\text{Simulations}}$$

$$+ \underbrace{\sum_{m=M_s+1,\ldots,M_s+M_e} \sum_{n=1}^{N_t} \sum_{j=1}^{N_e} \|\mathbf{u}(t_n, \mathbf{x}_j; \nu_m) - \mathbf{V}(t_n, \mathbf{x}_j; \nu_m)\|}_{\text{Experiments}}.$$

(5.9)

where $\mathbf{x}_j$ and $t_n$ in the second sum do not need to be full-field quantities. The DPM facilitates evaluations for whatever limited data points are available.

## 5.6 Relaxation of the divergence-free constraint

As discussed in Section 4, the divergence-free constraint in the LES equations (3.5) is not satisfied by filtered DNS data on the LES mesh. The standard approach, followed thus far, is to enforce incompressibility on the filtered target data as discussed in Section 5.2. This ensured consistency between the discretized LES and the filtered target data: the target data is a solution of the LES equations on the discrete grid.

We consider an alternate approach in this section. Instead of projecting the target data to be divergence-free, we relax the strict requirement that the LES solution be discretely divergence-free. This significantly reduces computational cost by eliminating the elliptic solve in the pressure-projection part of the algorithm. Of course, the filtered DNS data still satisfies the divergence-free constraint (3.2) on the DNS grid, so the model learns the divergence-free constraint from the evolution of the full-resolution DNS data rather than relying on its discrete enforcement on the LES grid. As such, this is also a demonstration of the DPM's ability to learn physics that is omitted from the PDE.

A spectrum for this DPM $(\mathcal{D}^\Delta(\bar{\mathbf{u}}) \neq 0)$ case is also plotted in Figure 6. Performance of this new model is comparable to the original DPM, better in some sense, particularly in matching low-wavenumber energy. However, the differences between the two deep learning models are small, particularly when compared to the improvement of both over the established LES models. The possibility of extending the non-divergence-free approach to more complex flows, in particular pressure-driven flows and reacting flows, is not automatically precluded. This example was included primarily as a further illustration of the capacity of the DPM framework to replace physics with learned models and closures. We do recognize that less explicitly represented physics is likely to diminish the capacity for extrapolation, and that some form of validation would be needed as for any model reduction, machine learning or otherwise. Our expectation, in this case, is some lost capacity to extrapolate to flows that include an important non-zero mean pressure gradient.

