

# Physics-informed machine learning approach for reconstructing Reynolds stress modeling discrepancies based on DNS data

Jian-Xun Wang, Jin-Long Wu, and Heng Xiao\*

*Department of Aerospace and Ocean Engineering, Virginia Tech, Blacksburg, Virginia 24060, USA*

(Received 6 July 2016; published 16 March 2017)

Turbulence modeling is a critical component in numerical simulations of industrial flows based on Reynolds-averaged Navier-Stokes (RANS) equations. However, after decades of efforts in the turbulence modeling community, universally applicable RANS models with predictive capabilities are still lacking. Large discrepancies in the RANS-modeled Reynolds stresses are the main source that limits the predictive accuracy of RANS models. Identifying these discrepancies is of significance to possibly improve the RANS modeling. In this work, we propose a data-driven, physics-informed machine learning approach for reconstructing discrepancies in RANS modeled Reynolds stresses. The discrepancies are formulated as functions of the mean flow features. By using a modern machine learning technique based on random forests, the discrepancy functions are trained by existing direct numerical simulation (DNS) databases and then used to predict Reynolds stress discrepancies in different flows where data are not available. The proposed method is evaluated by two classes of flows: (1) fully developed turbulent flows in a square duct at various Reynolds numbers and (2) flows with massive separations. In separated flows, two training flow scenarios of increasing difficulties are considered: (1) the flow in the same periodic hills geometry yet at a lower Reynolds number and (2) the flow in a different hill geometry with a similar recirculation zone. Excellent predictive performances were observed in both scenarios, demonstrating the merits of the proposed method.

DOI: [10.1103/PhysRevFluids.2.034603](https://doi.org/10.1103/PhysRevFluids.2.034603)

## I. INTRODUCTION

### A. RANS models as workhorse tool in industrial CFD

Computational fluid dynamics (CFD) simulations have been widely used in aerospace, mechanical, and chemical industries to support engineering design, analysis, and optimization. Two decades ago when large eddy simulations (LES) started gaining popularity with the increasing availability of computational resources, it was widely expected that LES would gradually displace and eventually replace Reynolds-averaged Navier-Stokes (RANS) equations in industrial CFD workflows for decades to come. In the past two decades, however, while LES-based methods (including resolved LES, wall-modeled LES, and hybrid LES/RANS methods) did gain widespread applications, and the earlier hope certainly did not diminish, the predicted time when these methods would replace RANS has been significantly delayed. This observation is particularly relevant in light of the recent discussions on the ending of the “Moore’s law era” with transistor sizes approaching their theoretical lower limit [1,2]. RANS solvers, particularly those based on standard eddy viscosity models (e.g.,  $k-\varepsilon$  [3],  $k-\omega$  [4,5], S-A [6], and SST  $k-\omega$  [7]), are still and will remain the dominant tool for industrial CFD in the near future. This is likely to be true even in mission critical applications such as aircraft design. Interestingly, even the advanced RANS models such as Reynolds stress transport models [8] and explicit algebraic Reynolds stress models [9] have not seen much development in the past few decades. These advanced models are computationally more expensive and less robust compared to the standard eddy viscosity RANS models. As such, it is still practically important to

---

\*Corresponding author: [hengxiao@vt.edu](mailto:hengxiao@vt.edu)

further develop the standard RANS models for industrial CFD applications. However, improving the predictive capabilities of these models is critical yet technically challenging.

### B. Progress and challenges in data-driven turbulence modeling

While traditional development of turbulence models has focused on incorporating more physics to improve predictive capabilities, an alternative approach is to utilize data. In the past few years, a number of data-driven approaches have been proposed. Researchers have investigated the use of both offline data (i.e., existing direct numerical simulation (DNS) data for flows different from that to be predicted [10–12]) and online data (streamed monitoring data from the flow to be predicted [13–15]). Dow and Wang [10] used DNS data from a plane channel flow to infer the full-field discrepancy in the turbulent viscosity  $\nu_t$  modeled by the  $k$ - $\omega$  model. To predict flows in channels with wavy boundaries, they modeled the (logarithmic) discrepancies of  $\nu_t$  in the new flows as Gaussian random fields, with the discrepancy field inferred above as mean. Duraisamy and co-workers [11,16] introduced a full-field multiplicative discrepancy term  $\beta$  into the production term of the transport equations of turbulent quantities (e.g.,  $\tilde{\nu}_t$  in the SA model and  $\omega$  in the  $k$ - $\omega$  models). They used DNS data to calibrate and infer uncertainties in the  $\beta$  term. It is expected that the inferred discrepancy field can provide valuable insights to the development of turbulence model and can be used to improve RANS predictions in similar flows. Xiao *et al.* [13] used sparse velocity measurements (online data) to infer the full-field discrepancies  $\Delta\tau_\alpha$  in the RANS-predicted Reynolds stress tensors, or more precisely the physical projections thereof (turbulent kinetic energy, anisotropy, and orientations). Throughout this paper it is understood that  $\tau_\alpha$  indicates the physical projections and not the individual components of the Reynolds stress tensor. Good performance was demonstrated on several canonical flows including flow past periodic hills and flow in a square duct [13].

All three approaches [10,11,13] discussed above can be considered starting points toward the same destination: the capability of predictive turbulence modeling by using standard RANS models in conjunction with offline data. To this end, the respective discrepancies terms ( $\Delta \log \nu_t$ ,  $\beta$ , and  $\Delta\tau_\alpha$ ) are expected to be extrapolated to similar yet different flows. These contributions are all relatively recent, and much of the research is still on-going. Duraisamy *et al.* [16,17] performed *a priori* studies to show the potential universality of their discrepancy term  $\beta$  among a class of similar flows, but their performances in *a posteriori* tests, i.e., using the calibrated discrepancy in one flow to predict another flow, have yet to be demonstrated. Dow and Wang [10] extrapolated the logarithmic discrepancies  $\Delta \log \nu_t$  calibrated in the plane channel flow to flows in channels with slightly wavy walls, where velocity predictions were made. Similarly, further pursuing the approach of Xiao *et al.* [13], Wu *et al.* [18] showed that the Reynolds stress discrepancy calibrated with sparse velocity data can be extrapolated to flows at Reynolds number more than an order of magnitude higher than that in the calibration case. The extrapolated discrepancy has led to markedly improved predictions of velocities and other quantities of interest (QoIs), showing the potential of the approach in enabling data-driven predictive turbulence modeling. However, an intrinsic limitation in the approach of Wu *et al.* [18] is that they inferred the functions  $f_\alpha^{(x)} : \mathbf{x} \mapsto \Delta\tau_\alpha$ , or simply denoted as  $\Delta\tau_\alpha(\mathbf{x})$ , in the space of *physical coordinates*  $\mathbf{x}$ . Therefore, strictly speaking they only demonstrated that the discrepancy  $\Delta\tau_\alpha$  can be extrapolated to flows in the same geometry at the same location. Consequently, their attempts of extrapolation to the flow in a different geometry (e.g., from a square duct to a rectangular duct) encountered less success. The approach of Dow and Wang [10] would share the same limitation since they built Gaussian random fields indexed by the physical coordinates  $\mathbf{x}$ .

### C. Motivation of the proposed approach

A natural extension that overcomes the key limitation in the calibration-prediction approach of Wu *et al.* [18] is to build such functions in a space of well-chosen features  $\mathbf{q}$  instead of physical coordinates  $\mathbf{x}$ . Despite its limitations, a key factor in the success of the original approach is that the

Reynolds stress discrepancies are formulated on its projections, such as the anisotropy parameters ( $\xi$  and  $\eta$ ) and orientation ( $\varphi_i$ ) of the Reynolds stresses, and not directly on the individual components. These projections are normalized quantities [18]. We will retain these merits in the current approach and thus use data to construct functions  $\Delta\tau_\alpha(\mathbf{q})$  instead of  $\Delta\tau_\alpha(\mathbf{x})$ . This extension would allow the calibrated discrepancies to be extrapolated to a much wide range of flows. In other words, the discrepancies of the RANS-predicted Reynolds stresses can be quantitatively explained by the mean flow physics. Hence, these discrepancies are likely to be universal quantities that can be extrapolated from one flow to another, at least among different flows sharing the same characteristics (e.g., separation). As such, discrepancies in Reynolds stress projections are suitable targets to build functions for.

With the function targets identified, two challenges remain: (1) to identify a set of mean flow features based on which the discrepancies functions  $\Delta\tau_\alpha(\mathbf{q})$  can be constructed and (2) to choose a suitable method for constructing such functions. Duraisamy and co-workers [17] identified several features and used a neural network to construct functions for the multiplicative discrepancy term. Ling and Templeton [12] provided a richer and much more complete set of features in their pioneering work, and they evaluated several machine learning algorithms to predict point-based binary confidence indicators of RANS models [12]. Ling *et al.* [19] further used machine learning techniques to predict the Reynolds stress anisotropy in jets in crossflow. Based on the success demonstrated by Ling and co-workers [12,19], we will use machine learning to construct the functions  $\Delta\tau_\alpha(\mathbf{q})$  in the current work. Specifically, we will examine a class of supervised machine learning techniques, where the objective of the learning is to build a statistical model from data and to make predictions on a response based on one or more inputs [20]. This is in contrast to unsupervised learning, where no response is used in the training or prediction, and the objective is to understand the relationship and structure of the input data. Unsupervised learning will be explored as an alternative approach in future works.

#### D. Objective, scope, and vision of this work

The objective of this contribution is to present an approach to predict Reynolds stress modeling discrepancies in new flows by utilizing data from flows with similar characteristics as the prediction flow. This is achieved by training regression functions of Reynolds stress discrepancies with the DNS database from the training flows.

In light of the consensus in the turbulence modeling community that the Reynolds stresses are the main source of model-form uncertainty in RANS simulations [5,21,22], the current work aims to improve the RANS modeled Reynolds stresses. In multiphysics applications the QoIs might well be the Reynolds stresses and/or quantities that directly depend thereon. In these applications the current work is significant by itself in that it would enable the use of standard RANS models in conjunction with an offline database to provide accurate Reynolds stress predictions. Moreover, the improvement of Reynolds stresses enabled by the proposed method is an important step towards a data-driven turbulence modeling framework. However, the Reynolds stresses corrected by the constructed discrepancy function from DNS databases cannot necessarily guarantee obtaining improved mean flow fields. There are a number of challenges associated with propagating the improvement of Reynolds stresses through RANS equation to the mean velocity field, which will be addressed in future works.

The rest of this paper is organized as follows. Section II introduces the components of the predictive framework, including the choice of regression inputs and responses as well as the machine learning technique used to build the regression function. Section III shows the numerical results to demonstrate the merits of the proposed method. Further interpretation of the feature importance and its implications to turbulence model development are discussed in Sec. IV. Finally, Sec. V concludes the paper.

