



Formulating turbulence closures using sparse regression with embedded form invariance

S. Beetham * and J. Capecelatro *Department of Mechanical Engineering, University of Michigan, Ann Arbor, Michigan 48105, USA*

(Received 28 March 2020; accepted 28 July 2020; published 28 August 2020)

A data-driven framework for formulation of closures of the Reynolds-average Navier-Stokes (RANS) equations is presented. In recent years, the scientific community has turned to machine learning techniques to translate data into improved RANS closures. While the body of work in this area has primarily leveraged neural networks (NNs), we alternately leverage a sparse regression framework. This methodology has two important properties: (1) The resultant model is in a closed, algebraic form, allowing for direct physical inferences to be drawn and naive integration into existing computational fluid dynamics solvers, and (2) Galilean invariance can be guaranteed by thoughtful tailoring of the feature space. Our approach is demonstrated for two classes of flows: homogeneous free shear turbulence and turbulent flow over a wavy wall. The model learned based upon the wavy wall configuration is then validated against flow over a backward-facing step. This work demonstrates similar performance to that of modern NNs but with the added benefits of interpretability, increased ease of use and dissemination, and robustness to sparse and noisy training data sets.

DOI: [10.1103/PhysRevFluids.5.084611](https://doi.org/10.1103/PhysRevFluids.5.084611)

I. INTRODUCTION

Simulation frameworks based on the Reynolds-averaged Navier-Stokes (RANS) equations [1–4] have been the most widely-used tool in industrial and large-scale applications of turbulent flows for the last several decades [5] and will remain a central tool for guiding design decisions well into the coming decades [6,7]. This is primarily driven by the wide range of length scales and timescales associated with turbulent flows of interest. Because of this, direct numerical simulations (DNS) that fully resolve all relevant scales are prohibitively costly. Instead, the RANS equations solve for mean flow quantities that are then used to assess global flow features of interest. A principal challenge associated with RANS is accurate modeling of the unresolved terms, which are denoted “unclosed” because they are not completely specified in terms of the unknowns (e.g., mean velocity, pressure, etc.).

With the rise of computational power and the accessibility of large, highly resolved data sets, the community has turned to machine learning techniques in recent years to distill this wealth of information into improved RANS models. As a consequence of the interest and prevalence of the use of machine learning in the turbulence modeling community, several thoughtful and thorough reviews have been published and the authors refer the interested reader to several of these works, including Brenner *et al.* [8], Duraisamy *et al.* [9], Holland *et al.* [10], and Duraisamy *et al.* [11].

Numerous studies in recent years have approached the RANS closure problem by leveraging a neural network (NN)-based framework. Ling *et al.* [12] used an invariant tensor basis integrated into a NN to model the Reynolds stress anisotropy tensor for turbulent duct flow as well as

*snverner@umich.edu

flow over a wavy wall. Galilean invariance, a critical model property, was ensured by nature of the invariant tensor basis as the inputs to the NN and demonstrated excellent agreement with DNS data as compared to traditional (linear and quadratic eddy viscosity) models. Following this work, many others have implemented similar strategies and employed a similar basis technique for ensuring invariance. Of these works, many have used flow through a periodically constricted channel or backward facing step as challenging tests of new modeling methodologies as these flows exhibit massive separation. This is notoriously difficult to accurately capture with traditional RANS closures [13]. A large body of works using NNs as the data-driven methodology to formulate closure models have demonstrated promising success [12,14–18].

Despite the demonstration of improved model performance, models based upon NNs have an important drawback. Due to the nature of the algorithm at the heart of NNs, the resultant model acts as a “black box” and cannot be expressed in a compact, algebraic form. This compromises interpretability, introduces difficulty in disseminating the learned model with end users and industries, and increases computational cost of the model in the context of a RANS solver (as compared with traditional algebraic closures). Further, a large number of NN approaches attempt to augment or correct existing models. However, this approach breaks down for more complex turbulent flows, such as disperse two-phase flows [19–23] or turbulent combustion [24,25], in which the fundamental assumption of an energy cascade breaks down due to production at the smallest scales. In these cases, existing closures adopted from single-phase flows are not appropriate, which precludes an augmentation modeling approach. For these reasons, the present study proposes an alternate method that allows for the development of physics-based, compact algebraic closures, thus affording interpretability, transportability, and efficiency.

Several studies have taken alternate approaches to NNs using symbolic methods in order to arrive at closed form, algebraic models. Gene expression programming [26–30] and random forest regression [31,32] have become increasingly popular methodologies. The early success of these works serves as motivation for the present work in which we present a methodology based upon sparse regression as an alternative to NNs for developing new RANS closures, with emphasis on the following key benefits:

Interpretability: Sparse regression produces an algebraic model with a limited number of terms, resulting in improved interpretability of underlying physics and better prediction of model behavior and stability outside the scope of training.

Galilean invariance: By careful construction of the feature space and structuring of the optimization cost functional, Galilean invariance of the resultant model is ensured.

Efficiency: Models developed using sparse regression are built using physics-based, functional terms and identify a *subset* of these terms that are most important for capturing physics. This is fundamentally different from a naive curve fit in which all possible terms are included in the model. Thus, forward computations are necessarily fewer compared to other techniques, such as NNs, which postulate a full rank model. Despite determining a simpler model form, we demonstrate comparable performance to NNs which were trained on larger data sets. Further, the resultant model is both simple and algebraic, making for a lighter and more efficient integration with existing solvers.

Beyond developing the methodology, its utility is demonstrated on two canonical cases: homogeneous free shear turbulence and turbulence through a periodically constricted channel. Within the context of homogeneous free shear turbulence, the sparse regression methodology is validated using a “toy” problem in which the training data set is synthetically generated using an existing model. Then sparse regression is used to recover this existing model. Subsequent cases are based upon DNS data and seek to uncover improved models in comparison with existing closures. Finally, the algorithm is given experimental data for training and this result is compared with those determined using full-field DNS data.

II. METHODOLOGY

The sparse regression approach expands upon the data-driven technique presented in Brunton *et al.* [33] for using temporally evolving data to “discover” nonlinear, dynamical systems. Rather than uncovering governing equations, this method is employed to identify robust, data-driven closure models. In this section, the sparse identification of nonlinear dynamics framework [33] is built upon by adapting it for the RANS closure problem and embedding invariance—a key property of any candidate RANS model.

It is first postulated that a tensor quantity of interest, \mathbb{D} , can be characterized by the linear combination of an invariant tensor basis, represented as \mathbb{T} , premultiplied by optimal coefficients, represented as $\hat{\beta}$,

$$\mathbb{D} = \mathbb{T} \hat{\beta}. \quad (1)$$

Using this postulated form of the model, the following objective function is minimized in order to determine the optimal coefficient vector, according to

$$\hat{\beta} = \min_{\beta} \|\mathbb{D} - \mathbb{T} \beta\|_2^2 + \lambda \|\beta\|_1, \quad (2)$$

where β represents intermediary realizations of the coefficient vector which may not necessarily be the optimal coefficient vector, $\hat{\beta}$. Here the L_2 and the L_1 norms are denoted by $\|\cdot\|_2^2$ and $\|\cdot\|_1$, respectively. The first term in the objective function is ordinary least squares, which regresses the coefficient vector to the trusted data, and the second term is a sparsity-inducing penalty on the coefficients. By choice of the L_1 penalty, the minimization of the objective function performs model selection by inducing sparsity (e.g., several of the terms of $\hat{\beta}$ are identically zero, indicating that the associated term in the invariant basis, \mathbb{T} , is not important. The interested reader can refer to Refs. [34–36] for further information.) Minimization of the cost function is performed using the open source iterative algorithm presented in Brunton *et al.* [33].

In order to obtain a model that is both compact and frame-invariant, consideration must be given to the construction of the trusted data vector, \mathbb{D} , and the invariant basis, \mathbb{T} . For compactness, \mathbb{D} , and as a consequence β , are restricted to column vectors. This ensures that the coefficients for each term in the model is a scalar, which guarantees the same model form regardless of orientation (e.g., if the coefficients were vectors or tensors, this would embed directionality into the coefficients and thereby enslave the model to the orientation in which it was learned).

In this work, \mathbb{D} is assembled by first assessing the symmetry of the problem. All nonzero, unique entries in the trusted data tensor are concatenated into a column vector. For example, as seen in Fig. 1, if \mathcal{D} is symmetric in the y and z directions and the only anisotropic contribution is in the x - y direction, then the full tensor is represented as $[\mathcal{D}_{11}, \mathcal{D}_{12}, \mathcal{D}_{22}]^T$. For each realization (e.g., in time) and for each configuration under consideration, these column vectors are vertically concatenated.

Finally, form (Galilean) invariance in the resultant model is guaranteed by assembling \mathbb{T} from an invariant tensor basis. The basis is crafted by using dimensional analysis to determine the relevant known tensor quantities that fully describe the physics under study. These tensors are then used to assemble a minimal integrity basis (see, e.g., Pope [37], Speziale *et al.* [38], or Ling *et al.* [12]), using the following arguments:

- (1) Any tensor can be represented by an infinite tensor sum of the form:

$$\mathcal{D}_{ij} = \sum_{n=1}^{\infty} G^{(n)} \mathcal{T}_{ij}^{(n)},$$

where G are coefficients that in the general sense may be functions of the invariants of the tensor basis $\mathcal{T}_{ij}^{(n)}$.

- (2) In some cases, the Cayley-Hamilton theorem can be leveraged to reduce the infinite tensor sum to a finite sum that still exactly represents the infinite sum. In cases where this is not possible, the basis is truncated once model improvement stagnates.

(which may be spatial, temporal or both) and the fluctuating portions of the velocity, respectively. Applying Reynolds averaging to the incompressible Navier-Stokes equations yields the RANS equations:

$$\frac{\partial \langle u_i \rangle}{\partial x_i} = 0, \quad (3)$$

$$\frac{\partial \langle u_i \rangle}{\partial t} + \langle u_k \rangle \frac{\partial \langle u_i \rangle}{\partial x_k} = -\frac{1}{\rho} \frac{\partial \langle p \rangle}{\partial x_i} + \frac{\partial}{\partial x_j} \left[\nu \left(\frac{\partial \langle u_i \rangle}{\partial x_j} + \frac{\partial \langle u_j \rangle}{\partial x_i} \right) - \langle u'_i u'_j \rangle \right]. \quad (4)$$

It is notable that the Reynolds averaging process yields a Reynolds stress term, $\langle u'_i u'_j \rangle$, which requires closure.

The strategy for closure of the Reynolds stress term generally falls into two categories: (1) an algebraic closure or (2) the inclusion of a transport equation for the Reynolds stresses. In this work, two flows serve as case studies for the implementation of the methodology described in Sec. II. The first case study (homogeneous free shear turbulence) will develop closures in the form of transport of the Reynolds stresses and the second (turbulence in a periodically constricted channel) will consider algebraic closure.

A. Homogeneous free shear turbulence

1. Problem statement

The flow configuration under consideration in this section is homogeneous free shear turbulence, in which an unbounded, three-dimensional fluid volume is subjected to a mean-velocity gradient that generates and sustains turbulence. After sufficient time, the Reynolds stresses reach a “self-similar” state, characterized by the anisotropy of the Reynolds stresses reaching stationarity in time [e.g., $\frac{d}{dt}(\langle u'_i u'_j \rangle / k) = 0$ with $k = \langle u'_k u'_k \rangle$ the turbulent kinetic energy (TKE)]. Consequently, Reynolds-averaged quantities are statistically one-dimensional (i.e., they depend only on time). It is this “self-similar” behavior that is of specific interest in formulating an improved RANS closure.