## 5.7 Computational cost

LES is used to reduce cost, so we provide a brief analysis of this. A general concern with deep learning models is high expense due to the many tensor operations to evaluate a large network. For the cases reported thus
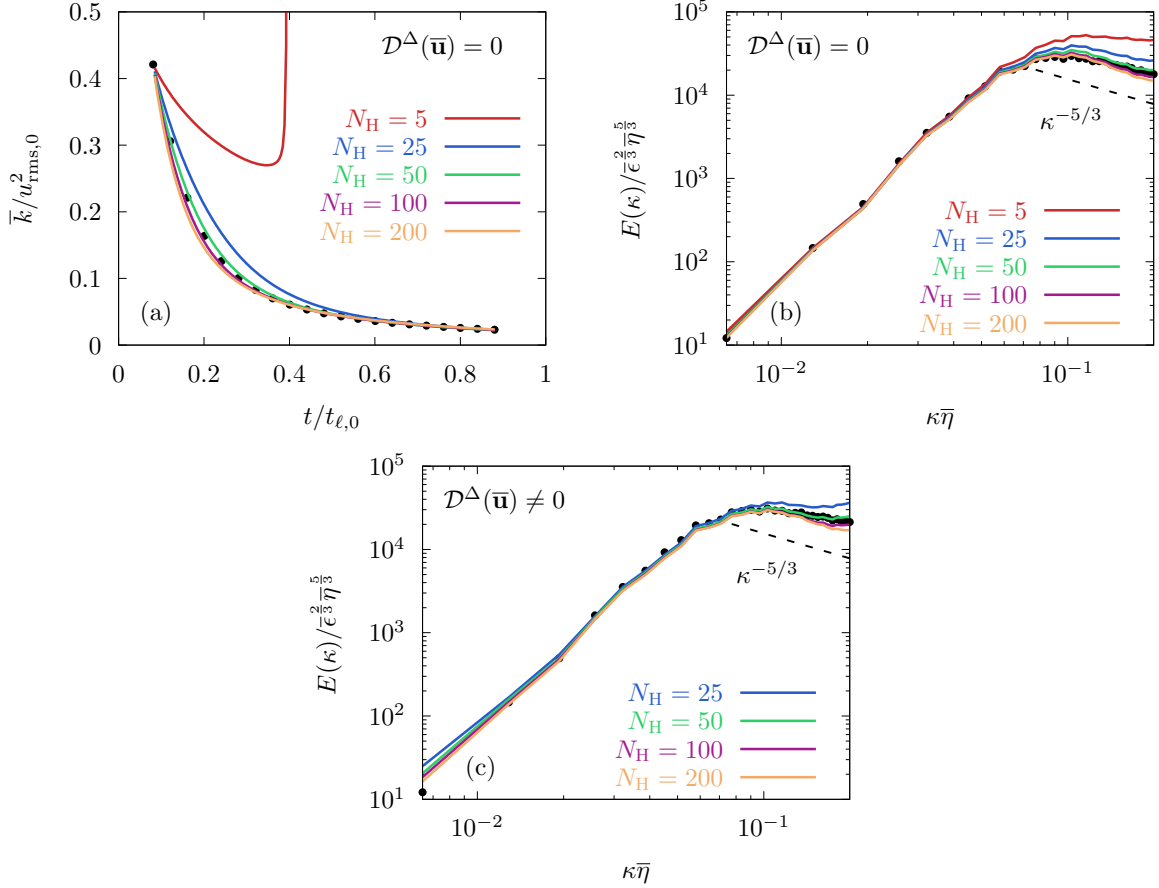
| $N_{\mathrm{H}}$ | $d_\theta$ | Time, CPU (s) | Time, GPU (s) | Speed-up |
|---|---|---|---|---|
| 5 | 4,278 | 1.83 | $4.98 \times 10^{-2}$ | 36.7 |
| 25 | 22,318 | 2.45 | $7.00 \times 10^{-2}$ | 35.0 |
| 50 | 47,118 | 3.35 | $7.75 \times 10^{-2}$ | 43.2 |
| 100 | 104,218 | 5.59 | $8.87 \times 10^{-2}$ | 63.0 |
| 200 | 248,418 | 10.49 | $1.53 \times 10^{-1}$ | 68.5 |

**Table 4:** Average time to evaluate the deep neural network (5.7) on a $N = 64^3$ mesh. The network size is given in terms of the number of hidden units per layer $N_{\mathrm{H}}$ and the total number of parameters $d_\theta$. The average wall-clock time is reported for a single AMD 6276 Interlagos CPU core and a single NVIDIA K20X GPU. The listed speed-up is the ratio of single-core CPU time to GPU time.

far, our model has 248,418 parameters, which includes a correspondingly large number of operations at each grid point each time step compared to the other models. Even in this case, the DPM, though based on a large network, is not uncompetitive. We address this empirically for the current implementation in two parts. First, we quantify the cost of evaluating the neural network—and the associated solution accuracy—as a function of the number of parameters. Second, for reference, we compare the measured full-simulation cost of our Python-native LES to that of the optimized *NGA* (Fortran) solver.

Table 4 lists the computational time required to evaluate the deep neural network (5.7) at all points in a $N = 64^3$ mesh for different network sizes. This accounts for almost $50\,\%$ of the total LES cost (see Table 5). Although the operation count is high, the cost of this evaluation benefits tremendously from the highly-optimized (*PyTorch*-based) implementation. Table 4 provides a measure of this by comparing the network-evaluation cost for a single CPU core (AMD 6276 Interlagos) and a single GPU accelerator (NVIDIA K20X). The efficiency of the GPU implementation of the neural network evaluation enables bringing seemingly large trained networks into current practice for prediction. Although cost increases with network size, the GPU-accelerated speedup also improves for the relevant range considered here.

The large $N_{\mathrm{H}} = 200$ network we have considered was selected to focus our study on the DPM rather than the neural-network design. Although it is able to leverage GPU acceleration to be practical at this large scale, smaller networks of similar design perform nearly as well, as shown in Figure 8. This compares DPM $\overline{k}$ decay and spectra for $N_{\mathrm{H}} \in \{5, 25, 50, 100, 200\}$ for both the divergence free $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$ and not strictly divergence-free $\mathcal{D}^\Delta(\overline{\mathbf{u}}) \neq 0$ variants. Even $N_{\mathrm{H}} = 25$ show better agreement with spectra than the Smagorinsky models (Figure 6). More hidden units (larger $N_{\mathrm{H}}$) increase accuracy, with the associated cost. We do not attempt to prove stability of the DPM models, which would likely be challenging. No stabilizing limiters were used in the present implementation, though if needed they could provide a practical means of ensuring stability, such as are typically used to stabilize dynamic Smagorinsky models [27]. For both $\mathcal{D}^\cdot(\mathbf{u}) = 0$ and $\mathcal{D}^\cdot(\mathbf{u}) \neq 0$ variants, deeper neural networks reduce numerical instability, with $N_{\mathrm{H}} \geq 50$