## II. METHODOLOGY

### A. Problem statement

The overarching goal of the current and companion works is a physics-informed machine learning (PIML) framework for predictive turbulence modeling. Here, “physics-informed” is to emphasize the attempt to account for the physical domain knowledge in every stage of machine learning. The problem targeted by the PIML framework can be formulated as follows: given high-fidelity data (e.g., Reynolds stresses from DNS or resolved LES) from a set  $\{\mathcal{T}_i\}_{i=1}^N$  of  $N$  training flows, the framework shall allow for using standard RANS turbulence models to predict a new flow  $\mathcal{P}$  for which data are not available. The flows  $\mathcal{T}_i$  for which high-fidelity simulation data are available are referred to as *training flows*, and the flow  $\mathcal{P}$  to be predicted is referred to as *test flow*. The lack of data in test flows is typical in industrial CFD simulations performed to support design and optimization. Furthermore, we assume that the training flows and the test flow have similar complexities and are dominated by the same characteristics, such as separation or shock–boundary layer interaction. This scenario is common in the engineering design process, where the test flow is closely related to the training flows. Ultimately, the envisioned machine learning framework will be used in scenarios where the training flows consist of a wide range of elementary and complex flows with various characteristics and the test flow has a subset or all of them. However, the latter scenario is much more challenging and is outside the scope of the current study. Considering that the proposed method is a completely new paradigm, we decided to take small steps by starting from the closely related flows and to achieve the overarching goal gradually.

### B. Summary of proposed approach

In the proposed approach we utilize training data to construct functions of the discrepancies (compared to the DNS data) in the RANS-predicted Reynolds stresses and use these functions to predict Reynolds stresses in new flows. The procedure is summarized as follows:

- (1) Perform baseline RANS simulations on both the training flows and the test flow.
- (2) Compute the feature vector field  $\mathbf{q}(\mathbf{x})$ , e.g., pressure gradient and streamline curvature, based on the RANS-predicted mean flow fields for all flows.
- (3) Compute the discrepancies field  $\Delta\tau_\alpha(\mathbf{x})$  in the RANS modeled Reynolds stresses for the training flows based on the high-fidelity data.
- (4) Construct regression functions  $f_\alpha : \mathbf{q} \mapsto \Delta\tau_\alpha$  for the discrepancies based on the training data prepared in Step 3.
- (5) Compute the Reynolds stress discrepancies for the test flow by querying the regression functions. The Reynolds stresses can subsequently be obtained by correcting the baseline RANS predictions with the evaluated discrepancies.

In machine learning terminology the discrepancies  $\Delta\tau_\alpha$  here are referred to as *responses* or *targets*, the feature vector  $\mathbf{q}$  as *input*, and the mappings  $f_\alpha : \mathbf{q} \mapsto \Delta\tau_\alpha$  as *regression functions*. A regression function  $f_\alpha$  maps the input feature vector  $\mathbf{q}$  to the response  $\Delta\tau_\alpha$ , and the term “function” shall be interpreted in a broad sense here. That is, depending on the regression technique used, it can be either deterministic (e.g., for linear regression) or random (e.g., Gaussian process) [20,23], and it may not even have an explicit form. In the case of random forests regression used in this work [24], the mapping does not have an explicit expression but is determined based on a number of decision trees.

In the procedure described above, after the baseline RANS simulations in Step 1, the input feature fields are computed in Step 2, the training data are prepared in Step 3, and the regression functions are constructed in Step 4. Finally, the regression functions are evaluated to make predictions in Step 5. It is worth noting that in each stage domain knowledge is incorporated, e.g., physical reasoning for identification of input features and consideration of realizability constraints of Reynolds stress in learning-prediction process. Each component is discussed in detail below. The choice of features

TABLE I. Nondimensional flow features used as input in the regression. The normalized feature  $q_\beta$  is obtained by normalizing the corresponding raw features value  $\hat{q}_\beta$  with normalization factor  $q_\beta^*$  according to  $q_\beta = \hat{q}_\beta / (|\hat{q}_\beta| + |q_\beta^*|)$  except for  $\beta = 3$ . Repeated indices imply summation for indices  $i, j, k$ , and  $l$  but not for  $\beta$ . Notations are as follows:  $U_i$  is mean velocity,  $k$  is turbulent kinetic energy (TKE),  $u'_i$  is fluctuation velocity,  $\rho$  is fluid density,  $\varepsilon$  is the turbulence dissipation rate,  $\mathbf{S}$  is the strain rate tensor,  $\boldsymbol{\Omega}$  is the rotation rate tensor,  $\nu$  is fluid viscosity,  $d$  is distance to wall,  $\boldsymbol{\Gamma}$  is unit tangential velocity vector,  $D$  denotes material derivative, and  $L_c$  is the characteristic length scale of the mean flow.  $\|\cdot\|$  and  $|\cdot|$  indicate matrix and vector norms, respectively.

Feature ( $q_\beta$ )	Description	Raw feature ( $\hat{q}_\beta$ )	Normalization factor ( $q_\beta^*$ )
$q_1$	Ratio of excess rotation rate to strain rate ( $Q$ criterion)	$\frac{1}{2}(\ \boldsymbol{\Omega}\ ^2 - \ \mathbf{S}\ ^2)$	$\ \mathbf{S}\ ^2$
$q_2$	Turbulence intensity	$k$	$\frac{1}{2}U_i U_i$
$q_3$	Wall-distance based Reynolds number	$\min(\frac{\sqrt{k}d}{50\nu}, 2)$	not applicable <sup>a</sup>
$q_4$	Pressure gradient along streamline	$U_k \frac{\partial P}{\partial x_k}$	$\sqrt{\frac{\partial P}{\partial x_j} \frac{\partial P}{\partial x_j} U_i U_i}$
$q_5$	Ratio of turbulent time scale to mean strain time scale	$\frac{k}{\varepsilon}$	$\frac{1}{\ \mathbf{S}\ }$
$q_6$	Cratio of pressure normal stresses to shear stresses	$\sqrt{\frac{\partial P}{\partial x_i} \frac{\partial P}{\partial x_i}}$	$\frac{1}{2}\rho \frac{\partial U_k^2}{\partial x_k}$
$q_7$	Nonorthogonality between velocity and its gradient [28]	$ U_i U_j \frac{\partial U_i}{\partial x_j} $	$\sqrt{U_l U_l U_i \frac{\partial U_l}{\partial x_j} U_k \frac{\partial U_k}{\partial x_j}}$
$q_8$	Ratio of convection to production of TKE	$U_i \frac{dk}{dx_i}$	$ \overline{u'_j u'_k} S_{jk} $
$q_9$	Ratio of total to normal Reynolds stresses	$\ \overline{u'_i u'_j}\ $	$k$
$q_{10}$	Streamline curvature	$ \frac{D\boldsymbol{\Gamma}}{Ds} $ where $\boldsymbol{\Gamma} \equiv \mathbf{U}/ \mathbf{U} $ , $Ds =  \mathbf{U} Dt$	$\frac{1}{L_c}$

<sup>a</sup>Normalization is not necessary as the Reynolds number is nondimensional.

and responses are presented in Secs. II C and II D, respectively, and the machine learning algorithm chosen to build the regression function is introduced in Sec. II E.

### C. Choice of mean flow features as regression input

As has been pointed out in Sec. I, mean flow features are better suited as input of the regression function than physical coordinates as they allow the constructed functions to predict flows in different geometries. Ling and Templeton [12] proposed a rich set of twelve features based on clear physical reasoning. The set of features used in the present study mostly follow their work, except that we excluded the feature ‘‘vortex stretching’’ (input 8 in Table II of Ref. [12]). This feature is present only in three-dimensional flows, but the test cases presented here are two-dimensional flow. We excluded two additional features related to linear and nonlinear eddy viscosities (features 6 and 12 in Ref. [12]). These quantities were specifically chosen for evaluating qualitative confidence indicators of RANS predictions and, in our opinion, are not suitable input for regression functions of Reynolds discrepancies. Finally, experiences in the turbulence modeling communities suggest that mean streamline curvature has important influences on the predictive performance of RANS models [25]. Therefore, curvature is included as an additional feature. The complete list of the mean flow features chosen as regression inputs in this work is summarized in Table I.

In choosing the mean flow features as regression inputs, we have observed a few principles in general. First, the input and thus the obtained regression functions should be Galilean invariant. Quantities that satisfy this requirement include all scalars and the invariants (e.g., norms) of vectors and tensors. An interesting example is the pressure gradient along streamlines (see feature  $q_4$  in Table I). While neither velocity  $U_k$  nor pressure gradient  $dP/dx_k$  (both being vectors) is Galilean-invariant by itself and thus is not a suitable input, their inner product  $U_k \frac{dP}{dx_k}$  is. Second, since the truth

of the mean flow fields in the test flows are not available, an input should solely utilize information of the mean flow field produced by the RANS simulations [12]. Therefore, all of the adopted features are based on RANS-predicted pressure  $P$ , velocity  $\mathbf{U}$ , turbulent kinetic energy  $k$ , and distance  $d$  to the nearest wall. Finally, to facilitate implementation and avoid ambiguity, only local quantities (i.e., cell- or point-based quantities in CFD solvers) of the flow field are used in the formulation of features, with the distance  $d$  to nearest wall being a notable exception. This principle is similar to that in choosing variables for developing turbulence models [25].

The interpretation of most feature variables are evident from the brief descriptions given in the table, but a few need further discussions. First, feature  $q_1$  ( $Q$  criterion) is based on the positive second invariant  $Q$  of the mean velocity gradient  $\nabla\mathbf{U}$ , which represents excess rotation rate relative to strain rate [26]. For incompressible flows, it can be computed as  $Q = \frac{1}{2}(\|\boldsymbol{\Omega}\|^2 - \|\mathbf{S}\|^2)$ , where  $\boldsymbol{\Omega}$  and  $\mathbf{S}$  are rotation rate and strain rate tensors, respectively;  $\|\boldsymbol{\Omega}\| = \sqrt{\text{tr}(\boldsymbol{\Omega}\boldsymbol{\Omega}^T)}$  and  $\|\mathbf{S}\| = \sqrt{\text{tr}(\mathbf{S}\mathbf{S}^T)}$ , with superscript  $T$  indicating tensor transpose and  $\text{tr}$  indicating trace. The  $Q$  criterion is widely used in CFD simulations as a post-processing tool to identify vortex structures for the visualization of flow structures [27]. Second, the wall-distance based Reynolds number  $\text{Re}_d = \sqrt{kd}/\nu$  in  $q_3$  is an indicator to distinguish boundary layers from shear flows. This is an important feature because RANS models behave very differently in the two types of flows. This quantity is frequently used in wall functions for turbulence models. Third, feature  $q_7$  defines the deviation from orthogonality between the velocity and its gradient [28], which indicates the deviation of the flow from parallel shear flows (e.g., plane channel flows). Since most RANS models are calibrated to yield good performance on parallel shear flows, this deviation is usually correlated well with large discrepancies in RANS predictions. However, since it only accounts for misalignment angle and not the velocity magnitude, in regions with near-zero velocities this quantity becomes the angle formed by two zero-length vectors and is thus mostly noise. Finally, we remark that most of the features in Table I including the  $Q$  criterion and wall-distance based Reynolds number are familiar to CFD practitioners.

Most of the features presented in Table I are formulated as ratios of two quantities of the same dimension, either explicitly ( $q_1, q_5, q_6, q_8, q_9$ ) or implicitly ( $q_2$  and  $q_3$ ). Hence, they are nondimensional by construction. Features  $q_4$  and  $q_7$  involve inner product of vectors or tensors, and thus they are normalized by the magnitude of the constituent vectors or tensors. Finally, feature  $q_{10}$  (streamline curvature) is normalized by  $1/L_c$ , where  $L_c$  is the characteristics length scale of the mean flow, chosen to be the hill height  $H$  (see Fig. 6) in the numerical examples.