As previously described, the Reynolds stresses in the RANS equations [Eq. (4)] require closure. In this example, we consider the transport of the Reynolds stresses, which are given exactly as

$$\begin{aligned} \frac{D \langle u'_i u'_j \rangle}{Dt} = & \underbrace{- \left[\langle u'_j u'_k \rangle \frac{\partial \langle u_i \rangle}{\partial x_k} + \langle u'_i u'_k \rangle \frac{\partial \langle u_j \rangle}{\partial x_k} \right]}_{\text{production, } \mathcal{P}_{ij}} - \underbrace{2\nu \left\langle \frac{\partial u'_i}{\partial x_k} \frac{\partial u'_j}{\partial x_k} \right\rangle}_{\text{dissipation, } \varepsilon_{ij}} + \underbrace{\left\langle \frac{p}{\rho} \left(\frac{\partial u'_i}{\partial x_j} + \frac{\partial u'_j}{\partial x_i} \right) \right\rangle}_{\text{redistribution, } \mathcal{R}_{ij}} \\ & - \underbrace{\frac{\partial}{\partial x_k} \langle u'_i u'_j u'_k \rangle}_{\text{turbulent convection}} + \underbrace{\nu \frac{\partial^2 \langle u'_i u'_j \rangle}{\partial x_k^2}}_{\text{viscous diffusion}} - \underbrace{\frac{\partial}{\partial x_k} \left(\frac{\langle u'_i p' \rangle}{\rho} \delta_{jk} + \frac{\langle u'_j p' \rangle}{\rho} \delta_{ik} \right)}_{\text{pressure transport}}, \end{aligned} \quad (5)$$

where D/Dt denotes the material derivative, ν is the kinematic viscosity, and ρ and p denote fluid density and pressure, respectively. In the case of homogeneous free shear turbulence, the domain is spatially homogeneous and consequently, spatial gradients of mean quantities are null. Thus, the transport of Reynolds stresses is reduced to

$$\frac{d \langle u'_i u'_j \rangle}{dt} = \underbrace{-[\langle u'_j u'_k \rangle \Gamma_{ik} + \langle u'_i u'_k \rangle \Gamma_{jk}]}_{\text{production, } \mathcal{P}_{ij}} - \underbrace{2\nu \left\langle \frac{\partial u'_i}{\partial x_k} \frac{\partial u'_j}{\partial x_k} \right\rangle}_{\text{dissipation, } \varepsilon_{ij}} + \underbrace{\left\langle \frac{p}{\rho} \left(\frac{\partial u'_i}{\partial x_j} + \frac{\partial u'_j}{\partial x_i} \right) \right\rangle}_{\text{redistribution, } \mathcal{R}_{ij}}, \quad (6)$$

where the shear rate tensor is given as $\Gamma_{ij} = \partial \langle u_i \rangle / \partial x_j$. Here the production term is closed, however the dissipation and redistribution tensors both require closure. In this work, new modeling efforts are directed toward the redistribution tensor and the dissipation tensor is closed using the standard

transport equation proposed by Hanjalic and Launder [39],

$$\frac{\partial \varepsilon}{\partial t} = C_{\varepsilon 1} \frac{\mathcal{P} \varepsilon}{k} - C_{\varepsilon 2} \frac{\varepsilon^2}{k}, \quad (7)$$

where $\mathcal{P} = \text{tr}(\mathcal{P}_{ij})/2$ and model constants are given by $[C_{\varepsilon 1}, C_{\varepsilon 2}] = [1.44, 1.92]$ [40].

2. Proof-of-concept: A synthetic data set

As an initial proof-of-concept for the sparse regression methodology described in Sec. II, a set of data is generated using a well-established closure for the redistribution tensor with the goal of recovering the known model. The closure utilized to generate the synthetic data set was proposed by Launder *et al.* [41] and is known as the LRR-IP model,

$$\mathcal{R}_{ij} = -C_R \frac{\varepsilon}{k} \left(\langle u'_i u'_j \rangle - \frac{2}{3} k \delta_{ij} \right) - C_2 \left(\mathcal{P}_{ij} - \frac{2}{3} \mathcal{P} \delta_{ij} \right), \quad (8)$$

where the constants are given as $[C_R, C_2] = [1.8, 0.6]$ [41]. This closure, embedded in the transport equation for the Reynolds stresses in Eq. (6) and the transport equation for dissipation given in Eq. (7) are solved for three shear rates ($\Gamma = \Gamma_{12} = [2.25, 11.24, 20.23]$). This results in one-dimensional (time-dependent) data for the Reynolds stresses for each shear rate.

Given the simple flow configuration, the redistribution tensor can be normalized by the viscous dissipation rate, ε , and characterized by a linear combination of the the following nondimensionalized, mean flow quantities:

- (1) Anisotropic stress tensor $b_{ij} = \frac{\langle u'_i u'_j \rangle}{2k} - \frac{1}{3} \delta_{ij}$
- (2) Mean rotation rate tensor $\hat{R}_{ij} = \frac{1}{2} \frac{k}{\varepsilon} \left(\frac{\partial \langle u_i \rangle}{\partial x_j} - \frac{\partial \langle u_j \rangle}{\partial x_i} \right)$
- (2) Mean shear rate tensor $\hat{S}_{ij} = \frac{1}{2} \frac{k}{\varepsilon} \left(\frac{\partial \langle u_i \rangle}{\partial x_j} + \frac{\partial \langle u_j \rangle}{\partial x_i} \right)$

such that

$$\mathcal{R}_{ij} = \varepsilon \mathbf{f}(b_{ij}, \hat{R}_{ij}, \hat{S}_{ij}), \quad (9)$$

where \mathbf{f} is a *form invariant* tensor-valued function, which, due to the linearity in b_{ij} , \hat{R}_{ij} , and \hat{S}_{ij} automatically satisfies

$$\mathbf{Q} \mathbf{f}(b_{ij}, \hat{R}_{ij}, \hat{S}_{ij}) \mathbf{Q}^T = \mathbf{f}(\mathbf{Q} b_{ij} \mathbf{Q}^T, \mathbf{Q} \hat{R}_{ij} \mathbf{Q}^T, \mathbf{Q} \hat{S}_{ij} \mathbf{Q}^T). \quad (10)$$

Here \mathbf{Q} is a Galilean rotation matrix, e.g., $\mathbf{Q} \mathbf{Q}^T = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ (where \mathbf{I} is the identity tensor) and $\det \mathbf{Q} = \pm 1$.

These bases have been extensively used in the literature (see, e.g., Refs. [38,42] for their derivation) and their usage is restricted to modeling equilibrium regimes (rather than the transient period). To briefly summarize, because the redistribution tensor, \mathcal{R}_{ij} , is symmetric and deviatoric, and its dependence on each of the bases is linear, each basis tensor must also satisfy these same properties. Further, the constraint of form invariance under coordinate transformation stipulates that \mathbf{f} must be an isotropic function of its arguments (i.e., b_{ij} , \hat{R}_{ij} , and \hat{S}_{ij}). Using these constraints guides the formulation of the minimal integrity basis [38] as shown in the leftmost column of Table I.

Using the data generated by solving Eqs. (6)–(14), the basis tensors are computed and the redistribution tensor is populated by taking the time derivative of the Reynolds stresses (using a sixth-order central difference scheme) and solving Eq. (6) for the redistribution tensor. Then these quantities are assembled into \mathbb{D} and \mathbb{T} as described in Sec. II. Note that since we are interested in modeling the self-similar regime, only data from this region are used for training and for assessing model error. After \mathbb{D} and \mathbb{T} are assembled, the cost functional defined by Eq. (2) is optimized for decreasing values of λ until reduction in model error is no longer achieved.

As shown in Table I, the methodology exactly returns the LRR-IP model used to generate the data set. In order to systematically challenge the robustness of the algorithm, *a posteriori*, artificial noise was added to the synthetic data set. The synthetic noise was normally distributed about the

TABLE I. Summary of model forms and associated error in the self-similar region of homogeneous free shear turbulence, with increasing amounts of artificial noise added to the synthetic data set.

Order in b_{ij}	$\mathcal{T}^{(n)}$	LRR-IP	Sparse regression			
			$\lambda = 0.1$ $\mathcal{N} = 0$	$\lambda = 0.1$ $\mathcal{N}(\mu, 0.1\mu)$	$\lambda = 0.1$ $\mathcal{N}(\mu, 0.2\mu)$	$\lambda = 0.5$ $\mathcal{N}(\mu, 0.3\mu)$
0	S_{ij}	0.8	0.8	0.8010	0.8020	0.803
	b_{ij}	-3.6	-3.6	-3.5761	-3.5522	-3.5282
1	$R_{il}b_{lj} + R_{jl}b_{li}$	1.2	1.2	1.2003	1.2007	1.2010
	$S_{il}b_{lj} + S_{jl}b_{li} - \frac{2}{3}S_{lm}b_{ml}\delta_{ij}$	1.2	1.2	1.2018	1.2036	1.2054
	$b_{ij}^2 - \frac{1}{3}b_{il}^2\delta_{ij}$	0	0	0	0	0
2	$S_{il}b_{lj}^2 + S_{jl}b_{li}^2 - \frac{2}{3}S_{lm}b_{ml}^2\delta_{ij}$	0	0	0	0	0
	$R_{il}b_{lj}^2 + R_{jl}b_{li}^2$	0	0	0	0	0
3	$b_{ik}^2 R_{kp} b_{pj} - b_{il} R_{lk} b_{kj}^2$	0	0	0	0	0
	$\epsilon^{b_{ij}}$	-	0.0	0.0076	0.015	0.023

mean of the synthetic data, denoted as $\mathcal{N}(\mu, \sigma)$, where σ is the standard deviation that is prescribed in terms of percentage of the mean value, μ . We consider $\sigma = 0.1, 0.2$, and 0.3μ . In each case, λ was reduced until the model error plateaued. Even in the case of the noisiest data provided, the learned model deviated from the expected LRR-IP model by only 2.3%, where error is defined by the L_2 norm

$$\epsilon = \frac{\|\mathbb{D} - \mathbb{T}\beta\|_2}{\|\mathbb{D}\|_2}. \quad (11)$$

This level of performance indicates the sparse regression methodology is robust to substantial noise in the training data without compromising the accuracy in learning the underlying physics. This is further demonstrated in Figs. 2(a)–2(c) where the models learned from the noisy data are shown against the LRR-IP model and the artificially noisy data. In all three cases, the learned model accurately describes the behavior of the LRR-IP model despite small amounts of error in the coefficients.

3. DNS-generated data: Can sparse regression improve upon existing models?

Here the same physical configuration is considered, albeit with the trusted data generated using DNS (see Fig. 3). The goal of using DNS-generated training data is to understand how the learned models differ from classically used models.

To generate the DNS data sets, NGA [43], a fully conservative, low-Mach number finite volume solver is used. A pressure Poisson equation is solved to enforce continuity via fast Fourier

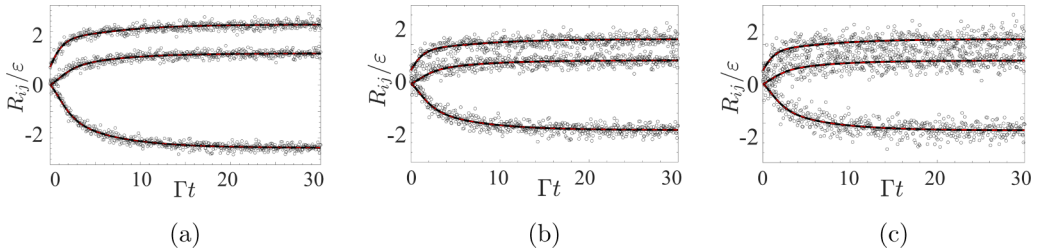


FIG. 2. Comparison between prescribed LRR-IP model (—) with the learned models (---) and artificially noisy data (○), for $\Gamma = 2.25$. (a) $\mathcal{N}(\mu, 0.1\mu)$, (b) $\mathcal{N}(\mu, 0.2\mu)$, and (c) $\mathcal{N}(\mu, 0.3\mu)$.

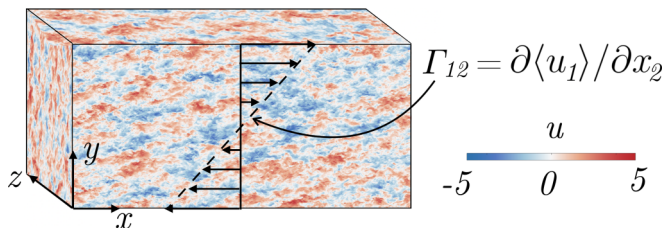


FIG. 3. Snapshot of the instantaneous velocity field in DNS homogeneous free shear turbulence at $\mathcal{S} = 11.2$.

transforms in all three periodic directions. The Navier-Stokes equations are solved on a staggered grid with second order spatial accuracy and time is advanced with second-order accuracy using the semi-implicit Crank-Nicolson scheme of Pierce [44]. Shear periodic boundary conditions are enforced using the recently developed algorithm of Kasbaoui *et al.* [45]. Turbulence in the domain is initialized using spectral methods in order to ensure consistency with Kolmogorov’s “ $-5/3$ ” spectrum [46,47].