**Figure 8:** Dependence on number of hidden units per layer $N_{\mathrm{H}}$ (lines) compared to filtered DNS (points): (a) decay of $\overline{k}(t)$ and (b) spectra for divergence-free ($\mathcal{D}^{\Delta}(\mathbf{u}) = 0$) models, and (c) spectra for $\mathcal{D}^{\Delta}(\mathbf{u}) \neq 0$ models. The filtered DNS fields in (a) and (b) are projected onto divergence-free manifolds. Spectra in (a) and (b) are at $t/t_{\ell,0} = 0.2$. In (c), the $N_{\mathrm{H}} = 5$ case has become unstable at this time, and the $N_{\mathrm{H}} = 25$ is showing signs of high-wavenumber divergence. All models are trained on $\mu/\mu_0 \in \{0.5, 1.0, 2.0\}$ and are tested on $\mu/\mu_0 = 1.5$.

| Model | Fortran (s/step) | Python (s/step) |
|---|---|---|
| No-model LES | 7.57 | 14.11 |
| Smagorinsky, $C_S = 0.18$ | 7.73 | 14.24 |
| Smagorinsky, Dynamic | 8.44 | — |
| DPM ($N_H = 200$) | — | 27.86 |

**Table 5:** Single-core CPU performance of Fortran and Python LES solvers for decaying isotropic turbulence on a $N = 64^3$ mesh. All computations were performed on a single AMD 6276 Interlagos CPU core. The dynamic Smagorinsky model was not implemented in Python, and the DPM was not implemented in Fortran.

| Model | Fortran (s/step) $N = 128^3$, 16-core | Python (s/step) $N = 64^3$, K20X GPU |
|---|---|---|
| Smagorinsky, $C_S = 0.18$ | 6.47 | — |
| Smagorinsky, Dynamic | 6.73 | — |
| DPM | — | 1.60 |
| DPM ($\mathcal{D}^\Delta(\overline{\mathbf{u}}) \neq 0$) | — | 0.31 |

**Table 6:** Single-node performance of Fortran and Python LES solvers for comparable solution accuracy. The Smagorinsky models are evaluated on a $N = 128^3$ mesh using the Fortran solver, and the deep learning models are evaluated on a $N = 64^3$ mesh using the Python solver. The Fortran solver used 16 CPU cores on a Cray XK7 compute node (AMD 6276 Interlagos CPUs; 313 GF peak performance), while the Python solver used an NVIDIA K20X GPU (1.31 TF peak performance).

required for long-time ($t \geq 10^{-3}$) stability for $\mathcal{D}^\cdot(\mathbf{u}) = 0$ and $N_H \geq 100$ without this constraint ($\mathcal{D}^\cdot(\mathbf{u}) \neq 0$).

Of course, significant inefficiency of the flow solver within the overall DPM model would also hide the network evaluation cost. To assess this, we compare the Python-based, *PyTorch*-coupled LES solver used against the optimized Fortran *NGA* solver. These data, computed by inserting system timers around the time-marching loop and advancing the solution 50 steps, are shown in Table 5. The Python solver is less than a factor of two slower.

Finally, we provide an estimate of application computational cost for comparable solution accuracies. We use the *NGA* Smagorinsky model on a $N = 128^3$ mesh to produce a comparably accurate solution to the DPM on a $N = 64^3$ mesh. Both solvers use the available computational resources on a Cray XK7 compute node containing 16 CPU cores and one NVIDIA K20X GPU. The Fortran solver is faster in single-core performance but, like many legacy solvers, is not readily GPU-accelerated. The network evaluation in the Python-native solver is easily GPU-accelerated using the *PyTorch* library. Performance is compared in Table 6. The Python-native DPM solution, even in divergence-free form, is approximately one-quarter the cost. Invoking the ability of the DPM to learn the divergence-free condition reduces the DPM-solution cost to approximately one-twentieth of the legacy-solution cost.

