# Uncertainty Quantification for Data-driven Turbulence Modelling with Mondrian Forests

Ashley Scillitoe[a,*], Pranay Seshadri[a,b], Mark Girolami[a,c]

[a]*Data-Centric Engineering, The Alan Turing Institute, London, UK*
[b]*Department of Mathematics, Imperial College London, London, UK*
[c]*Department of Engineering, University of Cambridge, Cambridge, UK*

## Abstract

Data-driven turbulence modelling approaches are gaining increasing interest from the CFD community. Such approaches generally aim to improve the modelled Reynolds stresses by leveraging data from high fidelity turbulence resolving simulations. However, the introduction of a machine learning (ML) model introduces a new source of uncertainty, the ML model itself. Quantification of this uncertainty is essential since the predictive capability of a data-driven model diminishes when predicting physics not seen during training. In this work, we explore the suitability of Mondrian forests (MF's) for data-driven turbulence modelling. MF's are claimed to possess many of the advantages of the commonly used random forest (RF) machine learning algorithm, whilst offering principled uncertainty estimates. An example test case is constructed, with a turbulence anisotropy constant derived from high fidelity turbulence resolving simulations. A number of flows at several Reynolds numbers are used for training and testing. MF predictions are found to be superior to those obtained from a linear and non-linear eddy viscosity model. Shapley values, borrowed from game theory, are used to interpret the MF predictions. Predictive uncertainty is found to be large in regions where the training data is not representative. Additionally, the MF predictive uncertainty is found to exhibit stronger correlation with predictive errors compared to an a priori statistical distance measure, which indicates it is a better measure of prediction confidence. The MF predictive uncertainty is also found to be better calibrated and less computationally costly than the uncertainty estimated from applying jackknifing to random forest predictions. Finally, Mondrian forests are used to predict the Reynolds discrepancies in a convergent-divergent channel, which are subsequently propagated through a modified CFD solver. The resulting flowfield predictions are in close agreement with the high fidelity data. A procedure for sampling the Mondrian forests' uncertainties is introduced. Propagating these samples enables quantification of the uncertainty in quantities of interest such as velocity or a drag coefficient, due to the uncertainty in the Mondrian forests' predictions. This work suggests that uncertainty quantification can be incorporated into existing data-driven turbulence modelling frameworks by replacing random forests with Mondrian forests. This would also open up the possibility of online learning, whereby new training data could be added without having to retrain the Mondrian forests.

*Keywords:* uncertainty quantification, supervised machine learning, turbulence modelling, dataset shift, Mondrian forests, machine learning interpretability.

## 1. Introduction

Turbulence is a key characteristic of fluid flows across many different industries. For aircraft, delaying the transition to turbulence over the wing surfaces can reduce drag [1], while large turbulent structures can increase the effectiveness of cooling systems in aircraft engines [2], both of which lead to reduced fuel

---

consumption. For wind turbines, turbulence can increase the power output, at the expense of increased fatigue loading [3]. In oil refineries, turbulence plays a key role in critical processes such as fluid catalytic cracking [4]. These examples, as well as many others, indicate the importance of being able to accurately predict the effects of turbulence for a range of industrial flows.

The characterisation of turbulence is challenging due to its chaotic nature, and the broad range of spatio-temporal scales involved. The continued growth of computing power has enabled the unsteady Navier-Stokes equations to be computed directly, so that all the scales of turbulence are simulated. Unfortunately, the extreme computational cost of these direct numerical simulations (DNS) limits their application to relatively simple flows. To mitigate the high computational cost, large eddy simulation (LES) techniques only simulate the larger energy containing turbulent eddies, whilst the smaller unresolved turbulent scales are modelled. LES type approaches are gaining popularity in many industrial applications, for example for turbo-machinery [5], wind turbines [6], and urban flows [7]. However, due to their high cost, such simulations are often simplified representations of real industrial flows. As discussed by Tucker and DeBonis [2], even with the ever-increasing computing power available, LES methods are unlikely to be routinely used in the design process of complex engineering systems over the next decade and beyond.

Instead, computations based on Reynolds-averaged Navier–Stokes (RANS) models are still the workhorse for predicting turbulent flows in many industries. RANS computations are orders of magnitude cheaper than LES techniques. However, RANS models are known to perform poorly in many flows of engineering relevance, including those with swirl, pressure gradients, and streamline curvature [8, 9]. Many researchers have attempted to use LES or DNS to better understand the physics of turbulence in various flows, with the aim of developing better turbulence models. As summarised in the recent review paper by Duraisamy et al. [10], data-driven turbulence modelling is emerging as a promising way to inform turbulence models with data in a more systematic way.

Many of the data-driven turbulence modelling strategies fall under what is known as supervised machine learning (ML). This involves learning a function based on *training data* consisting of a set of input-output pairs (see Sec. 2). The learned function can then be used to make predictions on the *test data*. Ling and Templeton [11] were one of the first to apply ML to turbulence modelling, using a random forest classifier to predict when RANS modelling assumptions would fail. Ling et al. [12] further used neural networks to predict Reynolds stress anisotropy. Singh et al. [13] used field inversion to determine functional discrepancies in existing RANS models, which are then reconstructed as functions of local flow features using the Adaboost ML algorithm. Wu et al. [14] investigated the use of random forests to predict the discrepancies of RANS modelled Reynolds stresses in separated flows. Kaandorp and Dwight [15] pursue similar lines, but they modify the random forest algorithm to accept a tensor basis in order to guarantee Galilean invariance.

These studies demonstrate the growing interest in applying ML techniques to turbulence modelling. However, for such methods to be employed in industry more trust is needed in the ML predictions. All ML algorithms have an inductive bias, which is the assumptions the learner uses to predict outputs given inputs that it has not encountered. The *no free lunch* theorem [16] states that if a bias is correct on some cases, it must be incorrect on equally many cases, thus in the context of data-driven turbulence modelling it is unlikely for one learner to be able to generalise to all unseen flows. This is evidenced in [11] and [14], where the authors report that their data-driven closures performed poorly on flows that were significantly different from the ones on which they were trained. It is desirable for the learning algorithm to be able to generalise well in situations with a significant *dataset shift*, where the test and training data have significantly different distributions. Yet, a more realistic and equally important requirement is that the algorithm can at least quantify the uncertainty in its own predictions. Then the user can at least know when the training data is suitable, and whether the ML predictions can be trusted. This is a pertinent issue for the field of data-driven turbulence modelling, since many of the available LES and DNS datasets are on simplified geometries at low Reynolds numbers[1]. But, we wish to know whether we trust ML models trained on these flows to make predictions on the more complex higher Reynolds number flows seen in industry.

The random forest algorithm is commonly used for turbulence modelling applications [11, 14, 15] for its robustness and high predictive accuracy. Jackknife re-sampling can be used to provide uncertainty estimates

---

[1]As noted in [9], the grid resolution required for DNS and LES increases dramatically with Reynolds number.

for random forest predictions [17]. However, these estimates do not correctly account for the uncertainty arising from differences between the training and test data. Ling and Templeton [11] and later Wu et al. [18] proposed the use of statistical distance metrics to measure co-variate shift between the test and training datasets. Although informative, these metrics don't provide uncertainty estimates for the ML predictions. A number of authors have explored the use of ML techniques with inbuilt uncertainty estimates. Parish and Duraisamy [19] used Gaussian processes (GPs) to predict the turbulence intermittency and correction terms for the turbulence transport equations. To predict the turbulence anisotropy tensor Geneva and Zabaras [20] used Bayesian neural networks (BNN), while Blauw and Dwight [21] used Bayesian additive regression trees (BART). All three approaches show promise, but also have disadvantages. GPs are known to deliver high quality uncertainty estimates, however they can be challenging to scale to large datasets since their cost scales cubically with dataset size [22]. Meanwhile, neural networks can be challenging and costly to train, and Bayesian neural networks are particularly computationally expensive [23]. For BART, Metropolis-Hastings type samplers commonly used to perform inference can struggle with large high dimensional datasets [24].

An alternative supervised ML technique is the Mondrian forest (MF), first conceived for classification [25], and extended to regression by Lakshminarayanan et al. [26]. This relatively new technique was demonstrated to be competitive in computational cost and accuracy to popular tree ensemble approaches such as random forests. Unlike random forests, the trees in a Mondrian forest are probabilistic; they use a hierarchical Gaussian prior, and the posterior parameters are efficiently computed using Gaussian belief propagation. Lakshminarayanan et al. claim that this endows MF's with principled uncertainty estimates, more akin to those of a GP, with predictions smoothly returning to the prior and exhibiting higher uncertainty far from training data. To the authors' knowledge the use of Mondrian forests for data-driven turbulence modelling is yet to be explored. The present paper aims to explore this, to understand whether Mondrian forests offer useful uncertainty estimates in addition to accurate predictions in data-driven turbulence frameworks. The paper is structured as follows: In Section 2, commonly used frameworks of supervised ML for turbulence modelling are introduced; the Mondrian forest algorithm is then introduced in Section 3, along with the procedures for generating training data and for propagating ML predictions through a CFD solver. In Section 4 random and Mondrian forests are used to predict a turbulence field variable, in order to understand their relative performances in more detail. Finally, in Section 5 an example of the use of Mondrian forests for a data-driven turbulence modelling task is explored.