Five cases are simulated for nondimensional shear rates $\mathcal{S} = 2\Gamma k_0/\varepsilon_0 = (2.3, 6.6, 11.2, 13.2, 20.2)$ on a grid of size $1024 \times 512 \times 512$, corresponding to a domain size of $2\pi \times \pi \times \pi$. Here k_0 and ε_0 denote the initial values of TKE and dissipation, respectively. The grid resolution ensures that the flow captures the dissipative scales. Each case is simulated to a nondimensional time of $\Gamma t \approx 25\text{--}30$ to ensure sufficient data in the self-similar region is captured. Of the five data sets, three are selected as training sets ($\mathcal{S} = 2.3, 11.2, 20.2$) from which a new model is learned. The remaining two data sets ($\mathcal{S} = 6.6, 13.2$) serve as validation sets in order to determine the optimal value of λ and therefore the optimal learned model.

In the same fashion as was described for the synthetic data set, the DNS data are organized into \mathbb{D} and the \mathbb{T} , and the cost functional is optimized for decreasing values of the sparsity parameter λ until all terms are populated in the learned model. The resulting models from this procedure are shown in Table II. As λ is decreased, additional terms are included in the learned model and the coefficients adjust accordingly. The four learned models are compared against existing models, the

TABLE II. Summary of learned and existing models and associated error in the self-similar region for homogeneous free shear turbulence.

Order	in b_{ij}	$\mathcal{T}^{(n)}$	Sparse regression						
			Rotta	LRR-IP	LRR-QI	$\lambda = 0.75$	$\lambda = 0.6$	$\lambda = 0.5$	$\lambda = 0$
0		S_{ij}	0	0.8	0.8	1.01	1.01	0.98	0.98
		b_{ij}	-3.6	-3.6	-3.0	1.27	1.31	1.45	1.46
1		$R_{il}b_{lj} + R_{jl}b_{li}$	0	1.2	1.31	1.53	1.56	1.49	1.48
		$S_{il}b_{lj} + S_{jl}b_{li} - \frac{2}{3}S_{lm}b_{ml}\delta_{ij}$	0	1.2	1.74	1.73	1.71	1.79	1.78
		$b_{ij}^2 - \frac{1}{3}b_{ii}^2\delta_{ij}$	0	0	0	5.22	4.64	7.02	6.71
2		$S_{il}b_{lj}^2 + S_{jl}b_{li}^2 - \frac{2}{3}S_{lm}b_{ml}^2\delta_{ij}$	0	0	0	0	0	0.57	0.56
		$R_{il}b_{lj}^2 + R_{jl}b_{li}^2$	0	0	0	0	0	0	0.13
		$b_{ik}^2 R_{kp} b_{pj} - b_{il} R_{lk} b_{kj}^2$	0	0	0	0	-0.65	2.08	2.45
3		Training error, $\epsilon_{\text{train}}^b$	0.68	0.26	0.26	0.090	0.092	0.078	0.073
		Testing error, ϵ_{test}^b				0.086	0.089	0.070	0.078

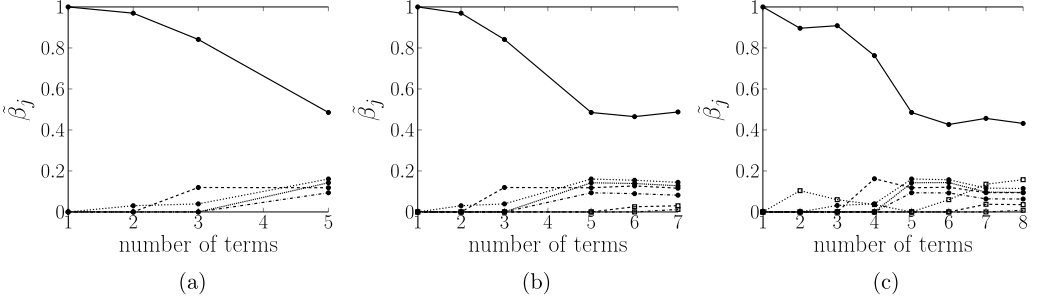


FIG. 4. As the number of terms in the model increases (by decreasing λ), terms that are most important for capturing key redistribution physics arise in the sparsest models and persist with prominent coefficients as terms are added. (a) Up to first order terms, (b) Up to second order terms, and (c) Up to third order terms.

Rotta [48] and LRR-IP and LRR-QI models [41] [Eqs. (12)–(15)], written in terms of the basis tensors as

$$\mathcal{R}_{ij}^{\text{Rotta}} = -2C_R b_{ij}, \quad (12)$$

$$\mathcal{R}_{ij}^{\text{LRR-IP}} = -2C_R b_{ij} + \frac{4}{3}C_2 S_{ij} + 2C_2 (R_{il} b_{lj} + R_{jl} b_{li}) + 2C_2 (S_{il} b_{lj} + S_{jl} b_{li} - \frac{2}{3} S_{lm} b_{ml} \delta_{ij}), \quad (13)$$

$$\begin{aligned} \mathcal{R}_{ij}^{\text{LRR-QI}} = & -2C'_R b_{ij} + \frac{4}{3}S_{ij} + \frac{2}{11}(10 - 7C'_2)(R_{il} b_{lj} + R_{jl} b_{li}) \\ & + \frac{6}{11}(2 + 3C'_2)(S_{il} b_{lj} + S_{jl} b_{li} - \frac{2}{3}S_{lm} b_{ml} \delta_{ij}). \end{aligned} \quad (14)$$

Here the coefficients are given as $[C_R, C_2] = [1.8, 0.6]$ and $[C'_R, C'_2] = [1.5, 0.4]$. The Rotta model assumes a linear relationship with the anisotropy tensor and thus models a linear return to isotropy. In contrast, the LRR-IP model includes nonlinear terms that are important for characterizing homogeneous *anisotropic* turbulence. In comparing these three models with four learned models of increasing complexity, it is observed that the least complex learned model, corresponding to $\lambda = 0.75$, already shows marked improvement over the highest performing existing models and reduces error in the anisotropic stress tensor from 26% to 9%.

Shown in Fig. 4, as λ is decreased and terms are added to the learned model, the normalized coefficients, $\tilde{\beta}_i = \beta_i / (\max \beta_i)$, change to accommodate contributions from additional terms. In Fig. 4(a), the sparse regression methodology is employed on a basis that is restricted to up to first order in b_{ij} . The basis is then expanded to second- and third-order terms in b_{ij} in Figs. 4(b) and 4(c), respectively. In each instance, the most prominent coefficient remains the largest contribution to the learned model, though its contribution decreases as subsequent terms are added. The order of prominence of the lesser contributing terms does not remain fixed once the number of terms in the model grows. This behavior has an insignificant effect on model performance and sensitivity and serves to demonstrate the relative lesser importance of these terms to the overall model performance as compared with the terms with larger contributions. Finally, it is observed in Figs. 4(b) and 4(c) that the dominant coefficient stagnates beyond a five-term model. This mirrors the reduction in overall model error as shown in Table II. The ability to identify which basis tensors are most important for modeling the flow and their relative sensitivity therein, is a unique benefit of the sparse regression methodology which allows for interpretability of the relationship between model form and flow physics.

A comparison of the training and validation errors give the clearest indication of when a learned model begins to exhibit symptoms of overfitting (see Fig. 5). While the learned closure predicts the redistribution tensor, \mathcal{R}_{ij} , the ultimate goal is to improve performance in predicting anisotropy in the Reynolds stresses, b_{ij} , making both measures of error relevant to assessing learned models. As shown in Fig. 5, the error in \mathcal{R}_{ij} decreases monotonically beyond a two-term model, however, we

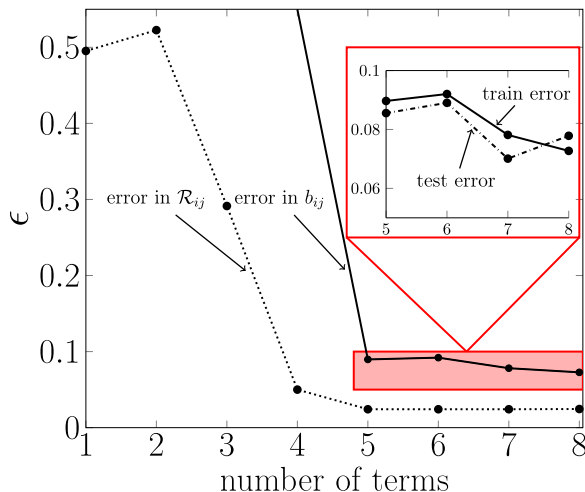


FIG. 5. Error in \mathcal{R}_{ij} and b_{ij} are shown for models of increasing complexity for homogeneous free shear turbulence. The inset figure delineates validation (test) and training error.

observe that five terms are required for stability in the transport equation for the Reynolds stresses. It is notable that while the LRR-IP and LRR-QI model are stable with four terms, the four-term model learned by sparse regression is not. This is likely due to a lack of a stability penalty in the cost function. Adding additional penalties to the cost function in order to enforce model stability is an area of active research. Testing and training errors are also compared in Fig. 5. As might be expected, the training error generally decreases as terms are added, but beyond seven terms in the learned model, an increase in testing error is observed. This is indicative of overfitting, thus making the seven-term model the ideal model that minimizes model error while maximizing accuracy of the model across different shear rates. As seen in Table II, the ideal learned model reduces error in predicting self-similar behavior by more than half as compared with the LRR-IP or LRR-QI models and more than eightfold as compared with the Rotta model.

In Fig. 6 the ideal seven-term model is compared with the highest performing existing model, the LRR-QI model. Both are plotted against the DNS values used for training [Figs. 6(a)–6(c)] and for validation [Figs. 6(d)–6(e)]. As previously discussed, it is observed that the learned model accurately captures the self-similar behavior (shown in gray shaded regions) of the normalized Reynolds stresses even in the testing cases which were not seen by the sparse regression method during training.

4. A note on noninertial frames of reference

If a *noninertial* frame is to be considered [42,49,50], one would need to modify the normalized, mean rotation rate tensor to include the rotation rate of the frame with respect to an inertial frame (Ω), i.e., $\hat{\mathbb{R}}_{ij} = \hat{R}_{ij} + \epsilon_{mji}\Omega_m$, where ϵ_{mji} denotes the permutation tensor. Additionally, Coriolis terms, $(\langle u_i u_k \rangle \epsilon_{mkj} \Omega_m + \langle u_j u_k \rangle \epsilon_{mki} \Omega_m)$, must be included in Eqs. (4) and (6).