We followed the procedure of Ling and Templeton [12] to normalize the features. Except for feature  $q_3$ , each element  $q_\beta$  in the input vector  $\mathbf{q}$  is normalized as

$$q_\beta = \frac{\hat{q}_\beta}{|\hat{q}_\beta| + |q_\beta^*|}, \quad \text{where } \beta = 1, 2, 4, \dots, 10, \quad (1)$$

where the summation on repeated indices is not implied,  $\hat{q}_\beta$  are raw values of the features, and  $q_\beta^*$  are the corresponding normalization factors. This normalization scheme limits the numerical range of the inputs within  $[-1, 1]$  and thus facilitates the regression. The normalization is not needed for feature  $q_3$  (wall-distance based Reynolds number) since it is already in a nondimensional form and in a limited range  $[0, 2]$ .

It can be seen that the choices of features and normalization factors heavily rely on physical understanding of the problem (turbulence modeling). That is, in the present data-driven modeling framework, the data are utilized only *after* physical reasoning from the modeler has been applied. This task can be a burden in certain applications. It is worth noting that the recent work of Ling *et al.* [29] aimed to relieve the modeler from such burdens by using a basis of invariants of tensors relevant in the specific application (e.g., strain rate  $\mathbf{S}$  in turbulence modeling). Their work has the potential to systematically construct the input features based on raw physical variables and thus makes data-driven modeling even “smarter.”

#### D. Representation of Reynolds stress discrepancies as responses

In Sec. IC the Reynolds stress discrepancies  $\Delta\tau$  have been identified as the responses for the regression functions. The response quantities should also be based on Galilean invariant quantities due to the same consideration as in the feature choice. As such, individual components of the Reynolds stresses or the discrepancies based thereon are not suitable, but those based on their eigenvalues or invariants are preferred. In turbulence modeling, the Lumley triangle has been widely used for the analysis of turbulence states related to realizability [30]. It is formulated based on the second and third invariants ( $II$  and  $III$ ) of the anisotropy tensor. Recently, Banerjee *et al.* [31] proposed an improved formulation in which the eigenvalues of the anisotropy tensor are mapped to barycentric coordinates as opposed to the variants  $II$  and  $III$  as in the Lumley triangle. An important advantage of their formulation is that the mapping to barycentric coordinates is linear, which is in contrast to the nonlinear mapping to invariants  $II$  and  $III$ . Therefore, barycentric coordinates provide a nondistorted visual representation of anisotropy and are easier for imposing realizability constraints. The formulation of discrepancy starts with the eigen-decomposition of the Reynolds stress anisotropy tensor  $\mathbf{A}$ :

$$\boldsymbol{\tau} = 2k\left(\frac{1}{3}\mathbf{I} + \mathbf{A}\right) = 2k\left(\frac{1}{3}\mathbf{I} + \mathbf{V}\Lambda\mathbf{V}^T\right) \quad (2)$$

where  $k$  is the turbulent kinetic energy, which indicates the magnitude of  $\boldsymbol{\tau}$ ;  $\mathbf{I}$  is the second-order identity tensor; and  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3]$  and  $\Lambda = \text{diag}[\lambda_1, \lambda_2, \lambda_3]$  with  $\lambda_1 + \lambda_2 + \lambda_3 = 0$  are the orthonormal eigenvectors and eigenvalues of  $\mathbf{A}$ , respectively, indicating its shape and orientation.

In the barycentric triangle, the eigenvalues  $\lambda_1, \lambda_2$ , and  $\lambda_3$  are mapped to the barycentric coordinates ( $C_1, C_2, C_3$ ) as follows:

$$C_1 = \lambda_1 - \lambda_2, \quad (3a)$$

$$C_2 = 2(\lambda_2 - \lambda_3), \quad (3b)$$

$$C_3 = 3\lambda_3 + 1, \quad (3c)$$

with  $C_1 + C_2 + C_3 = 1$ . As shown in Fig. 1, the barycentric coordinates of a point indicate the portion of areas of three subtriangles formed by the point and with edges of barycentric triangle. For example, a point located on the top vertex corresponds to  $C_3 = 1$  while a point located on the bottom edge has  $C_3 = 0$ . Similar to the Lumley triangle, all realizable turbulences are enclosed in the barycentric triangle (or on its edges) and have positive barycentric coordinates  $C_1, C_2$ , and  $C_3$ . The barycentric triangle has been used by Emory *et al.* [32] as a mechanism to impose realizability of Reynolds stresses in estimating uncertainties in RANS simulations.

Placing the triangle in a Cartesian coordinate system  $\boldsymbol{\xi} \equiv (\xi, \eta)$ , the location of any point within the triangle is a convex combination of those of the three vertices, i.e.,

$$\boldsymbol{\xi} = \boldsymbol{\xi}_{1c}C_1 + \boldsymbol{\xi}_{2c}C_2 + \boldsymbol{\xi}_{3c}C_3 \quad (4)$$

where  $\boldsymbol{\xi}_{1c}$ ,  $\boldsymbol{\xi}_{2c}$ , and  $\boldsymbol{\xi}_{3c}$  denote coordinates of the three vertices of the triangle. An advantage of representing the anisotropy of Reynolds stress in the barycentric coordinates is that it has a clear physical interpretation, i.e., the dimensionality of the turbulence state [33]. Typically, the standard-RANS-predicted Reynolds stress at a near-wall location is located close to the isotropic, three-component state (vertex 3C-I) in the barycentric triangle, while the true stress is near the two-component limiting states (bottom edge). Moreover, the spatial variations from the near-wall region to the shear flow region are indicated as arrows in Fig. 1. It is clear that the trend of spatial variation predicted by a standard RANS model is opposite to that of the actual trend.

The three mutually orthogonal, unit-length eigenvectors  $\mathbf{v}_1, \mathbf{v}_2$ , and  $\mathbf{v}_3$  indicate the orientation of the anisotropy tensor. They can be considered a rigid body and thus their orientation has three degrees of freedom, although they have nine elements in total. We use the Euler angle with the  $z$ - $x'$ - $z''$  convention to parametrize the orientation following the convention in rigid body dynamics [34]. That is, if a local coordinate system  $x$ - $y$ - $z$  spanned by the three eigenvectors of  $\mathbf{V}$  was initially

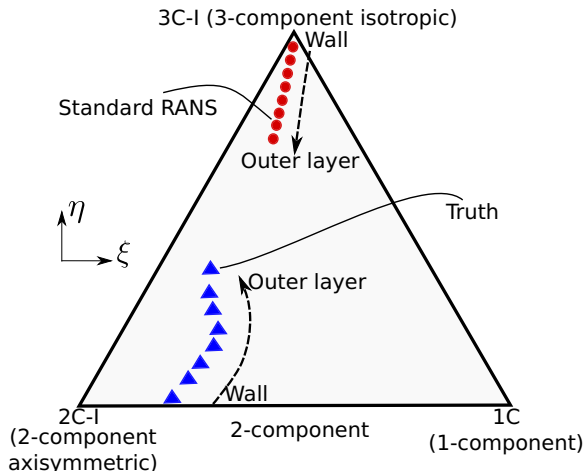


FIG. 1. The barycentric triangle that encloses all physically realizable states of Reynolds stress [31,33]. The position within the barycentric triangle represents the anisotropy state of the Reynolds stress. Typical mapped locations of near-wall turbulence states are indicated with predictions from standard RANS models near the isotropic state (vertex 3C-I), and the actual locations are indicated near the bottom edge (2C-I). The typical RANS-predicted trend of spatial variation from the wall to shear flow and the corresponding actual trend are indicated with arrows.

aligned with the global coordinate system ( $X$ - $Y$ - $Z$ ), the current configuration could be obtained by the following three consecutive intrinsic rotations about the axes of the local coordinate system: (1) a rotation about the  $z$  axis by angle  $\varphi_1$ , (2) a rotation about the  $x$  axis by  $\varphi_2$ , and (3) another rotation about its  $z$  axis by  $\varphi_3$ . The local coordinate axes usually change orientations after each rotation.

In summary, the Reynolds stress tensor is projected to six physically interpretable, Galilean invariant quantities representing the magnitude ( $k$ ), shape ( $\xi$ ,  $\eta$ ), and orientation ( $\varphi_1$ ,  $\varphi_2$ ,  $\varphi_3$ ). They are collectively denoted as  $\tau_\alpha$ . The actual values of these quantities can be written as baseline RANS predictions corrected by the corresponding discrepancy terms, i.e.,

$$\log_2 k = \log_2 \tilde{k}^{\text{rans}} + \Delta \log_2 k, \quad (5a)$$

$$\xi = \tilde{\xi}^{\text{rans}} + \Delta \xi, \quad (5b)$$

$$\eta = \tilde{\eta}^{\text{rans}} + \Delta \eta, \quad (5c)$$

$$\varphi_i = \tilde{\varphi}_i^{\text{rans}} + \Delta \varphi_i, \quad \text{for } i = 1, 2, 3. \quad (5d)$$

The discrepancies ( $\Delta \log_2 k$ ,  $\Delta \xi$ ,  $\Delta \eta$ ,  $\Delta \varphi_1$ ,  $\Delta \varphi_2$ ,  $\Delta \varphi_3$ , denoted as  $\Delta \tau_\alpha$  with  $\alpha = 1, 2, \dots, 6$ ) in the six projections of the Reynolds stress tensor are responses of the regression functions. We utilize data consisting of pairs of  $(\mathbf{q}, \Delta \tau_\alpha)$  from training flow(s) to construct the functions  $f_\alpha : \mathbf{q} \mapsto \Delta \tau_\alpha$ . It is assumed that the discrepancies in six quantities  $\Delta \tau_\alpha$  are independent, and thus separate functions are built for each of them. This simplification is along the same lines as that made in previous works [35].

### E. Random Forests for building regression functions

With the input (mean flow features  $\mathbf{q}$ ) and responses (Reynolds stress discrepancies  $\Delta \tau_\alpha$ ) identified above, a method is needed to construct regression functions from training data and to make predictions based on these functions. Supervised machine learning consists of a wide variety of such methods including  $K$ -nearest neighbors [36], linear regression and its variants (e.g., Lasso) [37],



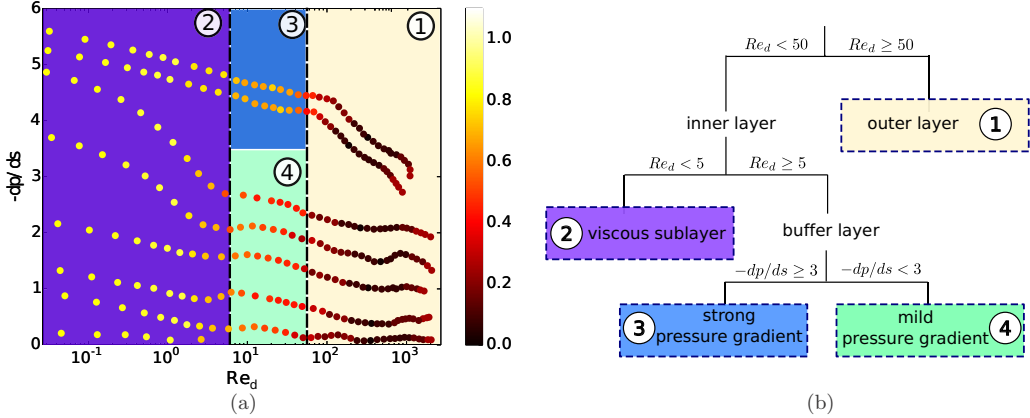


FIG. 2. Schematic of a simple regression tree in a two-dimensional feature space (pressure gradient along streamline  $dp/ds$  and wall-distance based Reynolds number  $Re_d$ ), showing (a) the stratification of feature space and (b) the corresponding regression tree built from the training data. The response is the discrepancy  $\Delta\eta$  in the barycentric triangle of the RANS-predicted Reynolds stress. When predicting the discrepancy for a given feature vector  $\hat{\mathbf{q}}$ , the tree model in (b) is traversed to identify the leaf, and the mean of the training data is taken as the prediction  $\Delta\eta(\hat{\mathbf{q}})$ .