## 2. Supervised machine learning for turbulence modelling

### 2.1. RANS modelling uncertainties

The difficulty of RANS modelling arises from the fundamental closure problem that is introduced when the Navier–Stokes equations are averaged with respect to time. The incompressible RANS momentum equation in the $i^{th}$ direction is

$$\langle u_j \rangle \frac{\partial \langle u_i \rangle}{\partial z_j} = \frac{\partial}{\partial z_j} \left[ -\frac{\langle p \rangle}{\rho} \delta_{ij} + \nu \left( \frac{\partial \langle u_i \rangle}{\partial z_j} + \frac{\partial \langle u_j \rangle}{\partial z_i} \right) + \frac{\tau_{ij}}{\rho} \right] + \langle g_i \rangle, \tag{1}$$

where $\langle \cdot \rangle$ indicates a time-averaged quantity and $\boldsymbol{z} = (z_1, z_2, z_3)$ is the spatial coordinate vector. The RANS equations are unclosed, with the turbulent (or Reynolds) stress term, $\tau_{ij}$, unknown. It is common for the Boussinesq linear eddy-viscosity hypothesis to be used

$$\tau_{ij} = -\rho \langle u'_i u'_j \rangle = \nu_t \left( \frac{\partial \langle u_i \rangle}{\partial z_j} + \frac{\partial \langle u_j \rangle}{\partial z_i} \right) - \frac{2}{3} k \delta_{ij}, \tag{2}$$

where $k = \langle u'_i u'_i \rangle / 2$ is the turbulent kinetic energy. However, the turbulent eddy viscosity $\nu_t$ is still unknown, and RANS models must be used to estimate this term. Xiao and Cinnella [27] divide the uncertainties arising from RANS modelling into two categories: (i) structural (or model-form) uncertainties, which arise from the structure of the RANS equations solved to provide $\nu_t$, or from the eddy-viscosity assumption in (2), and (ii) parametric uncertainties, which are caused by a lack of generalisation of closure coefficients in the RANS model. As discussed in the following section, machine learning can be used to inform both categories of uncertainties.

3

## 2.2. Machine learning to inform RANS

This paper deals only with the general framework of supervised machine learning. In this context, we define $\mathbf{x} \in \mathbb{R}^d$ to be a vector of $d$ non-dimensional flow-features obtained at a given location in a flow-field. These flow features, described further in Section 3.3, are intended to uniquely describe the RANS flow-field. Let $y \in \mathbb{R}$ be a scalar-valued output quantity of interest also obtained at each location within the flow-field. Assuming a flow-field is specified by $N$ spatial locations, the following matrix-vector notation is adopted to characterise the data

$$\boldsymbol{X}_N = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}, \quad \mathbf{y}_N = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \tag{3}$$

where the subscript denotes the $N$ *training samples*. Analogously, predictions of the spatially-varying scalar quantity of interest for $M$ different feature vectors are given by

$$\tilde{\boldsymbol{X}}_M = \begin{bmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_M^T \end{bmatrix}, \quad \tilde{\mathbf{y}}_M = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_M \end{bmatrix}, \tag{4}$$

where the superscript $\tilde{\cdot}$ indicates *testing samples*. These definitions will be important for the supervised ML exposition.

An example of supervised ML in the context of turbulence modelling is shown in Figure 1. A number of DNS and LES datasets are obtained, and since only well-validated LES data is used it is assumed to be well resolved. For the remainder of this paper the LES and DNS will both be referred to as near direct simulations (NDS). For each NDS dataset, a companion RANS solution is obtained with the same boundary conditions and geometry. Each NDS and RANS data pair are then pre-processed to form the training data set $\boldsymbol{X}_N, \mathbf{y}_N$, where $N$ samples are taken at matching locations between the RANS and NDS flows. It is the task of the ML algorithm to learn a functional mapping $f_{ML}$ between the training inputs $\boldsymbol{X}_N$ and training outputs $\mathbf{y}_N$, so that $\tilde{\mathbf{y}}_M$ can be predicted for a new *test* flow where only the RANS data inputs $\tilde{\boldsymbol{X}}_M$ is available. An important point to note here is that the NDS data is taken as the *truth* data, hence errors in this data are not accounted for in the uncertainty estimates discussed later on. Such errors are assumed to be small in the present work, since only well-validated NDS datasets are selected. However, numerous errors do exist in LES and DNS simulations [28], and accounting for these in data-driven modelling frameworks may be an important area of future research.
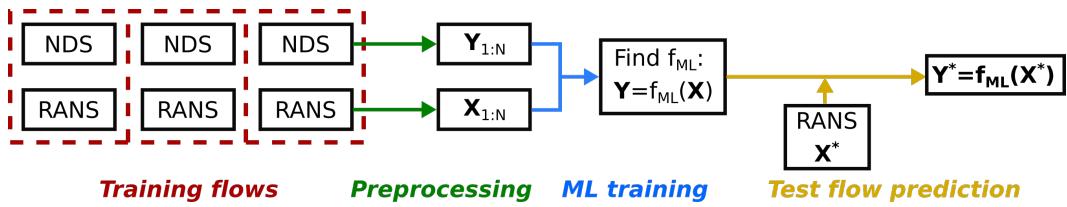


Figure 1: An example framework using supervised machine learning data-driven turbulence modelling.

Various authors have used the aforementioned framework to predict different turbulent quantities. Ling and Templeton [11] used a machine learning classifier, where the response $\mathbf{y}$ is a binary label which indicates whether the RANS model can be trusted in different flow regions. For example, since the Boussinesq linear eddy-viscosity assumption (Eq. 2) presumes a positive eddy-viscosity $\nu_t$, the response

$$y = \begin{cases} 0, & \text{if } \nu_t < 0 \\ 1, & \text{otherwise} \end{cases} \tag{5}$$

can be used to suggest regions where (2) can not be trusted. More recently, many authors [13, 19, 12, 15, 20, 14] have used machine learning regressors to predict the Reynolds stress discrepancy $\Delta\boldsymbol{\tau} = \boldsymbol{\tau}^{RANS} - \boldsymbol{\tau}^{NDS}$,

which can be used to correct $\boldsymbol{\tau}$ before it is injected back into a CFD solver. The raw discrepancy $\Delta\boldsymbol{\tau}$ is not Galilean invariant, so it is common to instead perform the eigenspace decomposition

$$\boldsymbol{\tau} = 2k\left(\frac{1}{3}\mathbf{I} + \mathbf{B}\right) = 2k\left(\frac{1}{3}\mathbf{I} + \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^T\right), \tag{6}$$

where $\mathbf{I}$ is the second order identity tensor, $\mathbf{B}$ is the normalised turbulent anisotropy tensor, and

$$\boldsymbol{V} = \begin{bmatrix} | & | & | \\ v_1 & v_2 & v_3 \\ | & | & | \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix}, \tag{7}$$

are the eigenvectors and eigenvalues of $\mathbf{B}$, respectively, indicating its shape and orientation. The eigenvalues are sorted such that $\lambda_1 \geq \lambda_2 \geq \lambda_3$. These eigenvalues can be transformed to the Barycentric coordinates $C_1$, $C_2$, and $C_3$, and then to Cartesian coordinates $\xi$ and $\eta$. Wu et al. [14] use machine learning to predict the discrepancy $\Delta\boldsymbol{\tau} = (\Delta\log k, \Delta\xi, \Delta\eta, h_1, h_2, h_3)$, where $h_1$, $h_2$ and $h_3$ are the first three unit quaternions used to represent the transformation from the RANS eigenvectors $\mathbf{V}^{RANS}$ to the target eigenvectors $\mathbf{V}^{NDS}$.

The prediction of $\Delta\boldsymbol{\tau}$ accounts for the overall uncertainty due to the RANS model. However, Wu et al. [29] show that the injection of Reynolds stresses directly into a CFD solver can sometimes lead to the RANS equations becoming ill-conditioned. Instead, field inversion can be used [19, 13] to find a field parameter $\beta(\boldsymbol{z})$ which improves the RANS model's predictions. For example, Singh et al. [13] use $\beta(\boldsymbol{z})$ to scale the production term in the $\omega$ equation of the $k$-$\omega$ RANS model:

$$\frac{D\omega}{Dt} = \beta(\boldsymbol{z})P\left(k, \omega, \langle\mathbf{u}\rangle\right) - D\left(k, \omega, \langle\mathbf{u}\rangle\right) + T\left(k, \omega, \langle\mathbf{u}\rangle\right) \tag{8}$$

This field parameter accounts only for the parametric uncertainties in the RANS model. The parameter is problem-specific since it is a function of the problem's coordinates $\boldsymbol{z}$, but machine learning can then be used to infer the relationship between $\beta(\boldsymbol{z})$ and the flow features $\boldsymbol{X}$.

## 3. Data-driven framework

In this section, the data-driven framework we explore in this paper is introduced. A python package implementing the framework is available from github.com/ascillitoe/mondrian_turbulence.