5. A note on constant coefficients

The analysis presented above considers the simple case in which model coefficients are constants. As previously discussed, the coefficients are permitted to theoretically depend nonlinearly on the principal invariants of the basis tensors. In the case of homogeneous free shear turbulence, model performance using constant coefficients performs well without the additional complexity of dependency on principal invariants. However, in situations in which this is not the case, this dependency can be added into \mathbb{T} by postulating functional forms of coefficients and appending

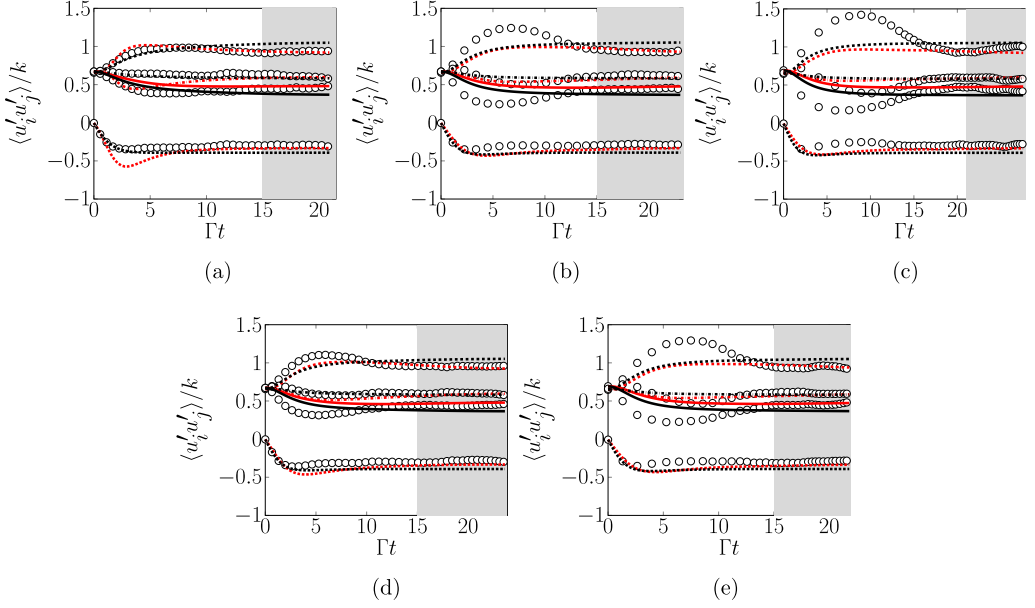


FIG. 6. Sparse regression (—) produces a more accurate model as compared with the most accurate traditional closure available (LRR-QI, —). DNS data are denoted by open circles, and the four lines correspond to the unique components of the normalized Reynolds stresses: $\langle u'u' \rangle$: - - -, $\langle u'v' \rangle$:, $\langle v'v' \rangle$: —, and $\langle w'w' \rangle$: -.-.-. The shaded portion denotes the self-similar region of the flow. (a) $S = 3.2$, (b) $S = 16.1$, (c) $S = 30.7$, (d) $S = 10.0$, and (e) $S = 20.1$.

them to each of the basis tensors. If there is a physics-based rationale for the functional dependency, this process can be prescribed by hand, or if the functional dependence is not in any way constrained, an algorithm such as gene expression programming [26–30] can be used to analytically determine complex coefficient dependencies on the principal invariants. This strategy is reserved for future work.

B. Turbulent flow through a periodically constricted channel

1. Problem statement

In this section, we consider the classical case of turbulent flow through a periodically constricted channel as shown in Fig. 7 and described in Breuer *et al.* [51]. As discussed in Sec. III, two main approaches are typically taken when developing closures for the Reynolds stresses. In Sec. III A the transport of the Reynolds stresses was addressed, and in this section algebraic closure of the Reynolds stresses will be developed.

In this strategy, the algebraic closure for the Reynolds stresses depends upon a model for the anisotropic stress tensor, such that $\langle u_i' u_j' \rangle = 2k(b_{ij} + \frac{1}{3}\delta_{ij})$. Further it has been well established that the model for b_{ij} depends upon \hat{S}_{ij} and \hat{R}_{ij} . Recalling from Sec. III A that these quantities are normalized by TKE, k , and dissipation of TKE, ε , this method requires the transport of both k and ε , which are given by

$$\frac{\partial k}{\partial t} + \frac{\partial(ku_i)}{\partial x_i} = \frac{\partial}{\partial x_j} \left[\left(v + \frac{v_t}{\sigma_k} \right) \frac{\partial k}{\partial x_j} \right] + \mathcal{P} - \varepsilon, \quad (15)$$

$$\frac{\partial \varepsilon}{\partial t} + \frac{\partial(\varepsilon u_i)}{\partial x_i} = \frac{\partial}{\partial x_j} \left[\left(v + \frac{v_t}{\sigma_\varepsilon} \right) \frac{\partial \varepsilon}{\partial x_j} \right] + C_{1\varepsilon} \frac{\varepsilon}{k} \mathcal{P} - C_{2\varepsilon} \frac{\varepsilon^2}{k}, \quad (16)$$

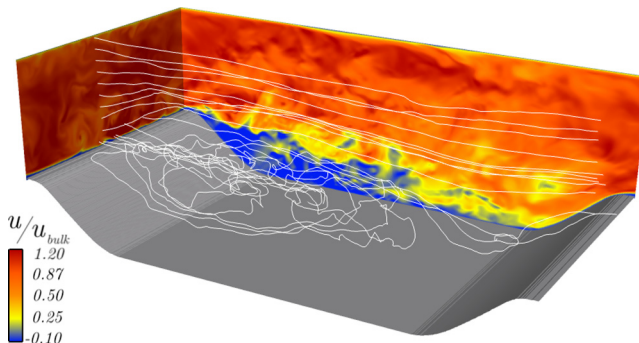


FIG. 7. Instantaneous streamwise velocity (color), with streamlines originating from $x = L_z/2$ (white lines).

where $[C_\mu, \sigma_k, \sigma_\varepsilon, C_{1\varepsilon}, C_{2\varepsilon}] = [0.09, 1.00, 1.30, 1.44, 1.92]$. The turbulent viscosity, ν_t , is given as $\nu_t = C_\mu k^2 / \varepsilon$ where $C_\mu = 0.09$ [52].

Using these equations along with a model for the anisotropic stress tensor, the RANS equations [Eqs. (3) and (4)] are closed. The aim of this study is to use sparse regression to develop an improved algebraic closure for the anisotropic stress tensor. As is commonly used in the literature, the configuration under consideration here is turbulent flow through a periodically constricted channel (see Fig. 7). This flow configuration is particularly challenging because the quantities of interest are statistically two-dimensional (with dependence on the streamwise and cross-stream directions) and the presence of the constriction generates massive separation in the flow.

The data set used for training was simulated using NGA, described in Sec. III A. The top and bottom walls apply a no-slip boundary condition and the bottom, constricted wall is enforced using a cut-cell immersed boundary method [53]. The geometry for the configuration under study matches the configuration described in Breuer *et al.* [51] with uniform grid spacing discretized by $[N_x, N_y, N_z] = [512, 380, 380]$. A Reynolds number of 2800 is considered, where $\text{Re} = u_{\text{bulk}} h / \nu$. The bulk velocity u_{bulk} is given by the mean velocity at the hill crest, and h is the hill height. After reaching a statistically stationary point, the DNS data were averaged in the cross-stream (z direction) and temporally for 44 flow-through times.

A linear eddy viscosity model (LEVM) is frequently used to close the Reynolds stresses that appear in the RANS equations [52]. This closure takes the form

$$b_{ij} = -C_\mu \hat{S}_{ij}, \quad (17)$$

which will serve for comparison purposes as the “existing” model.

As outlined in Sec. II, the basis on which to train the model must first be identified. As previously derived [37], a minimal integrity basis for the anisotropy tensor can be formulated using the normalized mean rotation and shear rate tensors, \hat{R}_{ij} and \hat{S}_{ij} , respectively. Since the anisotropy stress tensor is symmetric and deviatoric, each of $\mathcal{T}_{ij}^{(n)}$ must also have these properties. After formulating combinations of \hat{S}_{ij} and \hat{R}_{ij} with these properties, and owing to the Cayley-Hamilton theory, all symmetric and deviatoric tensors that are combinations of \hat{S}_{ij} and \hat{R}_{ij} can be formed as a linear combination of the 10 basis tensors shown in Table III [37].

Using this basis, the anisotropic stress tensor can be represented exactly as

$$b_{ij} = \sum_{n=1}^{10} G^{(n)} \mathcal{T}_{ij}^{(n)}(\hat{R}_{ij}, \hat{S}_{ij}). \quad (18)$$

In the case of statistically two-dimensional flows, as is the case here, the basis simplifies to only three tensors and the coefficients depend on at most only two invariants as shown in Table IV [37,42].

TABLE III. The 10 tensor bases that exactly describe the anisotropic stress tensor.

$\mathcal{T}_{ij}^{(1)} = \hat{S}_{ij}$	$\mathcal{T}_{ij}^{(6)} = \hat{R}_{ik}\hat{R}_{kl}\hat{S}_{lj} + \hat{S}_{ik}\hat{R}_{kl}\hat{R}_{lj} - \frac{2}{3}\hat{S}_{pk}\hat{R}_{kl}\hat{R}_{lp}\delta_{ij}$
$\mathcal{T}_{ij}^{(2)} = \hat{S}_{ik}\hat{R}_{kj} - \hat{R}_{ik}\hat{S}_{kj}$	$\mathcal{T}_{ij}^{(7)} = \hat{R}_{ik}\hat{S}_{kl}\hat{R}_{lp}\hat{R}_{pj} - \hat{R}_{ik}\hat{R}_{kl}\hat{S}_{lp}\hat{R}_{pj}$
$\mathcal{T}_{ij}^{(3)} = \hat{S}_{ik}\hat{S}_{kj} - \frac{1}{3}\hat{S}_{ik}\hat{S}_{kl}\delta_{ij}$	$\mathcal{T}_{ij}^{(8)} = \hat{S}_{ik}\hat{R}_{kl}\hat{S}_{lp}\hat{S}_{pj} - \hat{S}_{ik}\hat{S}_{kl}\hat{R}_{lp}\hat{S}_{pj}$
$\mathcal{T}_{ij}^{(4)} = \hat{R}_{ik}\hat{R}_{kj} - \frac{1}{3}\hat{R}_{ik}\hat{R}_{kl}\delta_{ij}$	$\mathcal{T}_{ij}^{(9)} = \hat{R}_{ik}\hat{R}_{kl}\hat{S}_{lp}\hat{S}_{pj} + \hat{S}_{ik}\hat{S}_{kl}\hat{R}_{lp}\hat{R}_{pj} - \frac{2}{3}\hat{S}_{qk}\hat{S}_{kl}\hat{R}_{lp}\hat{R}_{pq}$
$\mathcal{T}_{ij}^{(5)} = \hat{R}_{ik}\hat{S}_{kl}\hat{S}_{lj} - \hat{S}_{ik}\hat{S}_{kl}\hat{R}_{lj}$	$\mathcal{T}_{ij}^{(10)} = \hat{R}_{ik}\hat{S}_{kl}\hat{S}_{lp}\hat{R}_{pq}\hat{R}_{qj} - \hat{R}_{ik}\hat{R}_{kl}\hat{S}_{lp}\hat{S}_{pq}\hat{R}_{qj}$

Following the sparse regression methodology described in Sec. II, the DNS data set is formulated into \mathbb{D} and \mathbb{T} . However, instead of modeling b_{ij} directly, the anisotropic stress tensor is split into linear and nonlinear portions, denoted by b_{ij}^{\parallel} and b_{ij}^{\perp} , respectively. The linear portion will be taken as the standard LEVM and the nonlinear portion will be the subject of modeling efforts:

$$b_{ij} = b_{ij}^{\perp} + b_{ij}^{\parallel} \quad (19)$$

$$= b_{ij}^{\perp}(k, \varepsilon, \hat{S}_{ij}, \hat{R}_{ij}) - C_{\mu}\hat{S}_{ij}, \quad (20)$$

$$\mathcal{D}_{ij} = b_{ij}^{\perp} = b_{ij} + C_{\mu}\hat{S}_{ij}. \quad (21)$$

This strategy is employed based upon the recommendation of several works that have pointed out the ill-conditioning of the RANS equations [54]. These works suggest that separating the model into a linear portion (solved implicitly with the viscous terms in the RANS solver) and a nonlinear portion (solved explicitly) improves stability of the integrated RANS solver [17,28]. Further, since the standard LEVM model is used as the starting point for modeling, the basis is formulated using data from a forward solution in OpenFOAM [55] using the LEVM closure. Because the $k - \varepsilon$ equations contain models and are thereby a source of error in the “trusted” training data, these data must be used as a starting point for modeling.

Using this formulation, sparse regression is employed to discover an improved model. This effort results in both an *a priori* and an *a posteriori* analysis of the model. In the former analysis, the training data are used to evaluate the accuracy of the learned model within the context of predicting the anisotropy tensor. In the latter analysis, the learned model is implemented in OpenFOAM and the forward solution is compared against the trusted DNS data and the existing LEVM. Additionally, as an “upper end” metric, a look-up table was provided to the OpenFOAM RANS solver for the Reynolds stress terms that appear in both momentum and production in the $k - \varepsilon$ equations. This data set serves as the performance of an ideal model that exactly captures the behavior of the Reynolds stresses while highlighting the model errors associated with the $k - \varepsilon$ model equations themselves.

Two learned models are discovered using sparse regression, one with three terms ($\lambda = 0$, denoted Learned 1) and the second with two terms ($\lambda = 15$, denoted Learned 2). Both learned models take

TABLE IV. The reduced basis set for statistically two-dimensional flows.