Gaussian processes [23], tree-based methods (decision trees, random forests, bagging) [24], neural networks [38], and support vector machines [39], among others. A major consideration in choosing the regression method is the high dimensionality of the feature space, which is typically 10 or higher in our application. The curse of dimensionality makes such methods as  $K$ -nearest neighbors, linear regression, and Gaussian processes not suitable. A secondary consideration, which we believe is also important for turbulence modeling applications, is the capability to provide predictions with quantified uncertainties as well as physical insights (e.g., on the importance of each of the features and their interactions). After evaluating a number of existing machine learning techniques in light of these criteria, we identified *random forests* [24] as the optimal approach for our purposes, which is an ensemble learning technique based on decision trees.

In simple decision tree learning, a treelike model is built to predict the response variable by learning simple if-then-else decision rules from the training data. Decision trees have the advantage of being easy to interpret (e.g., via visualization) and implement. They are also computationally cheap. However, they tend to overfit the data and lack robustness. That is, a small change in the training data can result in large changes in the built model and its predictions. Random forests learning is an ensemble learning technique proposed by Ho [40] and Briemann [24,41] that overcomes these shortcomings of simple decision trees. Since these techniques are generally not familiar to readers in the fluid dynamics community, here we use an illustrative example in the context of turbulence modeling to explain the algorithm.

A simple decision tree model is illustrated in Fig. 2. For clarity we consider an input with only two features: pressure gradient  $dp/ds$  (normalized and projected to the streamline tangential) and wall-distance based Reynolds number  $Re_d = C_\mu^{1/4} d\sqrt{k}/\nu$ , as defined in Table I. It can be also interpreted as wall distance in viscous units. The response is the discrepancy  $\Delta\eta$  of the vertical coordinate in the barycentric triangle of the RANS-predicted Reynolds stress (see Fig. 1). During the training process, the feature space is *successively* divided into a number of boxes (leaves) based on the training data [shown as points in the  $dp/ds - Re_d$  plane in Fig. 2(a)]. In the simplest decision tree model used for regression, the feature space is stratified with the objective of minimizing the total in-leaf variances of the responses at each step, a strategy that is referred to as a greedy algorithm. After the stratification, a constant prediction model is built on each leaf. When predicting the response  $\Delta\eta$  for a given feature vector  $\mathbf{q}$ , the constructed tree model in Fig. 2(b) is traversed

to identify the leaf where  $\mathbf{q}$  is located, and the mean response on the leaf is taken as the prediction  $\Delta\eta(\mathbf{q})$ .

The tree model has a clear physical interpretation in the context of turbulence modeling. For example, it is well known that a standard isotropic eddy viscosity model has the largest discrepancy when predicting anisotropy in the viscous sublayer ( $Re_d \leq 5$ ). This is because the truth is located on the bottom, corresponding to a combination of one- or two-component turbulence, while a typical isotropic eddy viscosity model would predict an isotropic state located on the top vertex (see Fig. 1). In contrast, far away from the wall within the outer layer ( $Re_d > 50$ ), the RANS-predicted anisotropy is rather satisfactory. Therefore, the first two branches divide the space into three regions based on the feature  $Re_d$ : outer layer (region 1), viscous sublayer (region 2), and buffer layer (regions 3 and 4). In the buffer layer the pressure gradient plays a more important role than in the outer and viscous layers. Larger pressure gradients correspond to larger discrepancies, which can be explained by the fact that favorable pressure gradients (negative  $dp/ds$  values) tend to thicken the viscous sublayer [42], which leads to larger discrepancies in  $\eta$ . Therefore, a further division splits the buffer layer states to two regions in the feature space, i.e., those with strong (region 3) and mild (region 4) pressure gradients.

A simple regression tree model described above tends to overfit for a high-dimensional input space, i.e., yielding models that explain the training data very well but predict poorly for unseen data. In general the decision trees do not have the same level of predictive accuracy as other modern regression methods. However, by aggregating a large number of trees (ideally with minimum correlation), the predictive performance can be significantly improved and the overfitting can be largely avoided. In random forests an ensemble of trees is built with bootstrap samples (i.e., sampling with replacement) drawn from the training data [20]. Moreover, when building each tree, it utilizes only a subset of  $M \leq N_q$  randomly chosen features among the  $N_q$  features, which reduces the correlation among the trees in the ensemble and thus decreases the bias of the ensemble prediction.

Random forest regression is a modern machine learning method with predictive performance comparable to other state-of-the-art techniques [37]. In decision tree models the maximum depth of trees must be limited (e.g., by pruning the branches far from the root) to ensure a sufficient number of training points (e.g., 5) on each node. In contrast, in random forests, one can build each tree to its maximum depth by successive splitting the nodes until only one training data point remains on each leaf. While each individual tree built in this manner may suffer from overfitting and has large prediction variances, the use of ensemble largely avoids both problems. Moreover, random forest regression is simple to use with only two free parameters, i.e., the number  $N_{rf}$  of trees in the ensemble and the number  $M$  of selected features. In this work we used an ensemble of  $N_{rf} = 100$  trees and the a subset of features (i.e.,  $M = 6$ ) to build each tree. As a standard practice in statistical modeling, we performed cross-validations to optimize these parameters and performed sensitivity analysis to ensure that the predictions are not sensitive to the parameter choices.

### III. NUMERICAL RESULTS

Almost all industrial flows involve some characteristics (e.g., strong pressure gradient, streamline curvature, and separation) that break the equilibrium assumption of RANS model. Therefore, we have these challenges in mind when developing the data-driven approach. In this study, we focus on the cases where training and test flows have similar characteristics. Specifically, we evaluate the proposed method on two classes of flows: (1) fully developed turbulent flows in a square duct at various Reynolds numbers and (2) flows with massive separations. The flow in a square duct at Reynolds number  $Re = 3500$  and the flow in a channel with periodic hills at Reynolds number  $Re = 10595$  are chosen as the prediction (test) flows for the respective flow classes. The square duct flow has an in-plane secondary flow pattern induced by the normal stress imbalance, while the periodic-hill flow features a recirculation bubble, nonparallel shear layer and mean flow curvature.

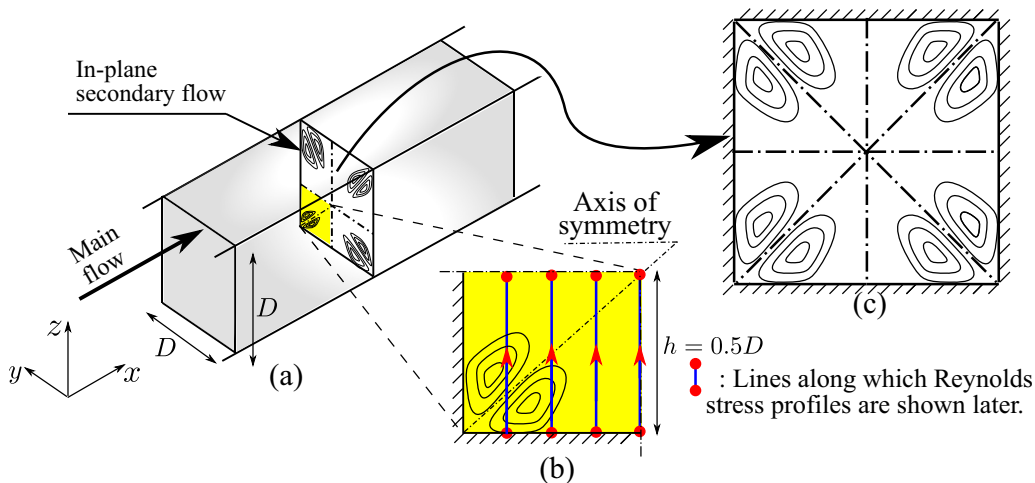


FIG. 3. Domain shape for the flow in a square duct. The  $x$  coordinate represents the streamwise direction. Secondary flows induced by Reynolds stress imbalance exist in the  $y$ - $z$  plane. Panel (b) shows that the computational domain covers a quarter of the cross section of the physical domain. This is due to the symmetry of the mean flow in both  $y$  and  $z$  directions as shown in panel (c).

All these characteristics are known to pose challenges for RANS based turbulence models, and thus large model-form discrepancies exist in the RANS-modeled Reynolds stresses. In the two test flows, the relative importance of Reynolds stress projections to the mean flow prediction are different. The Reynolds stress anisotropy plays an important role in obtaining the accurate secondary mean motion in the duct flow [43]. In contrast, the anisotropy is less important to predict the mean flow in the periodic-hill case, where the turbulent shear stress component is more essential to obtain an accurate mean velocity field [44]. Therefore, we use these two types of flows to highlight the improvements in the different Reynolds stress components that are important for the predictions of QoIs in the respective flow classes. In both cases, all RANS simulations are performed in an open-source CFD platform, OpenFOAM, using a built-in incompressible flow solver simpleFOAM [45]. Mesh convergence studies have been performed.

## A. Turbulent flows in a square duct

### 1. Case setup

The fully developed turbulent flow in a square duct is a challenging case for RANS-based turbulence models, since the secondary mean motion cannot be captured by linear eddy viscosity models (e.g.,  $k$ - $\epsilon$ ,  $k$ - $\omega$ ), and even the Reynolds stress transport models (RSTM) cannot predict it well [44]. In this test, we aim to improve the RANS-modeled Reynolds stresses of the duct flow at Reynolds number  $Re = 3500$  by using the proposed PIML approach. The training data are obtained from DNS simulations [46] of the duct flows in the same geometry but at lower Reynolds numbers  $Re = 2200, 2600,$  and  $2900$ . The DNS data of the prediction flow ( $Re = 3500$ ) are reserved for comparison and are not used for training. The geometry of this flow case is shown in Fig. 3. The Reynolds number is based on the edge length  $D$  of the square duct and the bulk velocity  $U_b$ . All lengths presented below are normalized by  $D/2$ .