### 3.1. Mondrian forests

This paper explores the possibility of using Mondrian forests to replace the random forests used as the functional mapping $f_{ML}$ in many data-driven turbulence modelling frameworks, such as that proposed by Wu et al. [14]. Random forests [30] are a popular supervised machine learning algorithm, which can often achieve good accuracy even with minimal tuning of their hyper-parameters. They consist of an ensemble of decision trees, akin to the one shown in Figure 2a. At each node in the tree, the data is split in a binary fashion, with the split locations chosen to optimise an appropriate criterion, such as mean squared error (MSE). Decision trees have low bias, but high variance (they over-fit to training data). By taking an average of each tree's prediction, random forests are able to reduce this variance.

Like random forests, a Mondrian forest [25, 26] is also an ensemble of trees. However, in this case, the trees are Mondrian trees. Figure 2b visualises a Mondrian tree, which is a restriction of a Mondrian process to a finite set of points. A Mondrian tree can be thought of as a decision tree, but with a number of important differences:

(i) Splits are committed only within the range of training data, hence the split represented by an internal node $j$ holds only within $B_j^x$ and not the whole of $B_j$. The data block $B_j^x$, shown in Figure 2b, represents the smallest rectangle enclosing the input data $\boldsymbol{X}$ at node $j$:

$$B_j^x = \left(l_{j1}^x, u_{j1}^x\right] \times \cdots \times \left(l_{jd}^x, u_{jd}^x\right], \tag{9}$$

where $l_{jd}^x$ and $u_{jd}^x$ denote the lower and upper bounds of the training data points in node $j$ along dimension $d$.
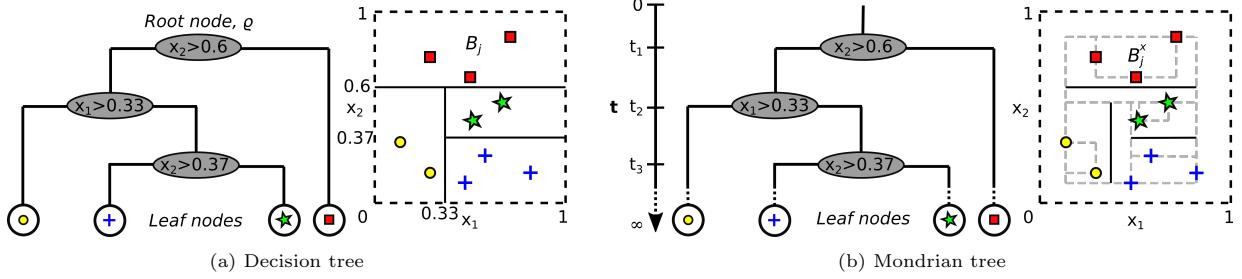
(a) Decision tree

(b) Mondrian tree

Figure 2: Example of classification trees trained on data with four classes over eight data points in $[0,1]^2$. Trees are shown on the left, and the partitions in feature space are shown on the right. Regression trees are similar, except the leaf nodes contain continuous values instead of classes.

(ii) The CART algorithm [31], often used to build decision trees, samples splits based on the MSE of the predicted response data $\mathbf{y}$. Whereas Mondrian tree splits are sampled independently of $\mathbf{y}$.

Unlike standard decision trees, Mondrian trees are probabilistic models, which determine $p_{\mathcal{T}}(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y}))$. This is a predictive distribution for the response $\tilde{y}$, tree $\mathcal{T}$, training dataset $(\boldsymbol{X}, \mathbf{y})$ and the test point $\tilde{\mathbf{x}}$. As with other Bayesian approaches, a *prior* must be chosen, and the *posterior* is obtained via marginalisation. The prediction of a Mondrian forest is then the average prediction from the $L$ number of Mondrian trees

$$p(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y})) = \frac{1}{L} \sum_{j}^{L} p_j(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y})). \tag{10}$$

This predictive posterior over $\tilde{y}$ is a mixture of Gaussians, and it is straightforward to calculate the predictive mean and variance from this. The predictive variance is said to provide principled uncertainty estimates. The bounded nature of the splits in Mondrian trees, noted in (i) above, means that Mondrian trees (and forests) exhibit higher uncertainty as $\tilde{\mathbf{x}}$ moves further away from the training data $(\boldsymbol{X}, \mathbf{y})$.

To obtain uncertainty estimates for random forests, Wager et al. [17] propose the use of jackknife re-sampling to provide confidence intervals for random forest predictions (see Appendix A.2.1 for more details). However, preliminary investigations (Appendix B.3) into using re-sampling for turbulence field predictions suggested the computational cost is prohibitively high, hence they are not explored in detail in this paper. For a more detailed exposition of Mondrian forests and random forests, see Appendix A.

*3.2. Database of turbulent flows*

Training and testing data are generated from the ten turbulent flow cases outlined in Table 1 and visualised in Figure 3. Each case is a well-validated near direct simulation (NDS), with the majority obtained from publicly accessible databases[2]. For each NDS case a companion RANS dataset is generated using the SU2 open source CFD solver (version 7.0.6) available from `github.com/su2code/SU2`. For all cases the incompressible solver is used, which spatially discretises the incompressible Navier-Stokes equations using a finite volume method on unstructured grids, and solves them in a coupled manner with a custom preconditioning approach [32]. For all cases, grid refinement is carried out to ensure mesh independence, and implicit time-stepping is used with an FGMRES linear solver. For turbulence modelling, the SST-2003 model [33] is used, which solves two equations for the turbulent kinetic energy $k$ and specific dissipation rate $\omega$

$$\frac{\partial k}{\partial t} + \frac{\partial \langle u_i \rangle k}{\partial z_i} = \frac{1}{\rho}\tilde{P}_k - \beta^* k\omega + \frac{\partial}{\partial z_i}\left[(\nu + \sigma_k \nu_t)\frac{\partial k}{\partial z_i}\right] \tag{11}$$

---

[2]`turbmodels.larc.nasa.gov`, `turbase.cineca.it` and `flow.kth.se/flow-database`.

$$\frac{\partial \omega}{\partial t} + \frac{\partial \langle u_i \rangle \omega}{\partial z_i} = \frac{1}{\rho} P_\omega - \beta \omega^2 + \frac{\partial}{\partial z_i} \left[ (\nu + \sigma_\omega \nu_t) \frac{\partial \omega}{\partial z_i} \right] + 2 (1 - F_1) \sigma_{\omega 2} \frac{1}{\omega} \frac{\partial k}{\partial z_i} \frac{\partial \omega}{\partial z_i} \qquad (12)$$

The production of $k$ term in (11) is limited by $\tilde{P}_k = \min\left(P_k, 20\beta^* \rho k \omega\right)$, where the raw production term is obtained from the Boussinesq hypothesis

$$P_k = \mu_t \frac{\partial \langle u_i \rangle}{\partial z_j} \left( \frac{\partial \langle u_i \rangle}{\partial z_j} + \frac{\partial \langle u_j \rangle}{\partial z_i} \right). \qquad (13)$$

The production term in the $\omega$ equation is then given by $P_\omega = \rho \alpha \tilde{P}_k / \mu_t$, where $\alpha$ is a term used to blend between the $k - \epsilon$ and $k - \omega$ models in the SST model (see [33] for more details).

For cases 1 and 3 a 1D inlet velocity profile is set at the inlet to match the NDS data, and a uniform static pressure outlet boundary condition is enforced to match the area-averaged static pressure from the NDS at the same location. For Cases 2, 4 and 5 streamwise periodicity is enforced, with the streamwise forcing calibrated to match the NDS mass flow rate. All cases are assumed to be fully turbulent, with no transition model used. All RANS cases are 2D, and the NDS cases are span-wise averaged to match the RANS cases. A locally adaptive CFL number is used, with the minimum set at 2 and the maximum at 20.

Table 1: Summary of flow cases used for training and testing.