# 6 Conclusion

Results demonstrate the promise of deep learning PDE models (DPMs) applied to problems with unresolved physics, demonstrating it on a case with limited high-fidelity numerical data for training. As formulated, however, the estimation method can also incorporate experimental data. It was demonstrated for learning the unclosed terms in the filtered Navier–Stokes equations (the LES sub-grid-scale stress tensor). The DPM is trained on filtered DNS data and its out-of-sample accuracy is studied.

The DPM recovered an accurate representation of the resolved turbulence. This result is an important first step in demonstrating the ability of the model to learn unrepresented (or unknown) physics. It outperforms established eddy-viscosity models, including the dynamic Smagorinsky model, in terms of reproducing the resolved kinetic energy decay rate and the resolved kinetic energy spectrum observed in the exact filtered DNS data. It is expected that additional optimizations and refinement of the design of the training regimen would further increase performance. Results also suggest that the adjoint PDE-based training of the DPM is able to correct for numerical discretization errors, which are widely known to affect the performance of LES calculations on coarse meshes. These discretization errors were interpreted as additional unclosed terms in the discrete governing equations on the coarse LES grid.

As a generalization toward learning additional physics, and an intriguing demonstration of a possible direction for turbulence simulation, a formulation was also considered that does not exactly enforce a discrete divergence-free constraint in LES. Rather, the physics of the pressure variable on the coarse LES grid is directly learned from the data. An advantage is reduced computational cost since it does not require a Poisson solver.

Extending to more complex Navier–Stokes turbulence flows is an obvious direction. An important challenge here will be to obtain and use training sets with sufficiently rich variability to enable extrapolation to new configurations.

The DPM design is also much broader than Navier–Stokes turbulence. We anticipate that it might be particularly suited to still more intricate sub-grid-scale modeling problems, such as in mixing, combustion, or additional physical mechanisms coupled with turbulence. The formulation is such that experimental data could be incorporated within the formulation even if the physical description is unclear. It is designed to leverage known and resolvable physics to the highest degree possible, not replace it.

# A    Additional Discretization Details

## A.1    Divergence-free projection of DNS fields

Since the filtered DNS solution does not satisfy the discrete divergence-free condition $\mathcal{D}^\Delta(\overline{\mathbf{u}}) = 0$ on the coarse LES mesh, as a pre-processing step we make the appropriate $\ell^2$ projection at each time $t$:

$$\mathbf{w}_t = \arg\min_{\mathbf{w}\in\mathcal{C}} \frac{1}{2} \left\| \overline{\mathbf{U}}_t - \mathbf{w} \right\|_2^2, \tag{A.1}$$

where $\overline{\mathbf{U}}_t = \left\{ \overline{\mathbf{u}}(t, i\Delta, j\Delta, k\Delta) \right\}_{i,j,k=1}^N$ is the filtered DNS data on the coarse $N^3$ grid at time $t$ and where

$$\mathcal{C} = \left\{ \mathbf{w} \in \mathbb{R}^{3\times N\times N\times N} : \quad \frac{w_{1,i+1,j,k} - w_{1,i,j,k}}{\Delta} + \frac{w_{2,i,j+1,k} - w_{2,i,j,k}}{\Delta} + \frac{w_{3,i,j,k+1} - w_{3,i,j,k}}{\Delta} = 0 \right.$$

$$\left. \text{and} \quad w_{m,N,j,k} = w_{m,1,j,k}, w_{m,i,N,k} = w_{m,i,1,k}, w_{m,i,j,N} = w_{m,i,j,1} \right\}. \tag{A.2}$$

Gradients of $\mathbf{w}$ in (A.2) have been derived using the same second-order, staggered-mesh central difference scheme as the flow solver [40].

That is, we estimate the $\mathbf{w}_t$ nearest to $\overline{\mathbf{U}}_t$ such that $\mathbf{w}_t$ satisfies the discretized divergence-free condition on the coarse LES grid. Using the method of Lagrange multipliers, the minimizer $\mathbf{w}$ of the constrained optimization problem (A.1) can be shown to satisfy

$$
\begin{aligned}
w_{1,i,j,k} &= \overline{U}_{1,i,j,k} - \frac{\lambda_{i+1,j,k} - \lambda_{i,j,k}}{\Delta}, \\
w_{2,i,j,k} &= \overline{U}_{2,i,j,k} - \frac{\lambda_{i,j+1,k} - \lambda_{i,j,k}}{\Delta}, \\
w_{3,i,j,k} &= \overline{U}_{3,i,j,k} - \frac{\lambda_{i,j,k+1} - \lambda_{i,j,k}}{\Delta},
\end{aligned} \tag{A.3}
$$