The baseline RANS simulations are performed for all training and test flows. The purpose is twofold: to obtain the mean flow feature fields  $\mathbf{q}(\mathbf{x})$  as inputs and to obtain the discrepancies of Reynolds stress by comparing with the DNS data. To enable the comparison, the high-fidelity data are interpolated onto the mesh of the RANS simulation. The Launder-Gibson RSTM [47] is adopted to perform the baseline simulations, since all the linear eddy viscosity models are not able to capture

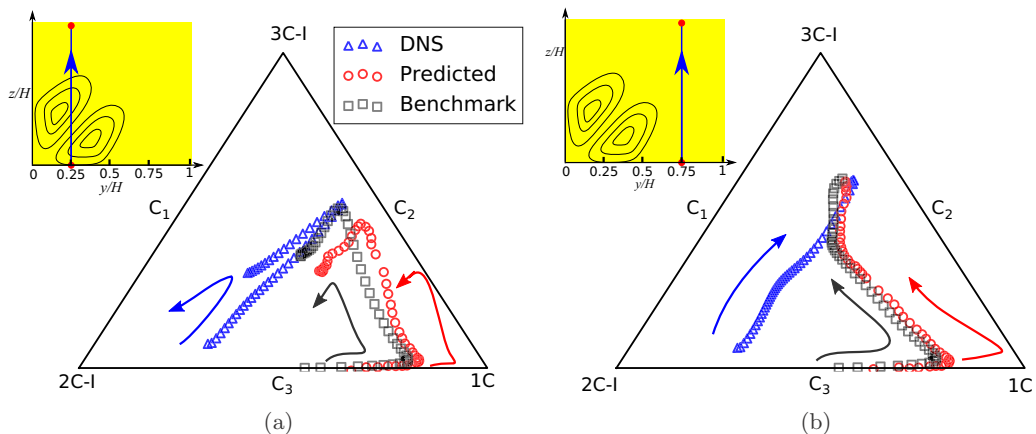


FIG. 4. Barycentric map of the predicted Reynolds stress anisotropy for the test flow ( $Re = 3500$ ), learned from the training flows ( $Re = 2200, 2600$ , and  $2900$ ). The prediction results on two streamwise locations at  $y/H = 0.25$  and  $0.75$  are compared with the corresponding baseline (RSTM) and DNS results in panels (a) and (b), respectively.

the mean flow features of the secondary motions. The  $y^+$  of the first cell center is kept less than 1 and thus no wall model is applied. As indicated in Fig. 3, only one quadrant of the physical domain is simulated owing to the symmetry of the mean flow with respect to the center lines along  $y$  and  $z$  axes. A no-slip boundary condition is applied on the walls, and a symmetry boundary condition is applied on the symmetry planes.

## 2. Prediction results

We first investigate the prediction performance on the Reynolds stress anisotropy tensor, since its accuracy is important to capture the secondary flow. Figure 4 shows PIML-corrected anisotropy in a barycentric triangle compared with baseline and DNS results. The comparisons are performed on two representative lines at  $y/H = 0.25$  and  $0.75$  on the in-plane cross section [Fig. 3(b)]. The two lines are indicated in the insets at the upper left corner of each panel. The arrows denote the order of sample points plotted in the triangle, which is from the bottom wall to the outer layer. The general trends of spatial variations of the DNS Reynolds stress anisotropies are similar on both lines. That is, from the wall to the outer layer, the Reynolds stress starts from the two-component limiting states (bottom edge of the triangle) toward three-component anisotropic states (middle area of the triangle). This trend is captured by the baseline RSTM to some extent, especially in the regions away from the wall. However, significant discrepancies still can be observed in the near-wall region. Very close to the wall, the DNS Reynolds stress is nearly the two-component limiting state. This is because the velocity fluctuations in the wall-normal direction are suppressed by the blocking of the bottom wall. Moreover, before approaching three-component anisotropic states, the DNS-predicted anisotropy first moves toward the one-component state (1C) as away from the wall. In contrast, the RANS-predicted anisotropy near the wall is closer to the two-component isotropic state (2C-I), and it approaches toward the three-component anisotropic state directly. Therefore, in the near-wall region there are large discrepancies between the RANS predicted Reynolds stress anisotropy and the DNS result, particularly in the horizontal coordinate  $\xi$ . By correcting the baseline RSTM results with the trained discrepancy function, the predicted anisotropy of Reynolds stress is significantly improved. For both lines, the predicted anisotropy (circles) agrees well with the DNS results (squares). Especially on the line  $y/H = 0.75$ , the PIML-predicted anisotropy is almost identical to the DNS data.

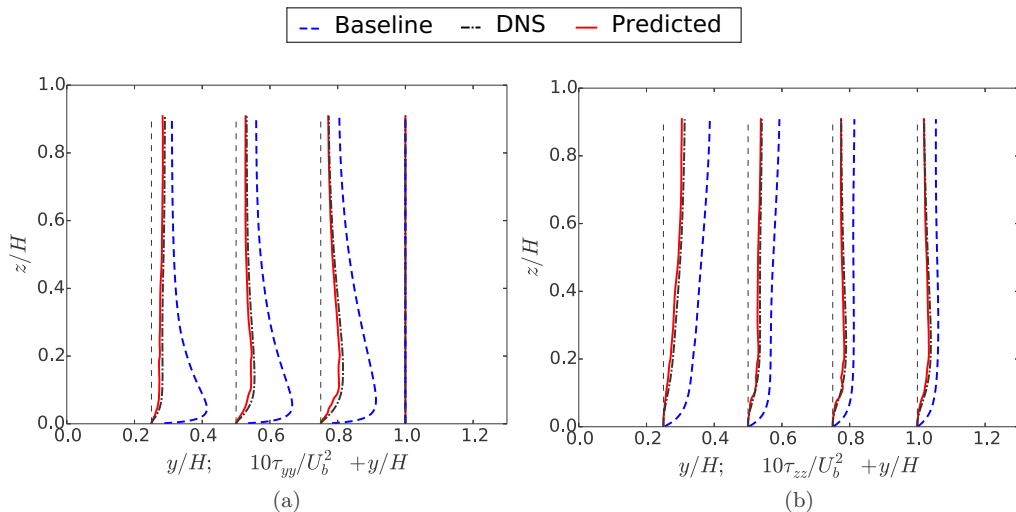


FIG. 5. Profiles of normal components (a)  $\tau_{yy}$  and (b)  $\tau_{zz}$  of corrected Reynolds stress with the discrepancy model. The profiles are shown at four streamwise locations  $y/H = 0.25, 0.5, 0.75,$  and  $1$ . Corresponding DNS and baseline (RSTM) results are also plotted for comparison.

Significant improvement of the PIML-predicted anisotropy can be seen from the barycentric maps shown in Fig. 4. Similar improvements have also been demonstrated in the other physical projections (TKE and orientations) of the PIML-corrected Reynolds stresses. Therefore, it is expected that the Reynolds stress tensor components should be also improved over the RSTM baseline. In the six tensor components, two normal stress components  $\tau_{yy}$  and  $\tau_{zz}$  are among the most important ones to the mean velocity field since the normal stress imbalance ( $\tau_{yy} - \tau_{zz}$ ) is the driving force of the secondary flow [43]. Figures 5(a) and 5(b) show the profiles of normal components  $\tau_{yy}$  and  $\tau_{zz}$  of the Reynolds stress reconstructed from the PIML-corrected physical projections. Corresponding baseline (RSTM) and DNS results are also plotted for comparison. Both  $\tau_{yy}$  and  $\tau_{zz}$  are overestimated by the RSTM over the entire domain. The discrepancy of the RSTM-predicted  $\tau_{yy}$  is large near the wall and decreases when moving away from the wall. In contrast,  $\tau_{zz}$  is significantly overestimated far from the wall but the discrepancy decreases toward the wall. As a result, the RSTM-predicted normal stress imbalance is markedly inaccurate, which leads to unreliable secondary mean flow motion. As expected, the PIML predictions nearly overlap with the DNS results for both  $\tau_{yy}$  and  $\tau_{zz}$  and show considerable improvements over the RSTM baseline. In fact, the improvements are observed in all the tensor components, which are omitted here for brevity. The results shown above demonstrate excellent performance of the proposed PIML framework by using RSTM as the baseline.

## B. Turbulent flows with massive separations

### 1. Case setup

The turbulent flow in a channel with periodic hills is another challenging case for RANS models due to the massive flow separations leeward of the hill. Here, we examine two training scenarios with increasing difficulty levels. In the first scenario the training flows have the same geometry as the test (prediction) flow but are different in Reynolds numbers. In the second scenario the training flows differ from the prediction case not only in Reynolds numbers but also in geometry.

Four training flows with DNS/LES data to build random forest regressors are summarized in Table II. In the first scenario two flows PH1400 and PH5600 are used for training, both of which are flows over periodic hills (same in geometry) at  $Re = 1400$  and  $Re = 5600$  (different in Reynolds numbers), respectively. For the second scenario, the training data are obtained from two different

TABLE II. Database of training flows to predict flow past periodic hills at  $Re = 10595$ . The Reynolds numbers are defined based on the bulk velocity  $U_b$  at the narrowest cross section in the flow and the crest/step height  $H$ .

Training flow scenario	Training flow & symbol	High fidelity data
Scenario I	Periodic hills, $Re = 1400$ (PH1400)	DNS by Breuer <i>et al.</i> [48]
	Periodic hills, $Re = 5600$ (PH5600)	DNS by Breuer <i>et al.</i> [48]
Scenario II	Wavy channel, $Re = 360$ (WC360)	DNS by Maaß <i>et al.</i> [49]
	Curved backward facing step, $Re = 13200$ (CS13200)	LES by Bentalieb <i>et al.</i> [50]

flows: one in a channel with a wavy bottom wall at  $Re = 360$  and one over a curved backward facing step at  $Re = 13200$ , indicated as flows WC360 and CS13200, respectively.