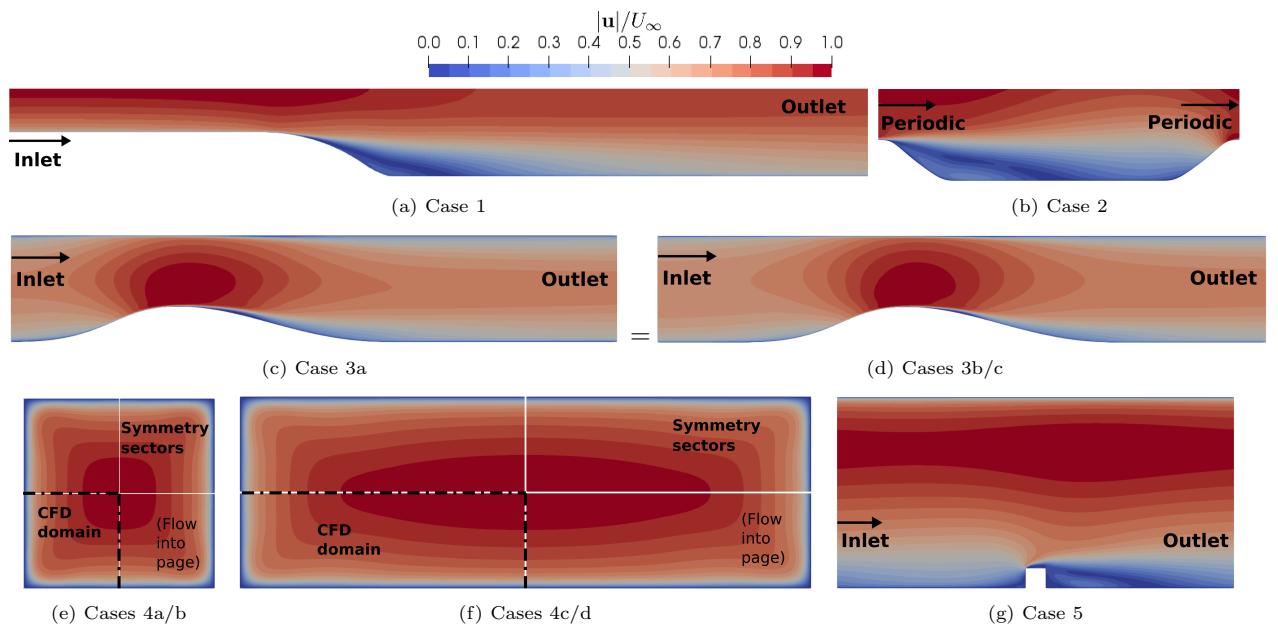| Case | Description | Re | LES/DNS | Ref. |
|------|-------------|-----|---------|------|
| 1 | Curved backwards step | $Re_\tau = 618$ | LES | [34] |
| 2 | Periodic hills | $Re_\tau = 160$ | LES | [35] |
| 3a/b/c | Convergent-divergent channel | $Re_\tau = 395/617/950$ | DNS | [36], [37] |
| 4a/b | Duct with aspect ratio $AR = 1$ | $Re_\tau = 180, 360$ | DNS | [38] |
| 4c/d | Duct with aspect ratio $AR = 3$ | $Re_\tau = 180, 360$ | DNS | [38] |
| 5 | Ribbed channel flow | $Re_\tau = 360$ | LES | [39] |



Figure 3: Contours of normalised velocity magnitude for the LES/DNS datasets listed in Table 1.

### 3.3. Pre-processing of CFD data

The RANS and NDS data is pre-processed according to the framework in Figure 1. To form the training dataset $\boldsymbol{X}_N, \mathbf{y}_N$, the mean flow data at each mesh point in the RANS cases is converted to a feature vector $\mathbf{x} \in \mathbb{R}^d$. To reduce the amount of data, the input data is filtered to remove non-unique data points (to the $4^{th}$ significant figure), and 20% of the points are randomly sub-sampled. After filtering and sub-sampling all ten cases together yield 25000 *observations*. At the same spatial locations the responses $\mathbf{y}$ are then calculated from the NDS data, with the quantity stored in $\mathbf{y}$ depending on what we wish to predict with the framework later on.

The $d$ flow features utilised in this paper are the same as in [14]. Following [11], the raw features $\alpha$ are non-dimensionalised by a local normalisation factor $\beta$

$$\hat{\alpha} = \frac{\alpha}{|\alpha| + |\beta|}, \tag{14}$$

where $\alpha$ and $\beta$ are given in Table 2. The strain-rate tensor $\boldsymbol{S}$ and rotation-rate tensor $\boldsymbol{\Omega}$ are included as it has been recognised that all algebraic Reynolds stress and eddy viscosity models can be written in the general form $\boldsymbol{\tau} = f(\boldsymbol{S}, \boldsymbol{\Omega})$ [40]. The pressure gradient is included since it is known that turbulence is also influenced by this, while the turbulent kinetic energy (TKE) gradient is added in an attempt to account for non-local turbulent transport effects. Finally, near wall viscous effects have an important influence on turbulence so a non-dimensional wall distance feature is included, along with features to describe the length and time-scales of turbulence.

Table 2: Non-dimensional features used to represent the flow fields. Notations are as follows: $\mathbf{u}$ is the mean velocity vector, $k = \overline{u_i u_i}/2$ is the turbulence kinetic energy (TKE), $\omega$ is the specific turbulent dissipation rate, $|\cdot|$ and $||\cdot||$ indicate vector and matrix norm.

| Description | Raw input, $\alpha$ | Normalisation, $\beta$ |
|---|---|---|
| Strain-rate tensor | $\boldsymbol{S}$ | $\omega$ |
| Rotation-rate tensor | $\boldsymbol{\Omega}$ | $||\boldsymbol{\Omega}||$ |
| Pressure gradient | $\nabla p$ | $\rho|D\mathbf{u}/Dt|$ |
| TKE gradient | $\nabla k$ | $\omega\sqrt{k}$ |
| Wall-distance based Reynolds number, $Re_d$ | $\min\left(\frac{\sqrt{k}d}{50\nu}, 2\right)$ | - |
| Turbulence intensity, $Ti$ | $k$ | $\frac{3}{2}|\mathbf{u}|^2$ |
| Ratio of turbulent to mean strain time-scales, $t_t/t_s$ | $\omega$ | $\frac{1}{||\boldsymbol{S}||}$ |

To achieve Galilean and rotational invariance of the machine learning function, the approach of [14] is adopted. The normalised gradients $\widehat{\nabla p}$ and $\widehat{\nabla k}$ are transformed into the anti-symmetric tensors $\widehat{\mathbf{A}}_p$ and $\widehat{\mathbf{A}}_k$, and the minimal integrity bases for the tensorial set $\{\widehat{\boldsymbol{S}}, \widehat{\boldsymbol{\Omega}}, \widehat{\mathbf{A}}_p, \widehat{\mathbf{A}}_k\}$ are calculated. Each of the 47 invariant bases is the trace of a tensor such as $\widehat{\boldsymbol{S}}^2$ and $\widehat{\boldsymbol{\Omega}}\widehat{\mathbf{A}}_p\widehat{\mathbf{A}}_k\widehat{\boldsymbol{S}}$. Combining the remaining 3 scalar features in Table 2 with the invariant bases leads to a 50 dimensional feature space $\mathbf{x} \in \mathbb{R}^{50}$. For further details see Appendix B in Ref. [14].

**Remark 1.** *By learning a mapping between $\mathbf{x}$ and $y$, we have assumed the Reynolds stress $\boldsymbol{\tau}$ can be described by purely local (point-wise) quantities $(k, \omega, \nu)$, and the local gradients $(\nabla\boldsymbol{u}, \nabla k, \nabla p)$ at a given point. This assumption is likely to be questionable in regions where the turbulence scales are large or transport of turbulence is strong. Additionally, since $y$ is derived from NDS data and $\mathbf{x}$ from RANS data, we have implicitly assumed the RANS and NDS mean velocity and pressure fields match. This assumption is likely to be invalid if more complex flow cases were used for training.*

As noted in point (ii) in 3.1, Mondrian forests choose splits independent of the response data $\mathbf{y}$. This means a Mondrian forest's predicitve accuracy can be degraded if features are present in $\boldsymbol{X}$ which are

irrelevant to the response variable $\mathbf{y}$. To check for the effect of irrelevant features, backward elimination was performed, with prediction accuracy monitored as the least important features were iteratively removed. This study was performed on Mondrian forest's trained to predict the $C_{aniso}$ constant, with further details given in Appendix B.2. Surprisingly, irrelevant features were not found to be an issue in this case, perhaps because all non-zero features are important and many of the 47 invariant terms are zero for two-dimensional flows. Nevertheless, a 12 dimensional subset of the original 50 features is used in the remainder of this paper:

$$\boldsymbol{x} = \left\{ \widehat{\boldsymbol{S}}^2, \widehat{\boldsymbol{\Omega}}^2, \widehat{\mathbf{A}}_p^2, \widehat{\mathbf{A}}_k^2, \widehat{\boldsymbol{\Omega}}\widehat{\mathbf{A}}_k, \widehat{\boldsymbol{\Omega}}\widehat{\mathbf{A}}_k\widehat{\boldsymbol{S}}^2, \widehat{\boldsymbol{\Omega}}^2\widehat{\mathbf{A}}_k\widehat{\boldsymbol{S}}, \widehat{\mathbf{A}}_k^2\widehat{\mathbf{A}}_p\widehat{\boldsymbol{S}}, \widehat{\mathbf{A}}_k^2\widehat{\boldsymbol{S}}\widehat{\mathbf{A}}_p\widehat{\boldsymbol{S}}^2, Re_d, \hat{k}, \hat{\omega} \right\}^T. \tag{15}$$

This reduced feature space allows for reduced training/prediction times and data storage requirements, with no appreciable loss in predictive accuracy. This pruning down of features via backward elimination constitutes a departure from the methodology pursued in [14].