where the pressure-like Lagrange multiplier $\lambda$ is a solution to the discrete Poisson equation

$$
\begin{aligned}
\frac{\overline{U}_{1,i,j,k} - \overline{U}_{1,i-1,j,k}}{\Delta} &+ \frac{\overline{U}_{2,i,j,k} - \overline{U}_{2,i,j-1,k}}{\Delta} + \frac{\overline{U}_{3,i,j,k} - \overline{U}_{3,i,j,k-1}}{\Delta} \\
&= \frac{\lambda_{i,j+1,k} - 2\lambda_{i,j,k} + \lambda_{i,j-1,k}}{\Delta^2} + \frac{\lambda_{i+1,j,k} - 2\lambda_{i,j,k} + \lambda_{i-1,j,k}}{\Delta^2} + \frac{\lambda_{i,j,k+1} - 2\lambda_{i,j,k} + \lambda_{i,j,k-1}}{\Delta^2},
\end{aligned} \tag{A.4}
$$

with periodic boundary conditions. This yields $\mathbf{w}_t$ used for the target $\mathbf{V}$ in (1.2).

# References

[1] S. Ghosal, An analysis of numerical errors in large-eddy simulations of turbulence, J. Comp. Phys. 125 (1) (1996) 187–206.

[2] A. G. Kravchenko, P. Moin, On the effect of numerical errors in large eddy simulations of turbulent flows, J. Comp. Phys. 131 (2) (1997) 310–322.

[3] F. K. Chow, P. Moin, A further study of numerical errors in large-eddy simulations, J. Comp. Phys. 184 (2) (2003) 366–380.

[4] J. Sirignano, K. Spiliopoulos, DGM: A deep learning algorithm for solving partial differential equations, J. Comp. Phys. 375 (2018) 1339–1364.

[5] B. A. Sen, S. Menon, Linear eddy mixing based tabulation and artificial neural networks for large eddy simulations of turbulent flames, Comb. and Flame 157 (2010) 62–74.

[6] B. A. Sen, E. R. Hawkes, S. Menon, Large eddy simulation of extinction and reignition with artificial neural networks based chemical kinetics, Comb. and Flame 157 (2010) 566–578.

[7] P. P. Popov, D. A. Buchta, M. J. Anderson, L. Massa, J. Capecelatro, D. J. Bodony, J. B. Freund, Machine learning-assisted early ignition prediction in a complex flow, Comb. and Flame 206 (2019) 451–466.

[8] M. Raissi, H. Babaee, P. Givi, Deep learning of PDF turbulence closure, Bulletin of the American Physical Society (2018).

[9] M. Raissi, Z. Wang, M. S. Triantafyllou, E. G. Karniadakis, Deep learning of vortex-induced vibrations, J. Fluid Mech. 861 (2019) 119–137.

[10] M. Raissi, P. Perdikaris, E. G. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, J. Fluid Mech. 378 (2019) 686–707.

[11] J. Ling, R. Jones, J. Templeton, Machine learning strategies for systems with invariance properties, J. Comp. Phys. 318 (2016) 22–35.

[12] J. Ling, A. Kurzawski, J. Templeton, Reynolds averaged turbulence modelling using deep neural networks with embedded invariance, J. Fluid Mech. 807 (2016) 155–166.

[13] J. Berg, K. Nystrom, Neural network augmented inverse problems for PDEs, arXiv preprint arXiv:1712.09685.

[14] Z. Wang, K. Luo, D. Li, J. Tan, J. Fan, Investigations of data-driven closure for subgrid-scale stress in large-eddy simulation, Phys. Fluids 30 (2018) 125101.

[15] S. L. Brunton, B. R. Noack, P. Koumoutsakos, Machine learning for fluid mechanics, Ann. Rev. Fluid Mech 52 (2020) 1.

[16] M. P. Brenner, J. D. Eldredge, J. B. Freund, A perspective on machine learning for advancing fluid mechanics, Phys. Rev. Fluids 4 (2019) 100501.

[17] J. R. Holland, J. D. Baeder, K. Duraisamy, Towards integrated field inversion and machine learning with embedded neural networks for RANS modeling, In: AIAA SciTech Forum (2019) 1884.

[18] A. J. Chorin, Numerical solution of the Navier–Stokes equations, Math. Comput. 22 (1968) 745–762.