A schematic of the flow geometry and RANS-predicted velocity contour for each case are presented in Fig. 6. The dimensions of each case are normalized with the respective hill heights  $H$ . Although the geometries of the training flows are different, all three flows share a similar characteristic as the test flow, i.e., separation on the leeward side of the hill or step. However,

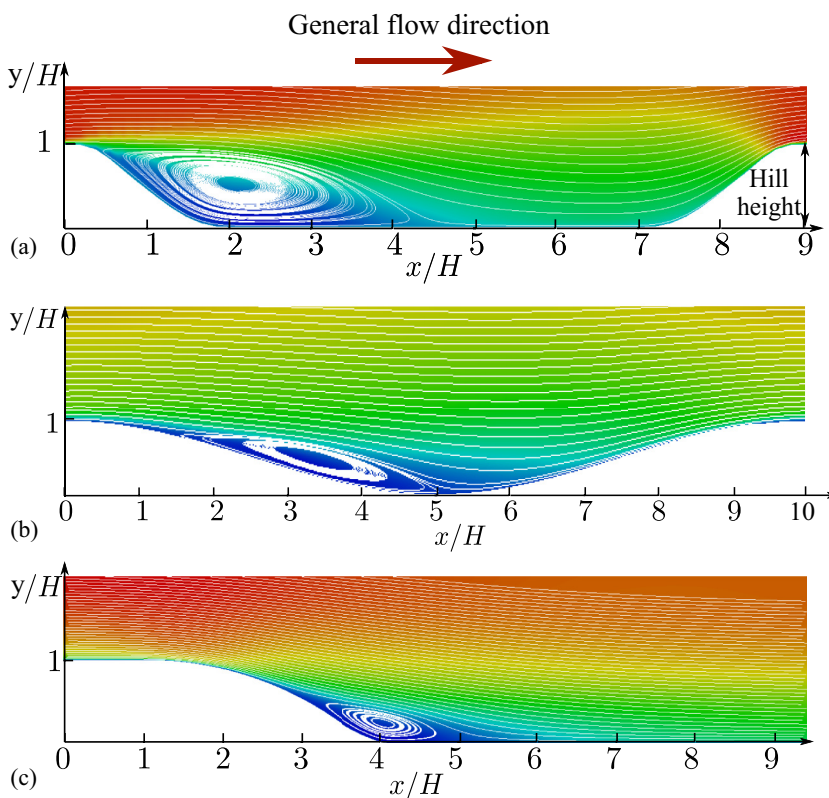


FIG. 6. Computational domain and velocity field of each case in the training flow database. The velocity contours and streamlines are obtained from the baseline RANS simulations. The dimensions of each case are normalized with the respective hill heights  $H$ . Note that periodic boundary conditions are applied on the flows in panels (a) and (b) in the streamwise direction, but not for the flow in panel (c). (a) Periodic hills ( $Re = 1400$  and  $5600$ , results from latter are shown), (b) wavy channel ( $Re = 360$ ), and (c) curved backward facing step ( $Re = 13200$ ).

the separation bubbles are different in size and shape. The flow over periodic hills has a stronger separation compared to the other two due to the steeper slope of the hill. Relatively mild separation can be observed in the flow over the wavy channel. For all cases, both high-fidelity data and RANS-predicted results are available. The high-fidelity data are obtained from DNS or resolved LES simulations, which have been reported in literature (see references in Table II). The DNS data for flows PH1400 and PH5600 are only available on vertical lines at eight streamwise locations  $x/H = 1, 2, \dots, 8$ . On the other hand, full-field high-fidelity results are available for flows WC360 and CS13200, but only the lower part of the channel is adequately resolved. Since the separated flow is of interest in this study, only the data in the separation region (i.e., the region below  $y/H = 1.2$ ) are included.

In this test, the performance of the proposed PIML framework is evaluated on standard RANS models. Specifically, the baseline RANS predictions are obtained by using the two-equation Launder-Sharma  $k-\varepsilon$  model [3]. The reason for choosing standard turbulence models here is because of two considerations. First, the standard RANS models are the dominant tools for industrial CFD applications, while other sophisticated RANS models have been rarely used. Therefore, it is more significant to improve the widely used standard RANS models. Second, we understand that improvement of the Reynolds stresses starting from a standard RANS model is challenging. Nonetheless, this challenging scenario also can better explore the capability of a machine learning approach.

## 2. Prediction results

The functional forms of discrepancies in the six physical projections of Reynolds stress are learned from the training flows, as mentioned in Sec. III B 1, and are used to correct the RANS-predicted Reynolds stress field of the test flow (PH10595). However, since the baseline RANS model used in this case is the standard eddy viscosity model, the Reynolds stress anisotropy cannot be accurately predicted. Therefore, the baseline RANS-predicted anisotropy is unphysical and is significantly different from the DNS result (see Fig. 1). Nonetheless, after the correction by using the discrepancy function learned from the training flows, the anisotropy of the test flow shows an excellent agreement with the DNS results [51]. The improvements are observed in the both training scenarios I and II, demonstrating that the discrepancy function even in the standard RANS-predicted anisotropy does exist and can be learned from the closely related flows based on the mean flow features  $\mathbf{q}$ . As mentioned above, in the periodic-hill flow, the correctness of Reynolds stress anisotropy is of little consequence to the prediction of the mean velocity, and the correct shear stress component and magnitude of the Reynolds stress are most important to obtain an accurate mean flow field. Therefore, the anisotropy prediction results are omitted here, and only the turbulence kinetic energy (TKE) and shear stress component of the PIML-corrected Reynolds stress are presented and discussed in detail.

The comparison of the TKE profiles of the baseline, DNS, and PIML-predicted results in the training scenario I are shown in Fig. 7. The TKE predicted by the baseline RANS model has notable discrepancies compared to the DNS result, particularly in the region with nonparallel free shear flow ( $y/H = 0.8$  to  $1.5$ ). The poor performance of the RANS model in such a region is typical in these flows [13]. The RANS model underestimates the turbulence intensity along the free shear at  $y/H = 1$ , especially near the leeward side of the hill ( $x/H = 1$  to  $2$ ). In the upper channel ( $y/H = 1.5$  to  $2.5$ ), the DNS TKE is slightly smaller than the baseline RANS prediction. The profiles of TKE corrected by the PIML-predicted discrepancy  $\Delta \log_2 k$  are significantly improved. The peaks along the streamwise free shear in the DNS profiles are well captured in the corrected results with the random forest prediction. It can be seen that the predicted TKE profiles (solid lines) nearly overlap with the DNS results (dashed lines). This clearly indicates that the TKE discrepancies can be learned from the data of the training flows.

It is also of interest to investigate the tensor components  $\tau_{ij}$  of Reynolds stress, which are more relevant for predicting velocities and other QoIs of the flow fields. For the plane shear flows, the turbulent shear stress  $\tau_{xy}$  is important to predict the velocity field. Figure 8 compares the

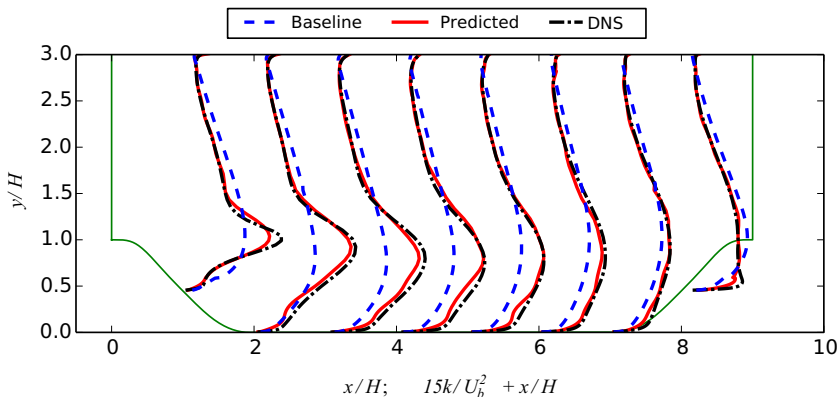


FIG. 7. Magnitude (turbulence kinetic energy) of the corrected Reynolds stress for the test flow (PH10595) learned from cases with with same geometry but at different Reynolds numbers (PH1400 and PH5600). The profiles are shown at eight streamwise locations  $x/H = 1, 2, \dots, 8$ . Corresponding baseline and DNS results are also plotted for comparison. The hill profile is vertically exaggerated by a factor of 1.3.

turbulence shear component  $\tau_{xy}$  of predicted Reynolds stress with the DNS. As expected, significant improvements are observed compared to the baseline results, which underestimate the peak of  $\tau_{xy}$  on the leeward hill side but overestimate it on the windward hill side. As shown in Fig. 8, the profiles of predicted  $\tau_{xy}$  agree well with the DNS results.

The results above demonstrate that the discrepancy function of Reynolds stress in its physical projections (i.e., magnitude, shape, and orientation) trained from the flows at  $Re = 2800$  and  $5600$  can be used to predict the Reynolds stress field of the flow at  $Re = 10\,595$ . Significant improvements are observed in the predicted Reynolds stress compared to the baseline RANS results. Although in this scenario the training and test flows are quite similar (with the same geometry), and the success of extrapolation has been demonstrated in physical space by Wu *et al.* [18], it should not be taken for granted that the accurate prediction is also guaranteed in feature space. Since the regressions are performed in the ten-dimensional feature space and there is no direct reference to the physical coordinates, success is not trivially expected *a priori*.

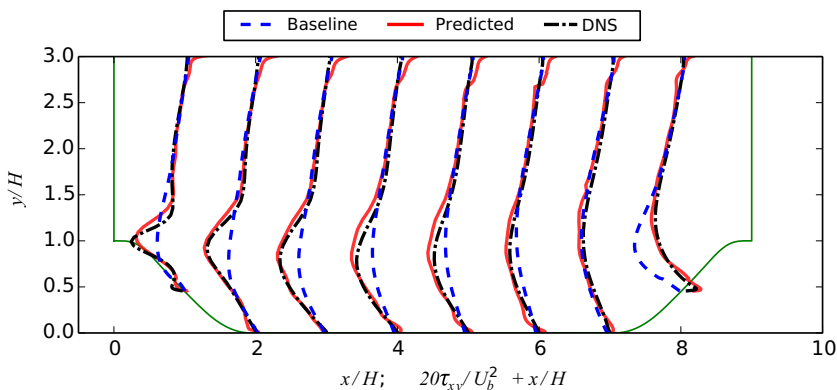


FIG. 8. Predicted turbulence shear stress for the test flow (PH10595) learned from the flows with the same geometry but at different Reynolds numbers (PH1400 and PH5600). The profiles are shown at eight streamwise locations  $x/H = 1, 2, \dots, 8$ . Corresponding baseline and DNS results are also plotted for comparison. The hill profile is vertically exaggerated by a factor of 1.3.



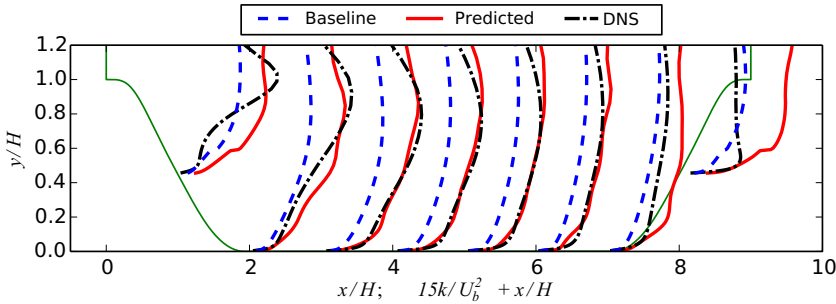


FIG. 9. Magnitude (turbulence kinetic energy) of the corrected Reynolds stress for the test flow (PH10595) learned from cases with different geometries and at different Reynolds numbers (WC360 and CS13200). The profiles are shown at eight streamwise locations  $x/H = 1, 2, \dots, 8$ . Corresponding baseline and DNS results are also plotted for comparison. The hill profile is vertically exaggerated by a factor of 2.4.

We investigate a more challenging scenario where the training flows have different geometries from the prediction case. This scenario is also more realistic in the context of using RANS simulation to support engineering design and analysis. Specifically, the data are more likely to be available for a few flows with specific Reynolds numbers and geometries, but predictions are needed for similar flows yet at different Reynolds numbers and with modified geometries.