### 3.4. Propagation of machine learning predictions

If one uses the ML framework to predict the discrepancy vector $\Delta\boldsymbol{\tau}$, it is desirable to propagate this through the RANS solver so that revised flow field predictions can be obtained. A number of authors propose strategies to achieve this, for example, see [12, 14, 15, 20]. We choose a variation of that proposed by Kaandorp and Dwight [15]. Here, to achieve favourable numerical convergence, the ML derived anisotropy tensor $\mathbf{B}_{ML}$ is under-relaxed against the original Boussinesq tensor $\mathbf{B}_{BL} := \nu_t\mathbf{S}$. We replace the Reynolds stress term in the open source SU2 CFD code with

$$\boldsymbol{\tau} = \frac{2}{3}k\mathbf{I} + 2k\left[(1 - \gamma_n)\mathbf{B}_{BL} + \gamma_n\mathbf{B}_{ML}\right], \tag{16}$$

where $\mathbf{B}_{ML}$ is reconstructed from the ML predicted eigenvalues and eigenvectors (recall Eqn. 6), with $\xi_{ML} = \xi_{RANS} + \Delta\xi$. The blending parameter $\gamma_n$ is defined with a linear ramp

$$\gamma_n = \gamma_{max} \min\left(1, \frac{n}{n_{max}}\right), \tag{17}$$

which iteratively converges to $n = n_{max}$, where $\gamma_n = \gamma_{max}$. After this a sufficient number of iterations are performed to achieve full convergence, i.e., $n > n_{max}$. As will be described later in Section 5, $n_{max}$ here is the user-defined number of iterations used for ramping up the blending of $\mathbf{B}_{ML}$. Further, Kaandorp and Dwight [15] replace the limited production term $\tilde{P}_k$ in (13) by the exact production term

$$P_k = -\tau_{jj}\frac{\partial \langle u_i \rangle}{\partial z_j}, \tag{18}$$

and the $\omega$ production term in (12) is then given by $P_\omega = \rho\alpha P_k/\mu_t$. Boundary conditions are easily satisfied in the above implementation since the anisotropy term in (16) is scaled by $k$, which equals zero at the wall for the low Reynolds number RANS model used here. If wall functions were to be used, a blending such as that proposed by Geneva and Zabaras [20] would be necessary. Since $k$ is allowed to vary here, we no longer require its discrepancy $\Delta \log k$, and so instead of requiring an ML prediction for the Reynolds stress discrepancy term $\Delta\boldsymbol{\tau} = (\Delta \log k, \Delta\xi, \Delta\eta, h_1, h_2, h_3)$, we only require a prediction for $\Delta\mathbf{B} = (\Delta\xi, \Delta\eta, h_1, h_2, h_3)$. The modified SU2 code described here is available at github.com/ascillitoe/SU2/tree/feature_DDRANS.

### 3.5. Propagation of Mondrian forest's predictive uncertainty

As mentioned in Section 3.1, the Mondrian forest framework outputs the mean and variance for the five quantities in $\mathbf{B}$. To propagate the variance estimates through the RANS solver, we need a way of generating spatially consistent samples that satisfy the mean and variances for these five quantities. To reiterate, these mean and variance values are spatially varying, and thus care must be taken to ensure that random samples adhering to these statistical quantities preserve the underpinning spatial correlations. The challenge, however, is that we do not have a parametric definition for this correlation.

9

One approach is to assume that a single discrepancy quantity, say, $\Delta k$ can be modelled as a spatial Gaussian random field, defined in terms of its spatial coordinates $\boldsymbol{z}$. Then, we use standard Gaussian process machinery to generate random realisations that adhere to the aforementioned mean $\boldsymbol{\mu}_{\Delta k}(\boldsymbol{z})$ and variance $\boldsymbol{\sigma}^2_{\Delta k}(\boldsymbol{z})$.

There are three key challenges that arise when doing so. First, one has to assume a certain covariance kernel. Here the squared exponential (radial basis function) is the default option, where the covariance function is set by the distance between two coordinates and not their spatial location within the domain. Second, within the Gaussian process framework, we need to invert (via Cholesky) a square matrix of dimension set by the number of mesh nodes, which can be prohibitive for large meshes. However, this can be abated by coarsening the flow-field. Third, and perhaps of greatest importance to highlight, is that many of our discrepancy fields have large variance values. Thus, when optimising using either maximum likelihood or Markov chain Monte Carlo, the Gaussian process model sought to explain the spatial distribution with a very flat mean and large variance posterior spatial model—violating the Mondrian forest yielded statistical output.

This motivated our heuristic, where we construct a custom covariance function of the form

$$\boldsymbol{K}_{\Delta k}(\boldsymbol{z}_i, \boldsymbol{z}_j) = \phi \sigma_{\boldsymbol{z}_i} \sigma_{\boldsymbol{z}_j}, \quad \text{with} \quad \phi = \frac{1}{\omega_1 + \omega_2 \exp\left(\|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2^2 - \omega_3\right)}. \tag{19}$$

Assuming the values of constants $\omega_1, \omega_2$ and $\omega_3$ are selected such that

$$0 \leq \phi \leq 1 \tag{20}$$

holds, (19) explicitly sets $\phi$ to be the correlation between any $\boldsymbol{z}_i$ and $\boldsymbol{z}_j$, whilst preserving the underlying variances. Although the subscript notation in (19) is shown for $\Delta k$'s covariance function, we use the same expression for all the remaining five spatial discrepancy fields. Random spatial samples for say $\Delta k$ can then be generated by sampling from the multivariate normal distribution identified by

$$\mathcal{N}\left(\boldsymbol{\mu}_{\Delta k}(\boldsymbol{z}), \boldsymbol{K}_{\Delta k}(\boldsymbol{z}, \boldsymbol{z}')\right). \tag{21}$$

We iterated over different values of the constants ensuring that (20) holds and identified $\omega_1 = 0.58$, $\omega_2 = 1.4$ and $\omega_3 = 1.2$ to offer sufficient spatial smoothing. It should be noted that these values yield a covariance matrix that is positive definite, which is another constraint governing these constants.

## 4. Predicting turbulence variables with Mondrian forests

Prior to examining the use of Mondrian forests for a data-driven turbulence model application, we compare them to random forests on a simple problem. The training response $\mathbf{y}$ is set to be the turbulence anisotropy constant proposed by Banerjee et al. [41],

$$C_{aniso} = -3\lambda_3, \tag{22}$$

where $\lambda_3$ is calculated from the NDS Reynolds stress fields. The constant $C_{aniso}$ describes the turbulent anisotropy, with $C_{aniso} = 0$ indicating isotropic turbulence and $C_{aniso} = 1$ indicating fully one component turbulence.

### 4.1. Convergent-divergent channel

The turbulence anisotropy constant $C_{aniso}$ obtained from the DNS data [36] for the convergent-divergent channel at $Re_\tau = 617$ is shown in Figure 4. As expected, the turbulence anisotropy is generally high near the upper and lower viscous walls. Near the lower wall, downstream of the bump ($5.5 < z_1 < 7$), a strong adverse pressure gradient leads to flow separation and a small recirculation region. The production of $\langle u_1' u_1' \rangle$ is an order of magnitude larger than $\langle u_2' u_2' \rangle$ in this region [36], leading to highly anisotropic turbulence. Near the upper surface at the same horizontal location, a weaker adverse pressure gradient is observed, implying that the flow is on the verge of separation, and turbulence anisotropy is also high. The turbulence returns towards isotropy in the center of the channel around $z_1 = 5.5$ due to the favourable pressure gradient (leading to flow acceleration) in this region.
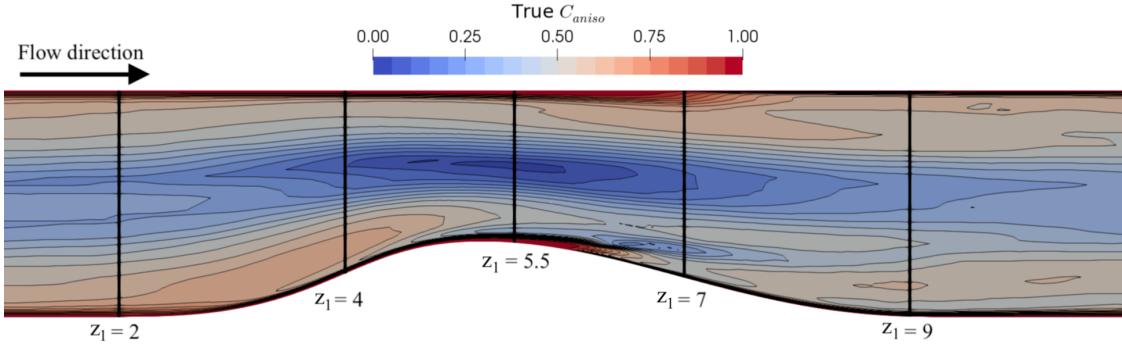
Figure 4: Contours of the *true* turbulent anisotropy constant $C_{aniso}$ for the convergent-divergent channel at $Re_\tau = 617$ (case 3b). The inflow is located at $z_1 = 0$, and $\mathbf{z}$ is non-dimensionlised by the channel half-height.