$\mathcal{T}_{ij}^{(1)}$	\hat{S}_{ij}
$\mathcal{T}_{ij}^{(2)}$	$\hat{S}_{ik}\hat{R}_{kj} - \hat{R}_{ij}\hat{S}_{kj}$
$\mathcal{T}_{ij}^{(3)}$	$\hat{S}_{ik}\hat{S}_{kj} - \frac{1}{3}\hat{S}_{ik}\hat{S}_{kl}\delta_{ij}$
λ_1	$\hat{S}_{ik}\hat{S}_{kl}$
λ_2	$\hat{R}_{ik}\hat{R}_{kl}$

TABLE V. Summary of learned model coefficients and *a priori* errors compared with the standard LEVM.

Model	C_1	C_2	C_3	ϵ^b	RMSE
LEVM	–	–	–	1.02	0.16
Learned 1	63.12	51.42	10.98	0.64	0.10
Learned 2	63.14	51.42	0	0.64	0.10

the form

$$b_{ij}^\perp = \left[\frac{1}{1000 + \lambda_1^3} \right] (C_1 \mathcal{T}^{(1)} + C_2 \mathcal{T}^{(2)} + C_3 \mathcal{T}^{(3)}) \quad (22)$$

and are detailed in Table V. In this expression, the strain damping factor $[1/(1000 + \lambda_1^3)]$ was selected following the formulation of Shih *et al.* [56]. While a detailed discussion and derivation can be found therein, this factor is resultant of constraining model coefficients to ensure realizability conditions.

The *a posteriori* analysis for the learned models includes an assessment of recirculation (Table VI) and velocity predictions for the training case (Fig. 9) and a test case at a higher Reynolds number [Fig. 11(a)]. These results are discussed in detail in Sec. III B 3.

TABLE VI. Summary of separation and reattachment locations for all models compared with DNS. Relative error with respect to the DNS values are shown in parentheses. Dashes indicate the data are either not reported or not observed. Note that models “Learned 3” and “Learned 4” are discussed in Sec. IV.

Configuration	Re	Model	Primary		Secondary	
			Separation	Reattachment	Separation	Reattachment
Periodic hills	2800	LEVM	0.43 (1.62)	3.64 (0.32)	–	–
		Learned 1	0.40 (1.53)	5.38 (0.005)	7.14 (0.04)	7.31 (0.008)
		Learned 2	0.40 (1.48)	5.38 (0.005)	7.14 (0.04)	7.31 (0.01)
		Learned 3	0.34 (1.14)	5.54 (0.04)	–	–
		Learned 4	0.71 (3.43)	4.98 (0.07)	7.15 (0.04)	7.50 (0.03)
	DNS	0.16	5.35	6.87	7.25	
	LEVM	0.40 (1.35)	3.01 (0.40)	–	–	
	Learned 1	0.40 (1.38)	5.38 (0.05)	7.14 (0.01)	7.30 (0.001)	
	Learned 3	0.34 (0.98)	5.41 (0.07)	–	–	
	Learned 4	0.41 (1.41)	4.55 (0.09)	–	–	
DNS [51]	0.18	5.41	–	–		
DNS [57]	0.17	5.04 ± 0.09	7.04	7.31		
Experiment [58]	–	4.83	–	–		
10 600	LEVM	0.41 (1.30)	3.78 (0.26)	–	–	
	Learned 1	0.40 (1.24)	5.10 (0.002)	–	–	
	Learned 3	0.33 (0.85)	5.59 (0.10)	–	–	
	Learned 4	0.40 (1.24)	4.60 (0.10)	–	–	
	LES [51]	0.19	5.09	–	–	
Experiment [58]	–	4.21	–	–		
Backward-facing step	5000	LEVM	0.49	4.52 (0.28)	–	–
		Learned 1	0.90	6.17 (0.02)	–	–
		DNS [59]	–	6.28	–	–
		Experiment [60]	–	6 ± 0.15	–	–
		Experiment [61]	–	6.51	–	–

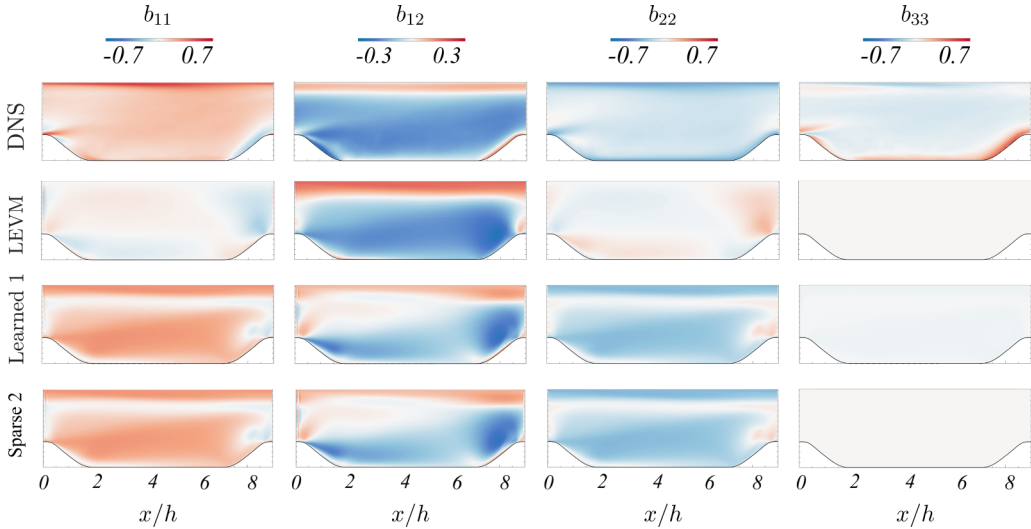


FIG. 8. Comparisons of the components of the Reynolds stress tensor computed using DNS, LEVM, and the two learned models.

2. *A priori* analysis

Each model developed can be assessed using the data with which it was trained. This represents an *a priori* assessment of the model, but has limitations as it does not take into account issues of stability or sensitivity that may be encountered within the context of a RANS solver. Further, since the forward solution is not computed here, all assessments of model accuracy are computed with respect to the anisotropic stress tensor.

Shown in Fig. 8, the standard LEVM does a reasonable job predicting the b_{12} component of the anisotropy tensor, but it struggles for the diagonal components. In all three cases, both the sign and magnitude are incorrect. The learned models, in contrast, capture the correct sign for the diagonal components and improve the magnitude inaccuracies present in the standard LEVM for the b_{12} component. However, for Learned 2, with the elimination of the third basis term, the prediction for b_{33} is also lost.

Using the L_2 norm as a metric for error, the learned models reduce model error in the anisotropic stress tensor by 41% with respect to LEVM.

3. *A posteriori* analysis

The true test of any model is its performance in the context of a forward solver. It is in this sense that model shortcomings become apparent, e.g., sensitivity or stability issues. Further, while the aim of Reynolds stress modeling is to improve accuracy in describing the stresses, the ultimate goal is that these models will improve predictions in the velocity field.

In order to assess the improvement of the learned model over the LEVM, the learned models were integrated into OpenFOAM, solved in conjunction with the $k - \varepsilon$ equations, and compared with the LEVM model and the $k - \varepsilon$ equations with a look-up table containing the DNS values for the anisotropy tensor. In each case, the RANS equations were solved on a two-dimensional grid of resolution $(N_x, N_y) = (200, 160)$ with the same physical dimensions as described in Breuer *et al.* [51] (and used for the DNS computations). Periodic conditions were imposed at the left and right faces, and “patch” conditions were imposed on the front and back faces to enforce a two-dimensional solution. The bottom and top walls were treated as no-slip and a forcing term was added such that the velocity at the top of the hill crest enforced the desired Reynolds number.

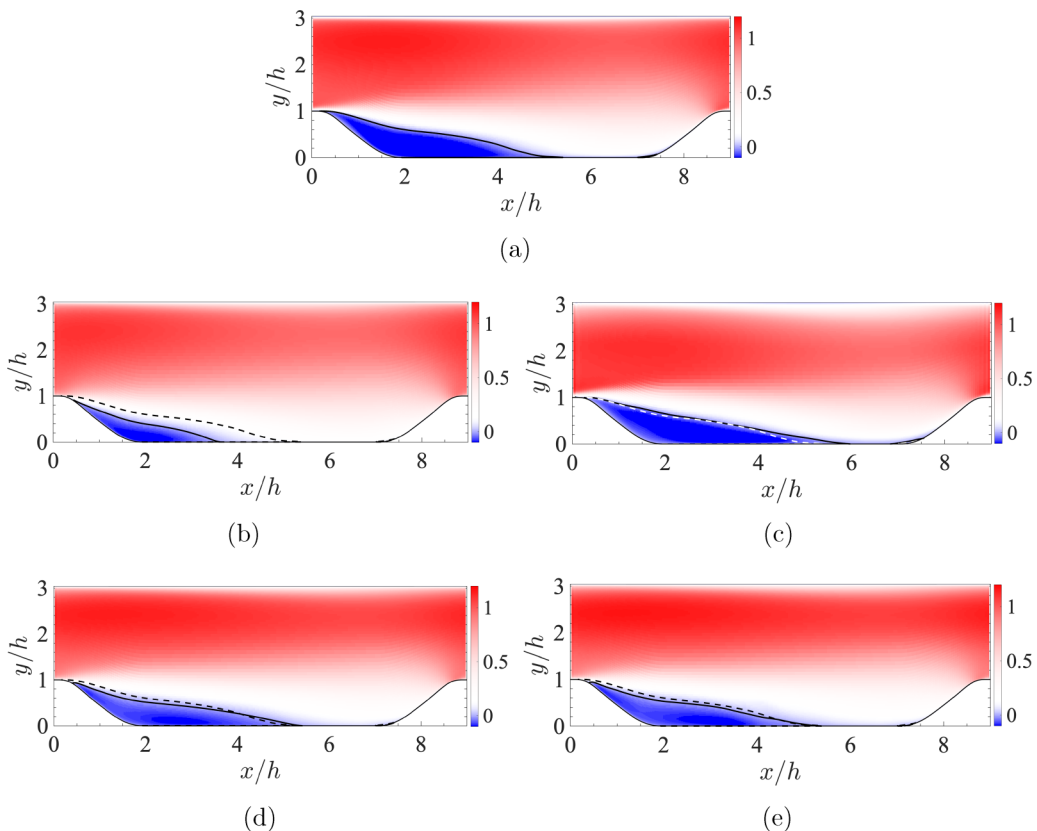


FIG. 9. Forward solutions of the mean, normalized velocity, $\langle u \rangle / u_{\text{bulk}}$, for the standard LEVM model, the two learned models, the lookup table for DNS values of b_{ij} and the DNS results. The solid line represents the region of recirculation, and the dashed line overlays where this region exists in the DNS data. (a) DNS, (b) LEVM, (c) Lookup table for b_{ij} , (d) Learned 1, and (e) Learned 2.

The mean velocity normalized by the bulk velocity, u_{bulk} , is shown in Fig. 9 and the detached regions are delineated by a black line. It is observed that LEVM underpredicts recirculation compared with the DNS results [Fig. 9(b)], while both learned models demonstrate marked qualitative improvement in velocity prediction. Quantitative measurements of separation and reattachment locations for both the primary and secondary recirculation regions are detailed in Table VI. The learned models predict both primary and secondary reattachment points within 1% of the DNS values, with exception of the primary separation point. In comparison, LEVM underpredicts the primary reattachment point by 32% compared with DNS and fails to predict existence of the secondary recirculation.