The comparison of the TKE profiles on eight lines is shown in Fig. 9. Note that only the domain below  $y/H = 1.2$  is investigated due to the lack of reliable high-fidelity training data in the upper channel region. This inadequacy of data quality can be exacerbated when the Reynolds stress is decomposed to its physical projections. Moreover, the flow separation is the phenomenon of concern in this study, and thus we only focus on the recirculation region. In Fig. 9, the random forest predicted TKE (solid lines) is significantly improved over the baseline results (dotted lines) and better agrees with the DNS profiles (dash-dotted lines). The agreement is particularly good in the region from the center of recirculation bubble ( $x/H = 2$ ) to the beginning of flow contraction ( $x/H = 6$ ). Nonetheless, the PIML-predicted TKE does not show any improvement and even deteriorates compared to the baseline results near the windward side of the hill ( $x/H > 7$ ), where the flow starts to contract. As shown in Fig. 9, the predicted TKE is markedly overestimated at  $x/H = 8$ . This is because the flow features in the contraction region ( $x/H > 7$ ) are not supported in the training set, since the contracted flow does not exist in the training flow CS13200 [Fig. 6(c)] and is much weaker in the training flow WC360 [Fig. 6(b)] due to the mild slope of wavy bottom in this geometry.

Finally, we compare the predicted turbulent shear stress  $\tau_{xy}$  with the DNS profiles in Fig. 10. Similar to the results of TKE, the PIML-predicted turbulent shear stress  $\tau_{xy}$  shows notable improvements in the recirculation region. However, deterioration occurs in the flow contraction region. At  $x/H = 7$  and 8, the magnitudes of turbulent shear stresses are overestimated with the correction based on the predicted discrepancies. This is consistent with the results observed in physical projections of Reynolds stress. Such a small region with abnormal Reynolds stress corrections (artificial peaks or bumps) can introduce large errors into the velocity predictions.

In general, the physical projections (i.e., magnitude, shape, and orientation) of Reynolds stress corrected by random forest predicted discrepancies are still significantly improved with the training flows in different geometries (WC360 and CS13200). The Reynolds stress is markedly improved in the separated flow region, but not in the contracted flow region. This is because the features in training flows cannot well support the predicted flow, and thus more extrapolations are expected. Although the improvement is less significant compared to that in scenario I, the random forest predictions in this more realistic scenario are still satisfactory, demonstrating the merits of the proposed PIML framework.

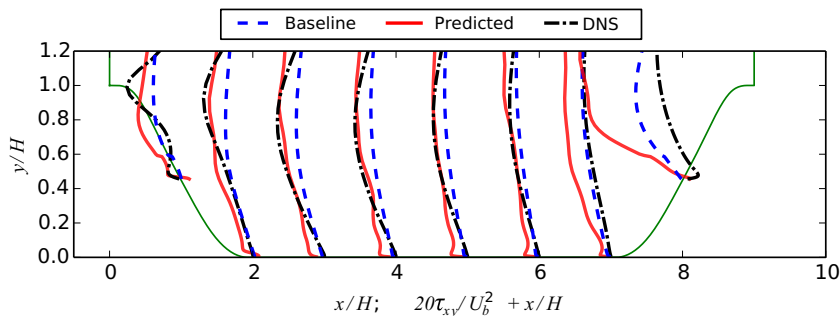


FIG. 10. Predicted turbulence shear stress for the test flow (PH10595) learned from the flows with different geometries and at different Reynolds numbers (WC360 and CS13200). The profiles are shown at eight streamwise locations  $x/H = 1, 2, \dots, 8$ . Corresponding baseline and DNS results are also plotted for comparison. The hill profile is vertically exaggerated by a factor of 2.4.

## IV. DISCUSSION

### A. Feature importance and insight for turbulence modeling

In addition to the predictive capability of the regression model, it is also important to interpret the functional relation between the mean flow features and the discrepancies of the RANS modeled Reynolds stresses. For example, it is useful to find the most important features for the Reynolds stress discrepancy in each of its physical projections (i.e., magnitude  $k$ , shape  $\xi, \eta$ , and orientation  $\varphi_i$ ), and how each of these features impacts the regression response. Identification of such a correlation or causal relationship enables modelers to improve the RANS turbulence models. The random forest regressor used in the proposed PIML framework can also shed light on this issue by calculating the feature importance, which is a measure to evaluate the relative importance of a feature variable for predicting response variables [24]. The bar plots of the importance of feature vector  $\mathbf{q}$  with respect to the discrepancies  $\Delta\eta$  and  $\Delta\log_2 k$  are shown in Figs. 11(a) and 11(b), respectively. For discrepancy  $\Delta\eta$  in the anisotropy, feature  $q_3$  (i.e., wall distance based Reynolds number  $Re_d$ ) is the most important one. As discussed in Sec. II C,  $Re_d$  is the wall distance normalized by the approximate viscous unit. Therefore, the result of feature importance is consistent with the PIML

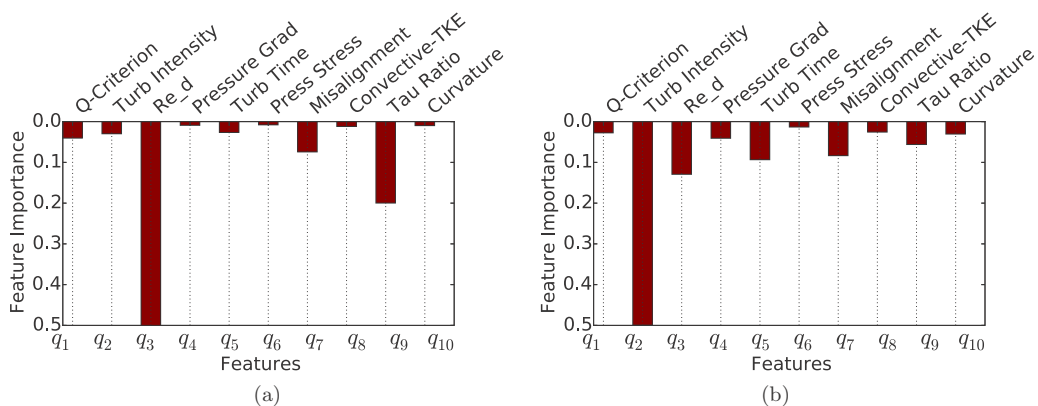


FIG. 11. Feature importance of random forest regressors (a) for  $\Delta\eta$  and (b) for  $\Delta\log_2 k$  for scenario I (i.e., training flows in the same geometry; see Table II). The features  $q_i$  ( $i = 1$  to 10) are denoted by their respective abbreviations. Turb Intensity denotes the turbulence intensity (feature  $q_2$ ), and Re\_d is the wall distance based Reynolds number (feature  $q_3$ ). For the full name list of features, see Table I.

prediction, which has shown that the discrepancy  $\Delta\eta$  is notably dependent on the distance away from the wall [51]. Figure 11(b) shows that the most important feature for predicting discrepancy  $\Delta \log_2 k$  of turbulence kinetic energy is feature  $q_2$ , turbulence intensity.

It is demonstrated that the random forest used in the proposed framework can interpret the relationship between the features and the response to a certain extent, although the feature importance has its limitation due to bias introduced under certain conditions [52,53]. In the machine learning community, improving interpretability of a random forest is an active research topic; e.g., several improvements of the importance measure have been proposed [53–55]. Moreover, in addition to calculating the feature importance, it is also helpful to examine the interactions among features, which have important implications for the interpretation of the regression models. As the base learner in random forest is a decision tree, which can capture the feature interactions, it is possible to further investigate the interacting relationship among mean flow variables. Breiman *et al.* have studied feature interaction in the random forest method [24], but more research is still ongoing. A better understanding of the physics behind the regression model for Reynolds stress discrepancies has a profound implication for RANS turbulence modeling. Therefore, to explore the correlation or causal relationship between the mean flow features and the discrepancies of RANS modeled Reynolds stress is an important and promising extension of the proposed framework.

### B. Success and limitation of the current framework

The objective of the proposed framework is to improve the baseline RANS-predicted Reynolds stresses of a flow where high-fidelity (e.g., DNS, LES, experimental) data are not available. The main novelty lies in using machine learning techniques to find the functional forms of Reynolds stress discrepancies with respect to mean flow features by learning from the existing offline database of the closely related flows. Numerical simulation results have demonstrated the feasibility and merits of the framework. Moreover, the excellent performance of the PIML predicted Reynolds stress in not only the anisotropy but also in the TKE and turbulent shear stress shows the fact that Reynolds stress discrepancies can be extrapolated even to complex flows sharing similar characteristics. This finding is noteworthy by itself.

The improvement of the RANS-predicted Reynolds stress is considered a viable and promising path toward obtaining better predictions of velocities and other quantities of interest. However, due to a few limitations of the current framework, the improvement of the propagated velocities from the corrected Reynolds stress field cannot be guaranteed. A small region with abnormal Reynolds stress corrections (e.g., nonsmoothness or artificial peaks) can introduce large errors to the velocity predictions. For example, the small wave-number variations in Reynolds stresses are visible in Fig. 10. These fluctuations, despite being small in amplitude, can lead to abnormal behaviors in the divergence term and thus in the predicted velocities. These abnormal predictions of Reynolds stress corrections can be caused by several factors. First, the features in certain regions of the prediction flow may not be well supported in the training flows, e.g. the contraction region of periodic-hill flow mentioned in Sec. III B. Second, the random forest regression used here only provides pointwise estimations but cannot consider the spatial information of the Reynolds stress field. Therefore, the smoothness of the prediction cannot be guaranteed. Finally, although the input feature space is constructed based on physical reasoning, it is possible that the input features are not rich enough, and thus the randomness in the ensemble of the trained decision trees is significant.

## V. CONCLUSION

In this work, we proposed a physics-informed machine learning approach to reconstruct Reynolds stresses modeling discrepancies by utilizing DNS databases of training flows sharing characteristics similar to those of the flow to be predicted. For this purpose, we formulated discrepancies of Reynolds stresses (or more precisely their magnitudes and the shape and orientation of the anisotropy) as target functions of mean flow features, and used modern machine learning techniques based on random

forest regression to learn the functions. The obtained functions are then used to predict Reynolds stress discrepancies in new flows. To evaluate the performance of the proposed approach, the method is tested by two classes of flows: (1) fully developed turbulent flows in a square duct at various Reynolds numbers and (2) flows with massive separations. In the separated flows, two training flow scenarios of increasing difficulties are considered: In the less challenging scenario, data from two flows in the same periodic hill geometry at lower Reynolds numbers ( $Re = 2800$  and  $5600$ ) are used for training. In a more challenging scenario, the training data come from separated flows in different geometries (wavy channel and curved backward facing step). In all test cases the corrected Reynolds stresses are significantly improved compared to the baseline RANS predictions, demonstrating the merits of the proposed approach. In the scenario where the training flows and the prediction flow have different geometries, the improvement is not as drastic as in the scenario with the same geometry. This is expected since the prediction involves more extrapolations in the feature space for this more challenging scenario. In other words, compared to the first scenario where the training and prediction flows have identical geometry, the prediction flow is less “similar” to the training flows in this scenario. The extent to which the training and prediction flows are “similar” to each other can be assessed *a priori* based on their respective RANS predicted mean flow field, and methods for such assessment are presented in companion publications [56,57].

As the inaccuracy in modeled Reynolds stresses is the dominant source of model-form uncertainty in RANS simulations, the proposed method for improving RANS-predicted Reynolds stresses is an important step towards the goal of enabling predictive capabilities of RANS models. Moreover, the random forests regression technique adopted in this work can provide physical insights regarding the relative importance of mean flow features that contributed to the discrepancies in the RANS predicted Reynolds stresses. This information can be used to assist future model development in that developers can devise models that are aware of and correctly respond to these flow features. However, a number of challenges need to be tackled before the improved Reynolds stresses can be used to predict more accurate quantities of interests that are needed in engineering design (e.g., draft and lift coefficients). This topic will be investigated in future research [58].