### 4.1.1. Initial Mondrian forest predictions

In Figure 5 the Mondrian forest (MF) predictions for $C_{aniso}$ and the *truth* (DNS) are plotted across the $z_1$-planes shown in Figure 4. When judging the Mondrian forest predictions here it is important to note that the Mondrian forest has not seen this flow case during training, so it is making predictions based on the physics it has learnt from the other flow cases. In Figure 5a, the Mondrian forest is trained only on the curved backward step (case 1, $Re_\tau = 618$) and periodic hills (case 2, $Re_\tau = 160$). Hence, a degree of extrapolation in feature space is likely here. Generally, the MF is predicting the physics seen in the DNS well, with high anisotropy predicted in the near-wall regions, and a further increase in anisotropy in the adverse pressure gradient, downstream of the bump (slices 3 and 4). Where predictions are less accurate, for example the return to isotropy region in the centre of the channel at slices 2-4, the Mondrian forest returns large predictive uncertainty values. Comparing the Random forest (RF) and MF predictions, they are generally about equal when measured against the Truth data. The RF is more accurate at the $z_1 = 4$ plane, but the MF predictions are superior at the $z_1 = 9$ plane. Overall, the MF appear to be slightly less noisy than the RF's predictions.

The linear eddy viscosity model (LEVM) in (2), and the quadratic terms which extend the LEVM to form a non-linear (NL) EVM, can be generalised as

$$
B_{ij} = \overbrace{-\frac{\nu_t}{k} S_{ij}}^{\text{LEVM}} + c_1 \left( S_{ik}S_{jk} - \frac{1}{3}S_{mk}S_{mk}\delta_{ij} \right)
$$
$$
\underbrace{+ c_2 \left( \Omega_{ik}S_{kj} + \Omega_{jk}S_{ki} \right) + c_3 \left( \Omega_{ik}\Omega_{jk} - \frac{1}{3}\Omega_{lk}\Omega_{lk}\delta_{ij} \right)}_{\text{Quadratic NLEVM}}.
$$

(23)

The values of $C_{aniso}$ obtained when the LEVM is used are shown in Figure 5, along with the values obtained when the quadratic NLEVM terms are included[3]. In most regions, the LEVM is seen to predict an insufficient level of turbulent anisotropy. Near the viscous walls, the NLEVM predicts high levels of turbulent anisotropy in closer agreement with the NDS, however it doesn't significantly improve predictions towards the centre of the channel. The clear superiority of the MF predictions here suggest that MFs offer the potential to augment LEVM and NLEVM based RANS models for such flows. As discussed by Hellsten and Wallin [43], many more elaborate RANS approaches exist, such as explicit algebraic Reynolds stress models (EARSM's) and Reynold stress transport models (RSTM's). However, Wu et al. [14] find a random forest based framework is able to offer substantial improvements over a RSTM model for duct and periodic hill flows, suggesting MF's could also offer improvements for more advanced RANS models.

---

[3]The constants $c_1, C_2, c_3$ are chosen to give Shih's quadratic non-linear eddy viscosity model [42].

11

(a) Only curved backward step and periodic hill cases (1 and 2) used for training
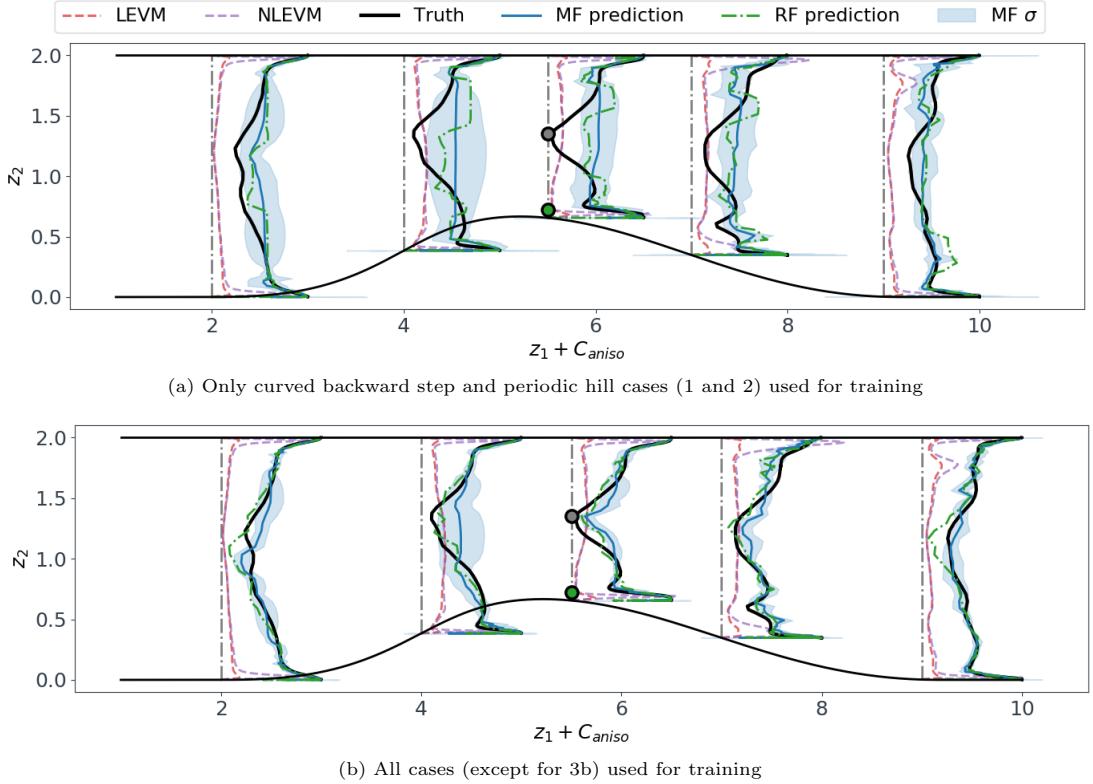


(b) All cases (except for 3b) used for training

Figure 5: Mondrian forest and random forest predicted turbulent anisotropy constant $C_{aniso}$ for the convergent-divergent channel at $Re_\tau = 617$ (case 3b), at the slices marked in Figure 4.

**Remark 2.** *The uncertainty bounds given by $\mu \pm \sigma$ do not entirely encompass the truth data here. However, this is expected, since we are plotting $\pm 1\sigma$, so would only expect $\approx 68\%$ of truth values to fall within these uncertainty bounds. Additionally, Lakshminarayanan et al. [26] examine the* probability calibration measures *for MF vs RF, and although MFs perform signficantly better than RFs with jackknifing, they still tend to be slightly over-confident with their uncertainty estimates—an artifact that could perhaps be fixed with calibration of their uncertainty estimates. That said, it should be noted these uncertainty estimates do not account for uncertainty due to the assumptions discussed in Remark 1.*

*4.1.2. More representative training data*

In Figure 5b all cases except the present test case (case 3b) are used for training. The effect of this additional training data is to improve the accuracy of the MF's mean predictions, whilst reducing its predictive uncertainty. This is especially noticeable in the return to isotropy region near the channel centre at slices 2-4. To understand why, the MF's predictions are interpreted using **SH**apley **A**dditive ex**P**lanation (SHAP) values, recently proposed by Lundberg and Lee [44]. This is an additive feature attribution method based on Shapley values from coalitional game theory. They are defined as Shapley values of a conditional expectation function of the original model [44]; effectively the $j$-th SHAP value $\phi_j$ represents the contribution to the prediction $\tilde{y}$, compared to the average prediction for the dataset, when conditioning on the $j$-th feature. Thus, the sum of the SHAP values satisfies

$$\sum_{j=1}^{D} \phi_j(\tilde{y}) = \tilde{y} - \mathbb{E}(\mathbf{y}). \tag{24}$$

To efficiently compute exact SHAP values for tree-based machine learning models, Lundberg et al. [45]

propose TreeSHAP, available from `github.com/slundberg/shap`[4]. This is used to compute SHAP values for predictions at the two locations marked by the circles in Figure 5. The seven SHAP values with greatest magnitude are plotted in Figure 6. Near the wall (Fig. 6a), the relatively small turbulent-to-strain timescale ratio ($t_t/t_s = 0.16$) decreases the predicted $C_{aniso}$ value compared to the average, whereas the axi-symmetric pressure gradient tensor value, $\widehat{\mathbf{A}}_p^2 = -0.092$, increases the predicted turbulent anisotropy. However, the sum of the SHAP values is relatively small here ($\sum \phi_i = 0.03$), meaning the predicted anisotropy is close to the average value ($\tilde{y} = 0.71$). In the center of the channel (Fig. 6b), the predicted anisotropy is considerably lower than the average mainly due to the high $Re_d$, but the other six features shown also contribute to the $C_{aniso}$ constant tending towards isotropy here.



$$f(\mathbf{x}) = \mathbb{E}(f(\mathbf{X})) + \sum \phi_i = 0.68 + 0.03 = 0.71$$

$$f(\mathbf{x}) = \mathbb{E}(f(\mathbf{X})) + \sum \phi_i = 0.68 - 0.58 = 0.10$$

(a) Near lower wall, $z_2 = 0.72$

(b) Mid-channel, $z_2 = 1.35$

Figure 6: The seven largest SHAP values at two different $z_2$ locations across $z_1 = 5.5$, for the convergent-divergent channel at $Re_\tau = 617$ (case 3b). Computed using the Mondrian forest trained on all cases (except for case 3b).