Examination of the momentum RANS equation [Eq. (4)] makes clear that $\langle u'u' \rangle$ and $\langle u'v' \rangle$ are the only Reynolds stress components that contribute to $\langle u \rangle$ and therefore to the prediction of recirculation. By examining b_{11} and b_{12} in Fig. 10, it can be seen that both components of anisotropy contribute to the prediction of the separation location, however the b_{12} component is most important for the prediction of reattachment. This can be seen by considering the areas of high gradients (since the contribution to the velocity field is in the form of the divergence of the Reynolds stresses), specifically the streamwise gradient for the 11-component and the vertical gradient for the 12-component of the model are important. As shown in Fig. 10, the second basis tensor, $\mathcal{T}_{ij}^{(2)}$, is the most important contribution for accurately describing b_{11} and b_{22} and the first basis tensor,

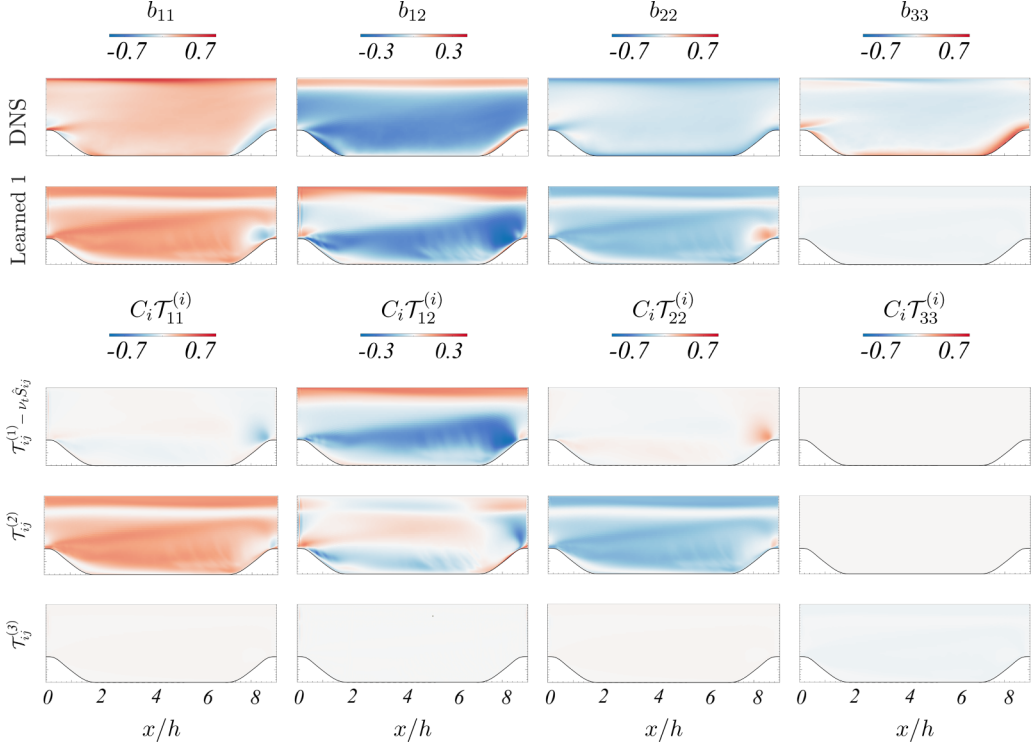


FIG. 10. Contributions to each component of the anisotropy tensor from each of the bases for the model Learned 1, compared with the DNS field.

$\mathcal{T}_{ij}^{(1)}$, is the most dominant contribution for modeling b_{12} . Thus $\partial \mathcal{T}_{11}^{(2)}/\partial x$ and $\partial \mathcal{T}_{12}^{(1)}/\partial y$ are the most important for accurately predicting the recirculation region. Finally, the third basis is critical for accurately describing the b_{33} component, though for this particular configuration (since z is a homogeneous direction), accuracy in this component is not required for predicting the statistically two-dimensional mean flow field.

4. Application outside the scope of training

To assess the range of application of the learned model (Learned 1), a forward RANS simulation was conducted in OpenFOAM for two configurations and Reynolds numbers outside the scope of its training: (1) the wavy wall configuration but at a higher Reynolds number and (2) flow over a backward facing step, where massive separation is also observed. In both of these configurations, DNS and/or experimental data are available to assess model performance. As in the previous section, the learned model's performance is compared against the LEVM.

Since full-field data are not available for these additional cases, model performance is determined based on prediction of the reattachment point in the flow. These results are summarized in Table VI.

The first out-of-scope configuration considered is the periodically constricted channel configuration as described in the previous section, but at a higher Reynolds number. More exhaustive details on the Reynolds number dependencies for this configuration can be found in Breuer *et al.* [51], but to briefly summarize, increasing the Reynolds number for this configuration results in differences in recirculation size as well as in separation and reattachment locations. The selection of $Re = 5600$ was chosen due to the existence of available DNS data. Here the learned model (Learned 1) was again implemented in OpenFOAM and compared against the openly available data set provided

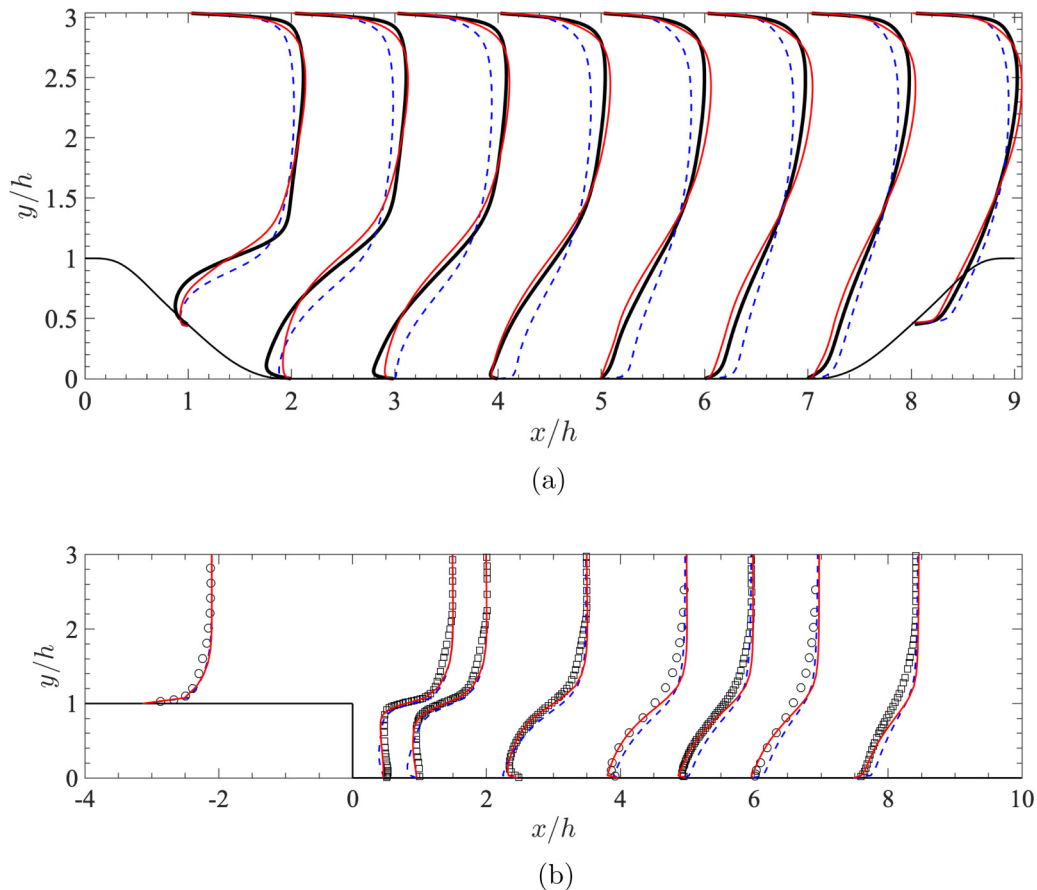


FIG. 11. Velocity profiles for flow through a periodically constricted channel at $Re = 5600$ (a) and flow over a backward-facing step at $Re = 5000$ (b). Learned I (—), LEVM (---), DNS of Ref. [51] (—), DNS of Ref. [59] (□), and experiments of Ref. [60] (○).

by Breuer *et al.* [51], in which only a primary recirculation region is observed. Improvement over LEVM is observed in this case, as shown in Fig. 11(a) and Table VI. The LEVM solution again underpredicts recirculation, while the learned model predicts the reattachment location within 5% of the DNS value reported in Breuer *et al.* [51] and within 5%–9% of the DNS value reported in Krank *et al.* [57]. The primary separation location is marginally over two times further in the streamwise direction as compared to both DNS results [51,57]. The learned model also predicts the small secondary recirculation region that is reported in Krank *et al.* [57]. Breuer *et al.* [51] does not observe this secondary recirculation, however this appears to be due to differences in numerical schemes and order of accuracy as compared with Krank *et al.* [57]. In this secondary region, the learned model predicts the separation and reattachment points within 1% and 0.1%, respectively, as compared with the DNS reported in Krank *et al.* [57].

The second configuration considered to assess model performance outside the scope of its training is turbulent flow over a backward-facing step for $Re = 5000$. For this configuration, model performance is assessed by comparison with reported DNS values [59] and experimental values [60,61] as shown in Fig. 11 and quantified in Table VI. The same configuration is used as described in these works, and the Reynolds number is defined using the step height ($h = 9.6$ mm). No-slip boundary conditions were enforced at the top and bottom walls in the RANS simulations, and a

TABLE VII. Summary of learned coefficients using sparse data, i.e., only the y -dependent data at the specified x/h location. Model error is reported for the *a posteriori* velocity, and separation and reattachment points are compared with DNS, LEVM, and the “Learned 1” model for $Re = 2800$.

Training (x/h)	Learned coefficients			Error $\epsilon^{(u)}$	Primary		Secondary	
	C_1	C_2	C_3		Separation	Reattachment	Separation	Reattachment
1	32.35	45.42	19.69	0.17	0.36 (1.24)	4.15 (0.22)	–	–
2	36.72	46.85	36.43	0.16	0.35 (1.23)	4.27 (0.20)	–	–
3	51.81	48.75	37.59	0.15	0.34 (1.17)	4.78 (0.11)	–	–
4	51.38	53.05	40.99	0.12	0.40 (1.52)	5.32 (0.01)	7.22 (0.05)	7.35 (0.01)
5	48.63	55.83	35.85	0.13	0.41 (1.56)	5.10 (0.04)	7.17 (0.04)	7.18 (0.01)
6	49.85	55.77	23.34	0.13	0.41 (1.56)	5.11 (0.04)	–	–
7	58.34	54.19	–4.04	0.13	0.41 (1.56)	5.14 (0.04)	–	–
8	112.56	51.42	–50.00	0.12	0.37 (1.31)	4.35 (0.18)	7.10 (0.03)	7.27 (0.002)
			Learned 1	0.12	0.40 (1.53)	5.38 (0.005)	7.14 (0.04)	7.31 (0.008)
			LEVM	0.17	0.43 (1.61)	3.64 (0.32)	–	–
			DNS	–	0.16 (–)	5.35 (–)	6.87 (–)	7.25 (–)

fixed velocity condition is enforced at the inlet and a zero gradient condition for the outflow. “Patch” conditions are implemented on the front and back surfaces to enforce a two-dimensional flow.

LEVM is found to underpredict the reattachment point between 23% and 32%, while the model trained on the periodic hill data in previous sections predicts reattachment within 5% as compared to the DNS and experimental values.

It is notable that the level of performance of the learned model outside the scope of its training, particularly in the case of the backward-facing step, is comparable to the out-of-scope performance of the Tensor Basis Neural Network developed by Ling *et al.* [12]. It is also relevant to point out that Ling *et al.* [12] trained on six cases to achieve this level of performance, while the present model was trained on only flow through a periodically constricted channel at $Re = 2800$.

5. Modeling with sparse data

Because the sparse regression methodology uses a physics-based constraint of the basis set and the L_1 penalty acts to regularize the model, far less data are required to achieve reasonable learned models as compared to approaches that evaluate a higher dimensional space of potential parameters, such as NNs. This is demonstrated in two contexts. First, a model is learned using only the y -dependent data along eight streamwise locations (see Table VII). Due to the grid spacing of the RANS simulation, only 160 data points were used for training of each case (compared with 32 000 when using the full data set in the previous section). As seen in Table VII, similar model performance is observed for models trained using data located at $x/h = (4-8)$ when compared with the model learned using the full data set. Secondary recirculation is predicted in three of these training sets. Interestingly, the model trained at $x/h = 8$, where recirculation is not present, is able to predict recirculation in both regions of flow separation. Additionally, it is notable that the model is insensitive to variation in coefficients, especially for the first and third terms.

To further assess the performance of sparse regression in using sparse data, subsets of data are randomly chosen throughout the domain and used as training data (see Fig. 12 for an exemplary subset of training data). Data sets ranging from 30 000 to 50 training points were assessed (see Table VIII). While the learned coefficients change as the data set is reduced, the *a priori* model error in the anisotropic stress tensor increases by only 8%. This suggests that sparse regression would make an excellent modeling construct for extremely sparse data sets, such as those available from experiments where obtaining a high level of resolution is challenging.