[19] P. Moin, K. Mahesh, Direct numerical simulation: a tool in turbulence research, Annu. Rev. Fluid Mech. 30 (1) (1998) 539–578.

[20] H. Choi, P. Moin, Effects of the computational time step on numerical solutions for turbulent flow, J. Comp. Phys. 113 (1994) 1–4.

[21] R. A. Clark, J. H. Ferziger, W. C. Reynolds, Evaluation of subgrid-scale models using a fully simulated turbulent flow, J. Fluid Mech. 91 (1) (1979) 1–16.

[22] J. Smagorinsky, General circulation experiments with the primitive equations I. The basic experiment, Mon. Weather Rev. 91 (3) (1963) 99–164.

[23] R. S. Rogallo, P. Moin, Numerical simulation of turbulent flows, Annu. Rev. Fluid Mech. 16 (1984) 99–137.

[24] M. Germano, U. Piomelli, P. Moin, W. H. Cabot, A dynamic subgrid-scale eddy viscosity model, Phys. Fluids 3 (7) (1991) 1760–1765.

[25] D. K. Lilly, A proposed modification of the Germano subgrid-scale closure method, Phys. Fluids 4 (1992) 633–635.

[26] J. Bardina, J. H. Ferziger, W. C. Reynolds, Improved subgrid-scale models for large-eddy simulation, AIAA Paper 80–1357 (1980).

[27] B. Vreman, B. Geurts, H. Kuerten, Large-eddy simulation of the turbulent mixing layer, J. Fluid Mech. 339 (1997) 357–390.

[28] A. Vollant, G. Balarac, C. Corre, A dynamic regularized gradient model of the subgrid-scale stress tensor for large-eddy simulation, Phys. Fluids 28 (2) (2016) 025114.

[29] M. Lesieur, O. Métais, New trends in large-eddy simulations of turbulence, Annu. Rev. Fluid Mech. 28 (1996) 45–82.

[30] C. Meneveau, J. Katz, Scale-invariance and turbulence models for large-eddy simulation, Annu. Rev. Fluid Mech. 32 (2000) 1–32.

[31] J. B. Perot, Discrete conservation properties of unstructured mesh schemes, Annu. Rev. Fluid Mech. 43 (1) (2011) 299–318.

[32] J. A. Langford, R. D. Moser, Optimal LES formulations for isotropic turbulence, J. Fluid Mech. 398 (1999) 321–346.

[33] T. Passot, A. Pouquet, Numerical simulation of compressible homogeneous flows in the turbulent regime, J. Fluid Mech. 181 (1987) 441–466.

[34] T. S. Lund, The use of explicit filters in large eddy simulation, Comput. Math. with Appl. 46 (4) (2003) 603–616.

[35] S. T. Bose, P. Moin, D. You, Grid-independent large-eddy simulation using explicit filtering, Phys. Fluids 22 (10) (2010) 105103.

[36] T. Ishihara, K. Morishita, M. Yokokawa, A. Uno, Y. Kaneda, Energy spectrum in high-resolution direct numerical simulations of turbulence, Phys. Rev. Fluids 1 (2016) 082403.

[37] O. Desjardins, G. Blanquart, G. Balarac, H. Pitsch, High order conservative finite difference scheme for variable density low Mach number turbulent flows, J. Comp. Phys. 227 (15) (2008) 7125–7159.

[38] J. F. MacArt, M. E. Mueller, Semi-implicit iterative methods for low Mach number turbulent reacting flows: Operator splitting versus approximate factorization, J. Comp. Phys. 326 (2016) 569–595.

[39] J. Kim, P. Moin, Application of a fractional-step method to incompressible Navier–Stokes equations, J. Comp. Phys. 59 (2) (1985) 308–323.

[40] F. H. Harlow, J. E. Welch, Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface, Phys. Fluids 8 (12) (1965) 2182–2189.

[41] B. Steiner, Z. DeVito, S. Chintala, S. Gross, A. Paszke, F. Massa, A. Lerer, G. Chanan, Z. Lin, E. Yang, A. Desmaison, A. Tejani, A. Kopf, J. Bradbury, L. Antiga, M. Raison, N. Gimelshein, S. Chilamkurthy, T. Killeen, L. Fang, J. Bai, Pytorch: An imperative style, high-performance deep learning library, NeurIPS 2019.

[42] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, Proceedings of the thirteenth international conference on artificial intelligence and statistics (2010) 249–256.