In Figure 7 kernel density estimates (KDE) show the distribution of training data in feature space, for the two training data sets used for the predictions in Figure 5. The six features identified as important for the two locations in question are shown, with $Re_d$ ignored since $Re_d = 2$ for all mid-channel locations. The near-wall test point ($z_2 = 0.72$, green circle) lies within both of the training data distributions. However, the mid-channel test point ($z_2 = 1.35$, grey circle) lies outside the first training data distribution (red contours), indicating that this training data isn't representative of the test point. The second training data set includes a training case (case 3c) with the same geometry as the test case, and the blue contours show that this training data better represents the mid-channel test point. The reduced extrapolation in feature space (for the six locally important features) explains why the predictive error and uncertainty are reduced in the return to isotropy region (fig. 5b), when more training data is added.

*4.2. Periodic hill*

The second case considered is the periodic hill (case 2), for which contours of $C_{aniso}$ from the DNS [35] are shown in Figure 8. The flow physics here is similar to the previous case, with high turbulence anisotropy observed in the shear layer where the flow separates at the hill crest ($z_1 < 1.5$). After the flow reattaches to the lower surface around $z_1 = 2$ the turbulence anisotropy decreases near the wall. When the windward side of the next hill is reached around $z_1 = 7.5$, the spanwise turbulent stress $\langle u_3' u_3' \rangle$ increases significantly compared to the $\langle u_1' u_1' \rangle$ and $\langle u_2' u_2' \rangle$ stresses near the wall, leading to high $C_{aniso}$ values here. This effect is caused by pressure–strain interactions in this region [35], and is interesting as many pressure-strain models assume that pressure–strain interactions isotropise the normal Reynolds stresses.

In Figure 9 the Mondrian forest predictions for $C_{aniso}$ and the *truth* are plotted across the $z_1$-planes shown in Figure 8. In Figure 9a, the Mondrian forest is trained only on the four duct flow cases (cases 4a/b/c/d). No pressure induced separation is present in the duct flows, hence these predictions represent a

---

[4]The code is slightly modified to be compatible with the Mondrian forests, and is made available at `github.com/ascillitoe/shap`.
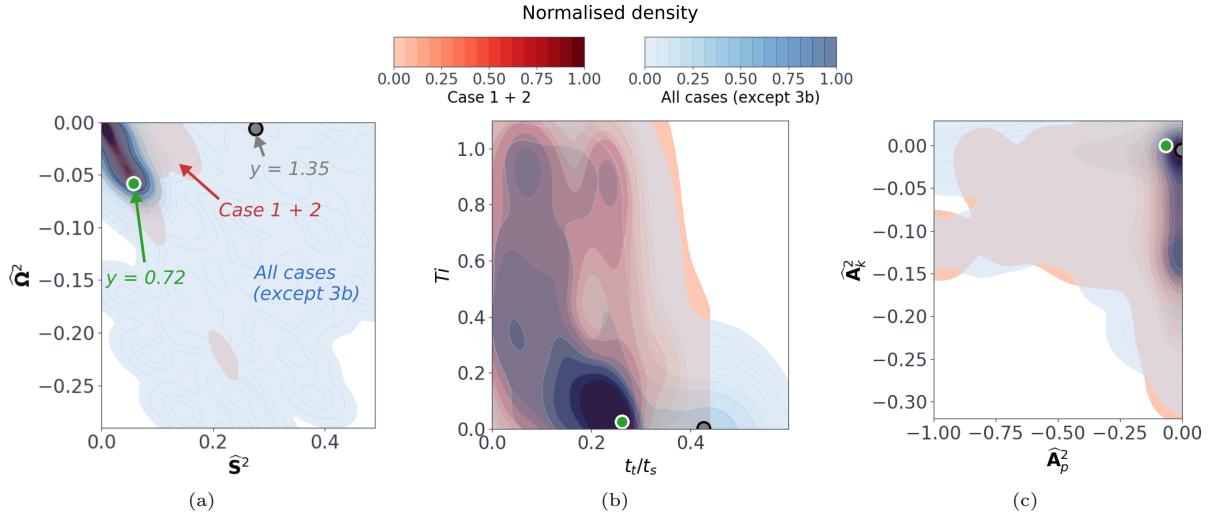
Figure 7: Bivariate kernel density estimation plots showing distributions of the two training data sets used for the convergent-divergent channel test case at $Re_\tau = 617$ (case 3b). The six most important features (after $Re_d$) are shown. The circular markers indicate the feature values at the two points plotted in Figure 6. Density is normalised by the maximum density.
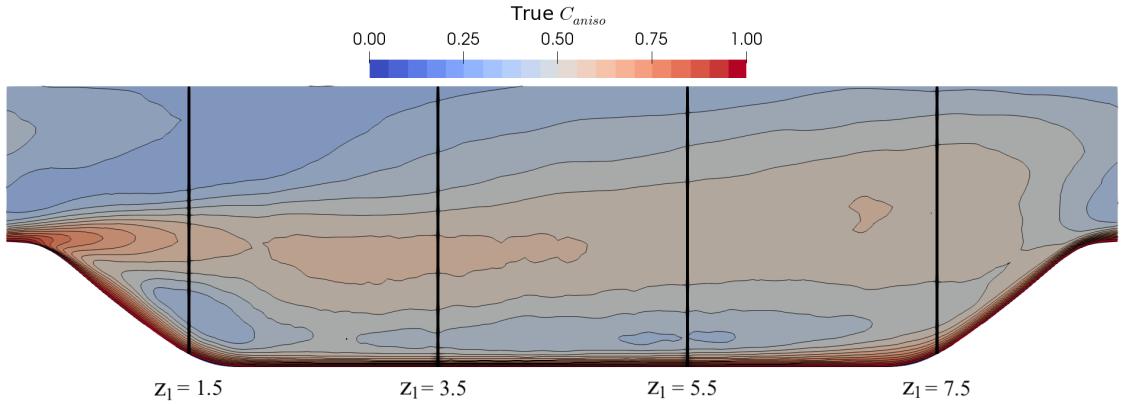


Figure 8: Contours of the *true* turbulent anisotropy constant $C_{aniso}$ for the periodic hill (case 2). The stream-wise periodic boundaries are located at $z_1 = 0$ and $z_1 = 9$, with $\boldsymbol{z}$ non-dimensionlised by the hill height.

considerable extrapolation. The anisotropy close to the lower wall is well captured, with small uncertainty bounds. However, away from the wall predictive errors are larger and the MF is less confident in its predictions. When the other flow cases (except for case 2) are added to the training data (Fig. 9b) the predictions are significantly more accurate, and the uncertainty bounds are much tighter. Again, the MF achieves significantly better accuracy compared to the LEVM and NLEVM predictions. Additionally, the MF and RF predictions are generally close here, with the MF predictions appearing less noisy.

### 4.3. Calibration of predictive uncertainty

In the preceding sections, the MF predictive mean $\mu$ and variance $\sigma^2$ are used to produce prediction intervals where, for example, the $1\sigma$ interval says we expect 68.27% of the true test values to lie within the interval $\mu \pm 1.0\sigma$. For these prediction intervals to be trusted, it is important to determine whether they are well calibrated. Following Lakshminarayanan et al. [26], *probability calibration curves* are plotted in Figure 10. For each $p\%$ (e.g. 90%), the prediction interval $\mu \pm q_p\sigma$ is calculated assuming Gaussian quantiles

$$q_p = \sqrt{2}\, \mathrm{erf}^{-1}(p). \tag{25}$$

14

(a) Only duct cases (4a/b/c/d) used for training
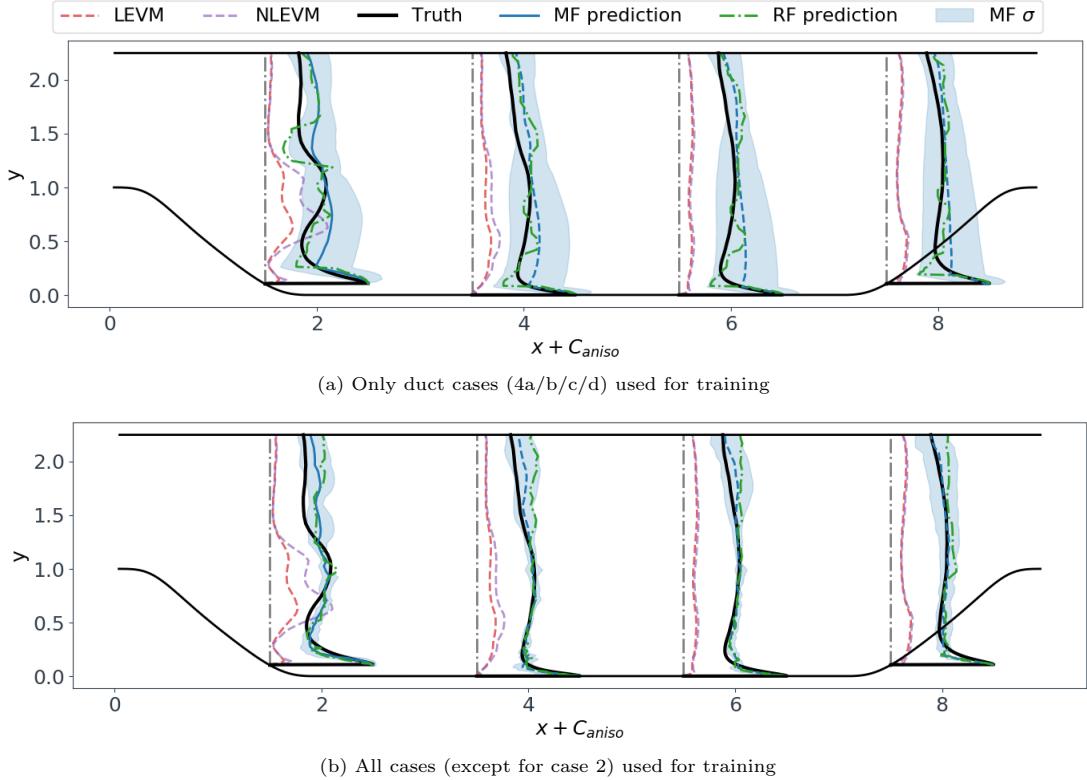