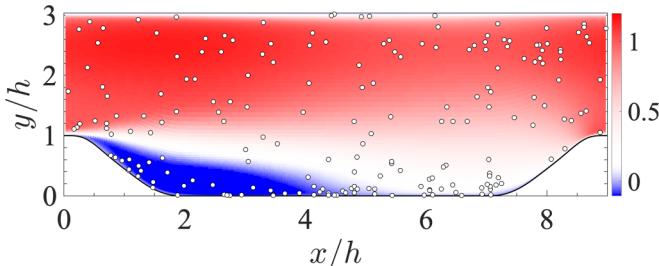


FIG. 12. Example of the random points used for training set corresponding to $n^{\text{train}} = 200$. Location of points used for training (\circ), $\langle u \rangle$ from DNS.

It is notable that the ability of the sparse regression methodology to produce reasonable models using sparse data sets stems primarily from the construction of the optimization itself as well as the number of degrees of freedom associated with the basis chosen. Specifically, since the L_1 penalty regularizes the model, this reduces the data requirement of the algorithm and the physics-based determination of the basis set tends to also reduce the degrees of freedom present in \mathbb{T} . This is particularly reduced in the case of constant coefficients, resulting in the ability to handle sparse sets. In the event that more complex dependencies are required of the coefficients, we find that the amount of data required to learn accurate models increases as $O(1)$ with the number of “trial” functional forms included in \mathbb{T} .

IV. TRAINING THE MODEL WITH EXPERIMENTAL DATA

For many practical systems, procuring highly resolved computational data (i.e., DNS or highly resolved LES) is not feasible. Thus, in these cases modeling efforts are directed toward experimental data, which are inherently both sparse and noisy. To this end, we demonstrate the ability of the sparse regression methodology to successfully model both sparse and noisy data by using the particle image velocimetry (PIV) data available from the Rapp and Manhart [58] experiments for turbulent flow through a periodically constricted channel. This work is the experimental analogy to Breuer *et al.* [51] and uses the same configuration described in Sec. III B.

We consider two cases ($\text{Re} = 5600$ and $10\,600$) for which highly resolved computational data (either DNS or LES) are available [51] in addition to the experimental data [58]. The reason for this is twofold. First, because the PIV measurements do not report $\langle w'w' \rangle$ or k , either of which is required

TABLE VIII. Summary of the learned coefficients for sparse, randomly sampled data using n^{train} training points. The error reported is the *a priori* error in the anisotropic stress tensor.

n^{train}	Coefficients			ϵ^b
	C_1	C_2	C_3	
$n_x \times n_y$	63.12	51.42	10.98	0.64
30 000	62.50	51.52	11.24	0.64
20 000	52.50	45.77	12.19	0.65
10 000	42.03	38.84	16.92	0.67
5000	37.35	36.37	19.69	0.69
1000	33.32	35.41	21.76	0.69
500	33.91	35.07	20.33	0.69
100	33.00	32.5	23.34	0.71
50	31.14	38.28	20.66	0.68

TABLE IX. Summary of the learned model coefficients trained using the three training data sets: full-field DNS at $Re = 2800$, sparse DNS/LES data at $Re = 5600$ and $10\,600$, and experimental data.

Model name	Training set	Learned coefficients		
		C_1	C_2	C_3
Learned 1	DNS (full field), $Re = 2800$	63.12	51.42	10.98
Learned 3	DNS/LES (sparse), $Re = 5600, 10\,600$ [51]	61.72	44.07	12.58
Learned 4	Experiment, $Re = 5600, 10\,600$ [58]	58.56	55.51	118.89

for determining the anisotropy tensor, we rely on using the k value from the computational data to estimate $\langle w'w' \rangle$ and complete the experimental data set. More exactly, we employ the relation, $\langle w'w' \rangle = 2k - \langle u'u' \rangle - \langle v'v' \rangle$, where k is taken from the computational data sets and $\langle u'u' \rangle$ and $\langle v'v' \rangle$ are supplied by the PIV measurements. This allows for the noise of the PIV measurements to be incorporated into the estimate for $\langle w'w' \rangle$. In the event that no highly resolved (e.g., DNS/LES) data are available to reconstruct the full tensor components, then approximations based on known configuration properties, such as symmetry or continuity, could be employed in order to supply missing information. Alternately, sparse regression could be employed on the incomplete Reynolds

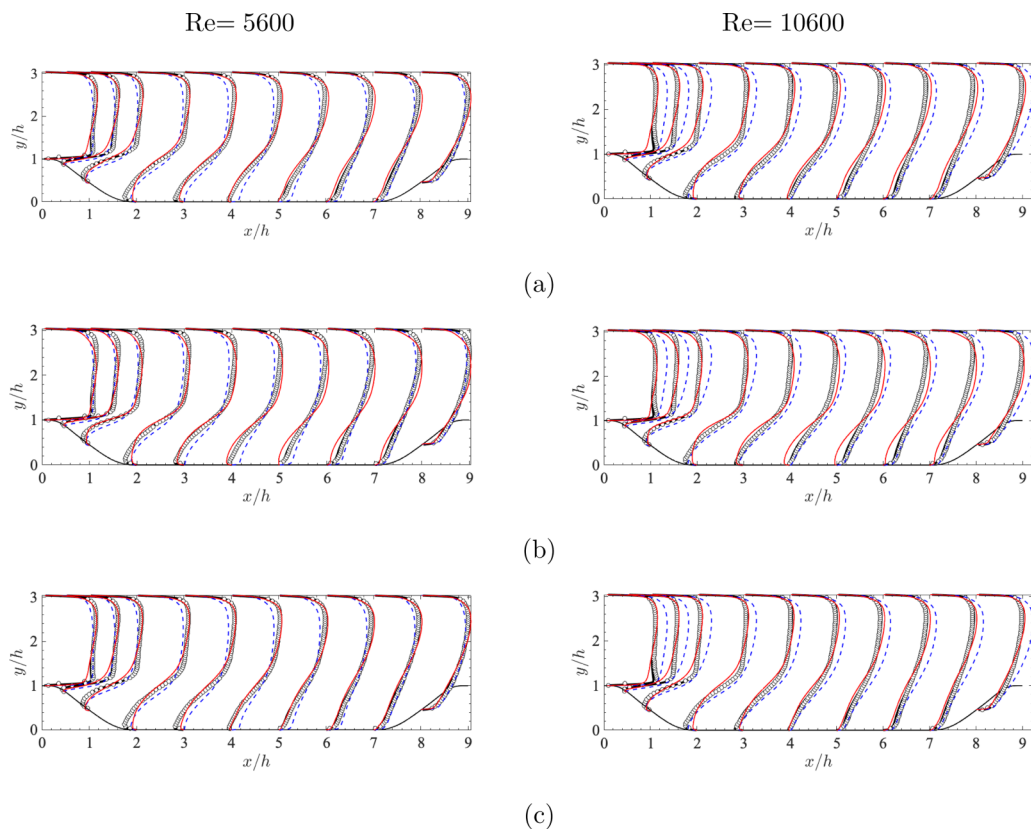


FIG. 13. Velocity profiles for flow through a periodically constricted channel at $Re = 5600$ (left column) and $Re = 10\,600$ (right column). Each plot shows the learned model denoted in the caption (—), LEVM (---), DNS/LES of Ref. [51] (—), experimental values from Ref. [58] (○). (a) Learned 1, (b) Learned 3, and (c) Learned 4.

stress, though this would place constraints on the applicability of the resultant model depending on the importance of the missing data to flow physics.

In addition to completing the experimental data set, the computational data also allow for the more systematic isolation of the effect of noise on modeling. To this end, we first train the model using the computational data (which is reported at 10 streamwise locations), interpolated to the same physical locations as the experimental data to ensure equivalent sparsity in the data set.

Following a similar procedure as was performed in Sec. III B 5, we find that a subset of the data (taken from all 10 streamwise locations and vertical locations corresponding to $y/h \geq 2$ and resulting in a data set of 591 points for $Re = 5600$ and 593 for $Re = 10600$) results in the most accurate model with respect to the *a priori* L_2 training error in b_{ij} . This model is termed “Learned 3” and shown in Table IX. The coefficients are within 15% of the “Learned 1” model, despite using a training set containing only 4% of the points as used in the full-field DNS training. The resulting recirculation predictions for $Re = (2800, 5600, 10600)$ are shown in Table VI and compared against the LEVM and Learned 1 models as well as the available computational and experimental results from the literature.

Next, the experimental data set [58] is used for training (using the same reduced sample locations as in training with the DNS/LES data set). This results in a model with differing coefficients, especially in the case of C_3 , though as reported previously, we observe that the resulting flow field is the least sensitive to $\mathcal{T}^{(3)}$ and thus this difference does not have large implications for the model’s ability to predict recirculation (see Table VI). The resultant mean velocity profiles for all three learned models describe in Table IX are compared against LEVM as well as the experimental and DNS values in Fig. 13. It is evident from these plots that the learned models improve prediction of the reattachment location over LEVM as well as the free flow in the remainder of the domain overall. In some regions, the LEVM outperforms the learned model, however, this is primarily a consequence of using the reattachment location to assess goodness of the learned model, rather than an L_2 norm of the full-field error (since it is not available for the sparse or experimental data). In summary, the sparse regression methodology is capable of providing algebraic closure of the terms appearing in the RANS equations with marked improvement over existing models, even when trained on sparse experimental measurements.

V. CONCLUSION

In this work, a turbulence closure modeling methodology has been proposed as an alternative to other machine learning techniques, such as NNs. This method is based upon sparse regression which uses an L_2 norm with an L_1 norm penalty cost functional to produce a compact, algebraic model. Further, the inputs to the optimization algorithm are specifically tailored in order to ensure form invariance. This is specifically accomplished by arranging the trusted and basis tensorial data into column vectors, thereby constraining coefficients to be invariant with respect to direction. By generating a model in this form, several important modeling properties can be achieved: form (or Galilean) invariance, interpretability, and ease of dissemination. Using two canonical cases, it was demonstrated that this technique produces results with model accuracies similar to that of modern NN methodologies, even when using a drastically reduced training data set.

Using homogeneous free shear turbulence as a preliminary example, sparse regression was able to return the LRR-IP model used to generate a synthetic data set, even when large amounts of noise were applied. Next, using DNS data for homogeneous free shear turbulence, sparse regression learned a model that reduces model error by 70% as compared to the existing LRR-IP and LRR-QI models.

In the case of turbulent flow through a periodically constricted channel, sparse regression uncovered a model that has comparable performance to a modern NN, however this performance can be achieved using a drastically minimal data set and the resultant model form is available in a compact, algebraic form. Interestingly, the resultant model takes on a simple form, with constant coefficients, suggesting that less complexity than is typically postulated in other methodologies

such as NNs is required for the accurate description of flow physics. Additionally, the learned model demonstrated significant improvements in performance as compared with LEVM for a much higher Reynolds number, and outside the scope of its training. Further, due to the ability of sparse regression to learn predictive models using minimal data sets and noisy data (as demonstrated in Sec. III A), it is an ideal candidate for translating experimental data, which may be both noisy and sparse, into accurate models.

Finally, sparse regression assumes complete generality and thus does not strictly require an existing model upon which to augment. This is an important property for other open areas of research, e.g., modeling multiphase turbulence [19–23], for which existing models are either unavailable or too inaccurate to reliably use as a baseline model upon which to build. Such an approach can also be applied to turbulent combustion, in which heat release due to chemical reactions give rise to “back scatter” and existing models based on an energy cascade fail to be predictive [24,25].

In future work, integration of gene expression programming (GEP) with sparse regression may be beneficial for cases in which complex algebraic dependencies upon the principal invariants become necessary (i.e., coefficients that are either constant or have simple dependencies on the invariants do not reduce model error sufficiently). In this event, sparse regression would be employed to determine the most important basis tensors, and then GEP could be used to determine the functional dependence of coefficients on the principal invariants.

ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship. We would also like to acknowledge the National Science Foundation for partial support from award CBET 1846054. The computing resources and assistance provided by the staff of the Advanced Research Computing at the University of Michigan, Ann Arbor, are greatly appreciated. Finally, the authors gratefully acknowledge Prof. M. Houssem Kasbaoui for the code used to generate homogeneous sheared turbulence data and Prof. C. Petty for his discussions on noninertial reference frames.