#### ACKNOWLEDGMENTS

We thank Dr. Julia Ling of Sandia National Laboratories and Dr. Eric G. Paterson of Virginia Tech for helpful discussions during this work. We also thank the anonymous reviewers for their comments, which helped improve the quality and clarity of the manuscript.

- 
- [1] M. M. Waldrop, The chips are down for Moore’s law, *Nat. News* **530**, 144 (2016).
  - [2] S. Kumar, Fundamental limits to Moore’s law, [arXiv:1511.05956](https://arxiv.org/abs/1511.05956).
  - [3] B. Launder and B. Sharma, Application of the energy-dissipation model of turbulence to the calculation of flow near a spinning disc, *Lett. Heat Mass Transfer* **1**, 131 (1974).
  - [4] D. C. Wilcox, Reassessment of the scale-determining equation for advanced turbulence models, *AIAA J.* **26**, 1299 (1988).
  - [5] D. C. Wilcox *et al.*, *Turbulence modeling for CFD*, 3rd ed. (DCW Industries, La Canada, CA, 2006), Vol. 2.
  - [6] P. R. Spalart and S. R. Allmaras, A one-equation turbulence model for aerodynamic flows, in *30th Aerospace Sciences Meeting and Exhibit* (AIAA, Reston, VA, 1992), p. 439.
  - [7] F. R. Menter, Two-equation eddy-viscosity turbulence models for engineering applications, *AIAA J.* **32**, 1598 (1994).
  - [8] B. Launder, G. J. Reece, and W. Rodi, Progress in the development of a Reynolds-stress turbulence closure, *J. Fluid Mech.* **68**, 537 (1975).

- [9] S. Wallin and A. V. Johansson, An explicit algebraic Reynolds stress model for incompressible and compressible turbulent flows, *J. Fluid Mech.* **403**, 89 (2000).
- [10] E. Dow and Q. Wang, Quantification of structural uncertainties in the  $k$ - $\omega$  turbulence model, in *52nd AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics and Materials Conference*, Denver, 2011 (AIAA, Reston, VA, 2011), paper 2011-1762.
- [11] E. J. Parish and K. Duraisamy, A paradigm for data-driven predictive modeling using field inversion and machine learning, *J. Comput. Phys.* **305**, 758 (2016).
- [12] J. Ling and J. Templeton, Evaluation of machine learning algorithms for prediction of regions of high Reynolds averaged Navier Stokes uncertainty, *Phys. Fluids* **27**, 085103 (2015).
- [13] H. Xiao, J.-L. Wu, J.-X. Wang, R. Sun, and C. Roy, Quantifying and reducing model-form uncertainties in Reynolds-averaged Navier–Stokes simulations: A data-driven, physics-informed Bayesian approach, *J. Comput. Phys.* **324**, 115 (2016).
- [14] J.-X. Wang and H. Xiao, Data-driven CFD modeling of turbulent flows through complex structures, *Int. J. Heat Fluid Flow* **62**, 138 (2016).
- [15] G. Iungo, F. Viola, U. Ciri, M. Rotea, and S. Leonardi, Data-driven RANS for simulations of large wind farms, *J. Phys. Conf. Ser.* **625**, 012025 (2015).
- [16] A. P. Singh and K. Duraisamy, Using field inversion to quantify functional errors in turbulence closures, *Phys. Fluids* **28**, 045110 (2016).
- [17] K. Duraisamy, Z. J. Zhang, and A. P. Singh, New approaches in turbulence and transition modeling using data-driven techniques, in *53rd AIAA Aerospace Sciences Meeting*, Kissimmee, Florida, 2015 (AIAA, Reston, VA, 2015), paper 2015-1284.
- [18] J.-L. Wu, J.-X. Wang, and H. Xiao, A Bayesian calibration-prediction method for reducing model-form uncertainties with application in RANS simulations, *Flow Turbul. Combust.* **97**, 761 (2016).
- [19] J. Ling, A. Ruiz, G. Lacaze, and J. Oefelein, Uncertainty analysis and data-driven model advances for a jet-in-crossflow, *J. Turbomach.* **139**, 021008 (2017).
- [20] J. Friedman, T. Hastie, and R. Tibshirani, *The Elements of Statistical Learning* (Springer, Berlin, 2001).
- [21] T. Oliver and R. Moser, Uncertainty Quantification for RANS Turbulence Model Predictions, in 62nd Annual Meeting of the APS Division of Fluid Dynamics, Minneapolis, MN, November 22–24, 2009.
- [22] S. Pope, *Turbulent Flows* (Cambridge University Press, Cambridge, UK, 2000).
- [23] C. E. Rasmussen, *Gaussian Processes for Machine Learning* (MIT Press, Boston, 2006).
- [24] L. Breiman, Random forests, *Mach. Learn.* **45**, 5 (2001).
- [25] P. R. Spalart, Strategies for turbulence modelling and simulations, *Int. J. Heat Fluid Flow* **21**, 252 (2000).
- [26] J. C. Hunt, A. A. Wray, and P. Moin, Eddies, streams, and convergence zones in turbulent flows, Center for Turbulence Research, Stanford University, Report No. 19890015184, 1988 (unpublished).
- [27] P. Chakraborty, S. Balachandar, and R. J. Adrian, On the relationships between local vortex identification schemes, *J. Fluid Mech.* **535**, 189 (2005).
- [28] C. Górlé, J. Larsson, M. Emory, and G. Iaccarino, The deviation from parallel shear flow as an indicator of linear eddy-viscosity model inaccuracy, *Phys. Fluids* **26**, 051702 (2014).
- [29] J. Ling, R. Jones, and J. Templeton, Machine learning strategies for systems with invariance properties, *J. Comput. Phys.* **318**, 22 (2016).
- [30] J. L. Lumley and G. R. Newman, The return to isotropy of homogeneous turbulence, *J. Fluid Mech.* **82**, 161 (1977).
- [31] S. Banerjee, R. Krahl, F. Durst, and C. Zenger, Presentation of anisotropy properties of turbulence, invariants versus eigenvalue approaches, *J. Turbul.* **8**, N32 (2007).
- [32] M. Emory, R. Pecnik, and G. Iaccarino, Modeling structural uncertainties in Reynolds-averaged computations of shock/boundary layer interactions, in *49th AIAA Aerospace Sciences Meeting including the New Horizons Forum and Aerospace Exposition*, Orlando, 2011 (AIAA, Reston, VA, 2011), paper 2011-0479.
- [33] M. Emory and G. Iaccarino, Componentality-based wall-blocking for RANS models, in *Annual Research Briefs* (Center for Turbulence Research, California, 2014), pp. 193–208.
- [34] H. Goldstein, *Classical Mechanics* (Addison-Wesley, Reading, MA, 1980), pp. 143–148.

- [35] B. Tracey, K. Duraisamy, and J. Alonso, Application of supervised learning to quantify uncertainties in turbulence and combustion modeling, in *51st AIAA Aerospace Sciences Meeting Including the New Horizons Forum and Aerospace Exposition*, Grapevine, TX, 2013 (AIAA, Reston, VA, 2013), paper 2013–0259.
- [36] N. S. Altman, An introduction to kernel and nearest-neighbor nonparametric regression, *Am. Stat.* **46**, 175 (1992).
- [37] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning*, Springer Texts in Statistics Vol. 112 (Springer, Berlin, 2013).
- [38] J. A. Anderson, *An Introduction to Neural Networks* (MIT Press, Cambridge, MA, 1995).
- [39] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods* (Cambridge University Press, Cambridge, MA, 2000).
- [40] T. K. Ho, The random subspace method for constructing decision forests, *IEEE Trans. Pattern Anal. Mach. Intell.* **20**, 832 (1998).
- [41] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and Regression Trees* (CRC, Boca Raton, FL, 1984).
- [42] W. Jones and B. Launder, The prediction of laminarization with a two-equation model of turbulence, *Int. J. Heat Mass Transfer* **15**, 301 (1972).
- [43] P. Bradshaw, Turbulent secondary flows, *Annu. Rev. Fluid Mech.* **19**, 53 (1987).
- [44] F. Billard, Development of a robust elliptic-blending turbulence model for near-wall, separated and buoyant flows, Ph.D. thesis, School of Mechanical Aerospace and Civil Engineering, The University of Manchester, Manchester, UK, 2011.
- [45] 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).
- [46] A. Pinelli, M. Uhlmann, A. Sekimoto, and G. Kawahara, Reynolds number dependence of mean flow structure in square duct turbulence, *J. Fluid Mech.* **644**, 107 (2010).
- [47] M. Gibson and B. Launder, Ground effects on pressure fluctuations in the atmospheric boundary layer, *J. Fluid Mech.* **86**, 491 (1978).
- [48] 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).
- [49] C. Maaß and U. Schumann, Direct numerical simulation of separated turbulent flow over a wavy boundary, in *Flow Simulation with High-Performance Computers II* (Springer, Berlin, 1996), pp. 227–241.
- [50] Y. Bentaléb, S. Lardeau, and M. A. Leschziner, Large-eddy simulation of turbulent boundary layer separation from a rounded step, *J. Turbul.* **13**, N4 (2012).
- [51] H. Xiao, J.-L. Wu, J.-X. Wang, and E. G. Paterson, Physics-informed machine learning for predictive turbulence modeling: Progress and perspectives, in *55th AIAA Aerospace Sciences Meeting*, Grapevine, TX, 2017 (AIAA, Reston, VA, 2017), paper 2017–1712.
- [52] A. Dobra and J. Gehrke, Bias correction in classification tree construction, in *Proceedings of the Eighteenth International Conference on Machine Learning* (Morgan Kaufmann, San Francisco, 2001), pp. 90–97.
- [53] C. Strobl, A.-L. Boulesteix, A. Zeileis, and T. Hothorn, Bias in random forest variable importance measures: Illustrations, sources and a solution, *BMC Bioinform.* **8**, 25 (2007).
- [54] M. Sandri and P. Zuccolotto, A bias correction algorithm for the Gini variable importance measure in classification trees, *J. Comput. Graph. Stat.* **17**, 611 (2008).
- [55] A. Altmann, L. Toloşi, O. Sander, and T. Lengauer, Permutation importance: A corrected feature importance measure, *Bioinformatics* **26**, 1340 (2010).
- [56] J.-L. Wu, J.-X. Wang, H. Xiao, and J. Ling, A priori assessment of prediction confidence for data-driven turbulence modeling, *Flow Turbul. Combust.* (to be published) [[arXiv:1607.04563](https://arxiv.org/abs/1607.04563)].
- [57] J.-L. Wu, J.-X. Wang, H. Xiao, and J. Ling, Visualization of high dimensional turbulence simulation data using t-SNE, in *19th AIAA Non-Deterministic Approaches Conference*, Grapevine, TX, 2017 (AIAA, Reston, VA, 2017), paper 2017–1770.
- [58] J.-X. Wang, J. Wu, J. Ling, G. Iaccarino, and H. Xiao, A comprehensive physics-informed machine learning framework for predictive turbulence modeling, [arXiv:1701.07102](https://arxiv.org/abs/1701.07102).