-
- [1] F. R. Menter, Two-equation eddy-viscosity turbulence models for engineering applications, *AIAA J.* **32**, 1598 (1994).
 - [2] F. R. Menter, M. Kuntz, and R. Langtry, Ten years of industrial experience with the SST turbulence model, *Turbul. Heat Mass Transfer* **4**, 625 (2003).
 - [3] O. Reynolds, On the dynamical theory of incompressible viscous fluids and the determination of the criterion, *Philos. Trans. R. Soc. London A* **186**, 123 (1895).
 - [4] H. Tennekes and J. L. Lumley, *A First Course in Turbulence* (MIT Press, Cambridge, MA, 1992).
 - [5] P. Moin and J. Kim, Tackling turbulence with supercomputers, *Sci. Am.* **276**, 62 (1997);
 - [6] R. H. Bush, T. Chyczewski, K. Duraisamy, B. Eisfeld, C. L. Rumsey, and B. R. Smith, Recommendations for future efforts in RANS modeling and simulation, in *AIAA SchiTech 2019 Forum, 7-11 January 2019, San Diego, California* (AIAA, Reston, VA, 2019), Paper No. AIAA 2019-0317.
 - [7] J. Slotnick, A. Khodadoust, J. Alonso, D. Darmofal, W. Gropp, E. Lurie, and D. Mavriplis, CFD vision 2030 study: A path to revolutionary computational aerosciences, technical report for Langley Research Center under contract NNL08AA16B (2014).
 - [8] M. P. Brenner, J. D. Eldredge, and J. B. Freund, Perspective on machine learning for advancing fluid mechanics, *Phys. Rev. Fluids* **4**, 100501 (2019).
 - [9] K. Duraisamy, G. Iaccarino, and H. Xiao, Turbulence modeling in the age of data, *Annu. Rev. Fluid Mech.* **51**, 357 (2019).

-
- [10] J. R. Holland, J. D. Baeder, and K. Duraisamy, Towards integrated field inversion and machine learning with embedded neural networks for RANS modeling, in *AIAA Scitech 2019 Forum, 7-11 January 2019, San Diego, California* (AIAA, Reston, VA, 2019), Paper No. AIAA 2019-1884.
- [11] K. Duraisamy, Z. J. Zhang, and A. P. Singh, New approaches in turbulence and transition modeling using data-driven techniques, in *Proceedings of the 53rd AIAA Aerospace Sciences Meeting, 5-9 January 2015, Kissimmee, Florida* (AIAA, Reston, VA, 2015), pp. 1–14.
- [12] J. Ling, A. Kurzawski, and J. Templeton, Reynolds averaged turbulence modeling using deep neural networks with embedded invariance, *J. Fluid Mech.* **807**, 155 (2016).
- [13] C. L. Rumsey, Successes and challenges for flow control simulations (invited), in *AIAA, 4th Flow Control Conference, 23–26 June 2008, Seattle, WA* (American Institute of Aeronautics and Astronautics, Inc., Seattle, Washington, 2008), pp. 1–26.
- [14] F. Köhler, J. Munz, and M. Schäfer, Data-driven augmentation of RANS turbulence models for improved prediction of separation in wall-bounded flows, in *AIAA SchiTech 2020 Forum, 6-10 January 2020 Orlando, FL* (AIAA, Reston, VA, 2020), Paper No. AIAA 2020-1586.
- [15] B. Parmar, E. Peters, K. E. Jansen, A. Doostan, and J. A. Evans, Generalized non-linear eddy viscosity models for data-assisted Reynolds stress closure, in *AIAA SchiTech 2020 Forum, 6-10 January 2020 Orlando, FL* (AIAA, Reston, VA, 2020), Paper No. AIAA 2020-0351.
- [16] E. Rajabi and M. R. Kavianpour, Intelligent prediction of turbulent flow over backward-facing step using direct numerical simulation data, *Eng. Appl. Comput. Fluid Mech.* **6**, 490 (2012).
- [17] L. Weishuo and F. Jian, Iterative framework of machine-learning based turbulence modeling for Reynolds-averaged Navier-Stokes simulations, [arXiv:1910.01232](https://arxiv.org/abs/1910.01232) [physics.flu-dyn].
- [18] H. Xiao, J.-L. Wu, S. Laizet, and L. Duan, Flows over periodic hills of parameteri geometries: A dataset for data-driven turbulence modeling from direct simulations, *Comput. Fluids* **200**, 104431 (2020).
- [19] S. Beetham and J. Capecelatro, Biomass pyrolysis in fully-developed turbulent riser flow, *Renewable Energy* **140**, 751 (2019).
- [20] J. Capecelatro, O. Desjardins, and R. O. Fox, Strongly-coupled gas-particle flows in vertical channels. Part II: Turbulence modeling, *Phys. Fluids* **28**, 1 (2006).
- [21] J. Capecelatro, O. Desjardins, and R. O. Fox, Numerical study of collisional particle dynamics in cluster-induced turbulence, *J. Fluid Mech.* **747**, R2 (2014).
- [22] J. Capecelatro, O. Desjardins, and R. O. Fox, On fluid-particle dynamics in fully developed cluster-induced turbulence, *J. Fluid Mech.* **780**, 578 (2015).
- [23] R. O. Fox, On multiphase turbulence models for collisional fluid–particle flows, *J. Fluid Mech.* **742**, 368 (2014).
- [24] H. Pitsch, Large-eddy simulation of turbulent combustion, *Annu. Rev. Fluid Mech.* **38**, 453 (2006).
- [25] D. Veynante and L. Vervisch, Turbulent combustion modeling, *Prog. Energy Combust. Sci.* **28**, 193 (2002).
- [26] C. Lav, R. D. Sandberg, and J. Philip, A framework to develop data-driven turbulence models for flows with organized unsteadiness, *J. Comput. Phys.* **383**, 148 (2019).
- [27] M. Schmelzer, R. Dwight, and P. Cinnella, Data-driven deterministic symbolic regression of nonlinear stress-strain relation for RANS turbulence modelling, *2018 Fluid Dynamics Conference, Atlanta, GA, June 25–29* (AIAA, Reston, VA, 2018), Paper No. 2018-2900.
- [28] M. Schmelzer, R. Dwight, and P. Cinnella, Discovery of algebraic Reynolds-stress models using sparse symbolic regression, *Flow Turbul. Combust.* **104**, 579 (2020).
- [29] J. Weatheritt and R. D. Sandberg, Improved junction body flow modeling through data-driven symbolic regression, *J. Ship Res.* **63**, 283 (2019).
- [30] Y. Zhao, H. D. Akolekar, J. Weatheritt, V. Michelassi, and R. D. Sandberg, Turbulence model development using CFD-driven machine learning, *J. Comput. Phys.* **411**, 109413 (2020).
- [31] J.-X. Wang, J. Wu, and H. Xiao, Physics informed machine learning approach for reconstructing Reynolds stress modeling discrepancies based on DNS data, *Phys. Rev. Fluids* **2**, 1 (2017).
- [32] J. Wu, H. Xiao, and E. Paterson, Physics-informed machine learning for augmenting turbulence models: A comprehensive framework, *Phys. Rev. Fluids* **3**, 1 (2018).

- [33] S. L. Brunton, J. L. Proctor, and J. N. Kutz, Discovering governing equations from data by sparse identification of nonlinear dynamical systems, *Proc. Natl. Acad. Sci. USA* **113**, 3932 (2016).
- [34] C. M. Bishop, *Pattern Recognition and Machine Learning* (Springer Science+Business Media, NY, 2006).
- [35] R. Tibshirani, Regression shrinkage and selection via the LASSO, *J. R. Stat. Soc. B* **58**, 267 (1996).
- [36] H. Zou and T. Hastie, Regularization and variable selection via the elastic net, *J. R. Stat. Soc. B* **67**, 301 (2005).
- [37] S. B. Pope, A more general effective-viscosity hypothesis, *J. Fluid Mech.* **72**, 331 (1975).
- [38] C. G. Speziale, S. Sarkar, and T. B. Gatski, Modelling the pressure-strain correlation of turbulence: An invariant dynamical systems approach, *J. Fluid Mech.* **227**, 245 (1991).
- [39] K. Hanjalic and B. E. Launder, A Reynolds stress model of turbulence and its application to thin shear flows, *J. Fluid Mech.* **52**, 609 (1972).
- [40] B. E. Launder, Phenomenological modeling: Present . . . and future? in *Whither Turbulence? Turbulence at the Crossroads*, edited by J. L. Lumley (Springer, NY, 1990), pp. 439–485.
- [41] B. E. Launder, G. J. Reese, and W. Rodi, Progress in the development of a Reynolds-stress turbulence closure, *J. Fluid Mech.* **68**, 537 (1975).
- [42] T. B. Gatski and C. G. Speziale, On explicit algebraic stress models for complex turbulent flows, *J. Fluid Mech.* **254**, 59 (1993).
- [43] O. Desjardins, G. Blanquart, G. Balarac, and H. Pitsch, High order conservative finite difference scheme for variable density low Mach number turbulent flows, *J. Comput. Phys.* **227**, 7125 (2008).
- [44] C. D. Pierce, Progress-variable approach for large-eddy simulation of turbulent combustion, Ph.D. thesis, Stanford University (2001).
- [45] M. H. Kasbaoui, R. Patel, D. Koch, and O. Desjardins, An algorithm for solving the NavierStokes equations with shear-periodic boundary conditions and its application to homogeneously sheared turbulence, *J. Fluid Mech.* **833**, 687 (2017).
- [46] T. Passot and A. Pouquet, Numerical simulation of compressible homogeneous flows in the turbulent regime, *J. Fluid Mech.* **181**, 441 (1986).
- [47] K. R. Sreenivasan, On the universality of the Kolmogorov constant, *Phys. Fluids* **7**, 2778 (1995).
- [48] J. C. Rotta, Statistische Theorie nichthomogener Turbulenz, *Physics* **129**, 547 (1951).
- [49] T. Jongen, G. Mompean, and T. B. Gatski, Accounting for Reynolds stress and dissipation rate anisotropies in inertial and noninertial frames, *Phys. Fluids* **10**, 674 (1998).
- [50] C. G. Speziale, Turbulence modeling in noninertial frames of reference, *Theor. Comput. Fluid Dyn.* **1**, 3 (1989).
- [51] M. Breuer, N. Peller, C. Rapp, and M. Manhart, Flow over periodic hills—Numerical and experimental study in a wide range of Reynolds numbers, *Comput. Fluids* **38**, 433 (2009).
- [52] S. B. Pope, *Turbulent Flows* (Cambridge University Press, Cambridge, 2000).
- [53] J. Capecehatro and O. Desjardins, An Euler–Lagrange strategy for simulating particle-laden flows, *J. Comput. Phys.* **238**, 1 (2013).
- [54] J. Wu, H. Xiao, R. Sun, and Q. Wang, Reynolds-averaged Navier-Stokes equations with explicit data-driven Reynolds stress closure can be ill-conditioned, *J. Fluid Mech.* **869**, 553 (2019).
- [55] H. G. Weller, G. Tabor, H. Jasak, and C. Fureby, A tensorial approach to computational continuum mechanics using object-oriented techniques, *Comput. Phys.* **12**, 620 (1998).
- [56] T. H. Shih, J. Zhu, and J. Lumley, A realizable Reynolds stress algebraic equation model, NASA Technical Memorandum no. 105993 (1993).
- [57] B. Krank, M. Kronbichler, and W. A. Wall, Direct numerical simulation of flow over periodic hills up to $Re_H = 10$, 595, *Flow Turbul. Combust.* **101**, 521 (2018).
- [58] Ch. Rapp and M. Manhart, Flow over periodic hills: An experimental study, *Exp. Fluids* **51**, 247 (2011).
- [59] H. Le, P. Moin, and J. Kim, Direct numerical simulation of turbulent flow over a backward-facing step, *J. Fluid Mech.* **330**, 349 (1997).
- [60] S. Jovic and D. M. Driver, Backward-facing step measurements at low Reynolds number, $Re_h = 5000$, NASA Technical Memorandum no. 108807 (1994).
- [61] N. Kasagi and A. Matsunaga, Three-dimensional particle-tracking velocimetry measurement of turbulence statistics and energy budget in a backward-facing step flow, *Intl. J. Heat Fluid Flow* **16**, 477 (1995).