(b) All cases (except for case 2) used for training

Figure 9: Mondrian forest and random forest predicted turbulent anisotropy constant $C_{aniso}$ for the periodic hill (case 2), at the slices marked in Figure 8.

The percentage $\mathcal{Q}\%$ of true test points $\tilde{y}$ that lie within each prediction interval $\mu \pm q_p \sigma$ is then measured. If the model is perfectly calibrated, $\mathcal{Q}\%$ would equal $p\%$ for all $p$. The MF calibration curves for cases 2 and 3b are plotted in Figure 10, and RF jackknife uncertainty estimates are included for comparison. For both cases, the MF uncertainty estimates are well calibrated, with the MF calibration curves lying close to the ideal dashed line. This means that for a given $p\%$ interval, approximately $p\%$ of the test predictions are expected to lie within the interval. On the other hand, the RF jackknife uncertainty estimates are poorly calibrated, with the confidence intervals[5] displaying significant under-confidence. This suggests the random forest confidence intervals are not suitable to be used as prediction intervals. Additionally, as mentioned previously, preliminary (studies(Appendix B.3) indicated that jackknife re-sampling is prohibitively expensive for use in the framework explored here.

### 4.4. Comparison to distance based prediction confidence

A number of other approaches have been suggested to supplement data-driven turbulence modelling frameworks with measures of prediction confidence. Ling and Templeton [11] suggest using the Mahalanobis distance, while Wu et al. [18] suggest a measure based on KDE. Both are statistical *distance measures*, which are used to measure the distance between a test point and a distribution of training data. The Mahalanobis distance is defined as the distance between a point $\boldsymbol{x}$ in feature space and the mean of the training points $\boldsymbol{\mu}$,

$$D_m = \sqrt{(\boldsymbol{x} - \boldsymbol{\mu})^T \Sigma^{-} 1 (\boldsymbol{x} - \boldsymbol{\mu})}, \tag{26}$$

---

[5]The jackknife estimate returns confidence intervals not prediction intervals for the random forest (see Sec. Appendix A.2.1).

(a) Case 2: Periodic hill

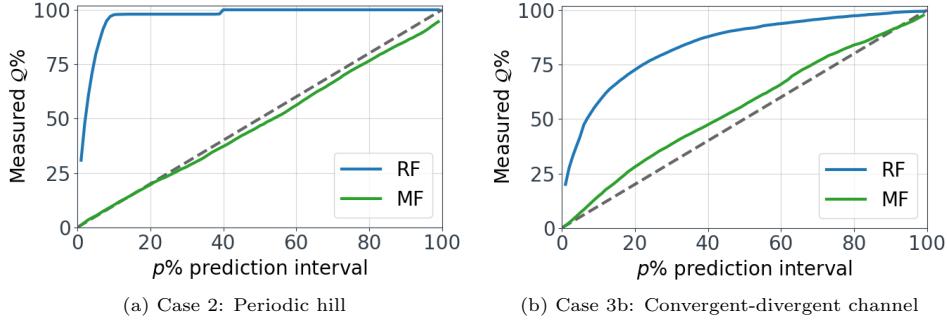(b) Case 3b: Convergent-divergent channel

Figure 10: Probability calibration curves for random forest and Mondrian forest predictions.

where $\Sigma$ is the covariance matrix of the training points. The normalised Mahalanobis distance for the point $\boldsymbol{x}$ is then defined as $\hat{D}_m = 1 - \psi_{D_m}$, where $\psi_{D_m}$ is the fraction of training points with a larger raw Mahalanobis distance than the point $\boldsymbol{x}$. It follows that a prediction at a point where $\hat{D}_m = 0$ involves no extrapolation, whereas $\hat{D}_m = 1$ involves very high extrapolation.

To assess how well $\hat{D}_m$ performs as a measure of prediction confidence, *violin plots* are used to inspect the correlation between $C_{aniso}$ prediction error and $\hat{D}_m$. The normalised Mahalanobis distance $\hat{D}_m$ is discretised into ten bins, and a kernel density plot[6] is generated for each bin. Figures 11a and 11b consist of violin plots[7] of a random forest's predictions for the convergent-divergent channel at $Re_\tau = 617$ (case 3b). Comparing Figures 11a and 11b it is apparent that when there is more extrapolation (Fig. 11b) a greater proportion of the test points have high predictive errors and large $\hat{D}_m$ values. Within the test data there is some correlation between error and $\hat{D}_m$. However, there are many data points where $\hat{D}_m$ is large while the error is small.

In Figures. 11c and 11d, the same comparison is shown for the Mondrian forest's predictive uncertainty. Correlation between the MF's predictive error and its predictive uncertainty $\sigma$ is generally much better here, compared to the correlation between error and $\hat{D}_m$. This implies the Mondrian forest's $\sigma$ is a more reliable measure of uncertainty compared to $\hat{D}_m$. However the MF is not infallible, and there are a small number of outliers where the MF is overconfident, with small $\sigma$ yet high error. The addition of an outlier filter like that proposed by Kaandorp and Dwight [15] might help remove the small number of outliers seen here.
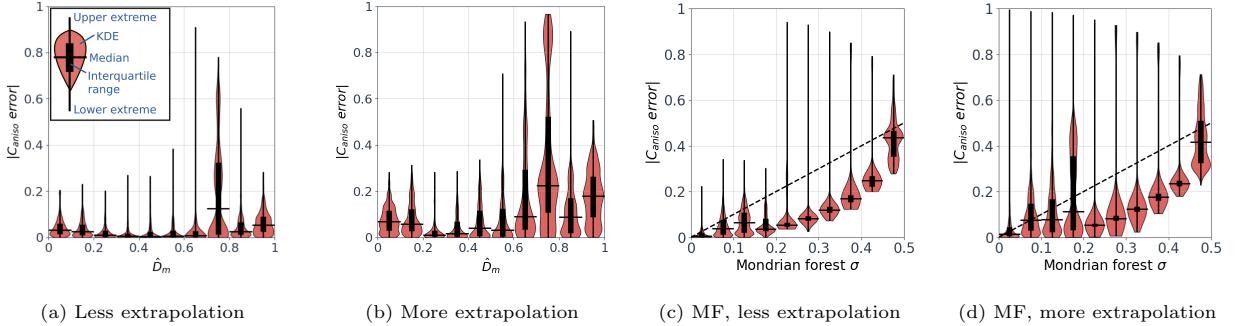


(a) Less extrapolation

(b) More extrapolation

(c) MF, less extrapolation

(d) MF, more extrapolation

Figure 11: Violin plot of error in random forest predicted turbulent anisotropy constant $C_{aniso}$ versus normalised Mahalanobis distance $\hat{D}_m$ ((a) and (b)), and Mondrian forest error versus predictive uncertainty ((c) and (d)), for the convergent-divergent channel at $Re_\tau = 617$ (case 3b). For (a) and (b), all the cases except for 3b are used as training data. For (c) and (d), the convergent-divergent channel cases at $Re_\tau = 395$ and 950 (3a and 3c) are also removed from the training data, implying a greater extent of extrapolation from the training data.

---

[6]Scott's rule is used to estimate kernel bandwidth.

[7]For the RF in this section, $L = 1280$ is chosen with no maximum depth to offer slightly more converged $\sigma$ estimates.

## 5. Data-driven turbulence modelling

The study of a Mondrian forest's predictions for the anisotropy constant suggests that Mondrian forests are capable of predicting a turbulent field variable with comparable predictive accuracy to random forests. The predictive uncertainty was also shown to be valuable in informing us whether we can trust the predictions. However, in practice, the end user is usually interested in other flow quantities such as velocity, or lift and drag coefficients. Therefore, a more useful application of the machine learning predictions is to augment an existing RANS solver. It may also be desirable to propagate the Mondrian forests' uncertainties through the RANS solver, to obtain uncertainties in the flow quantities of interest.

To explore the above, we train five Mondrian forests on the five variables in the anisotropy tensor discrepancy vector $\Delta \mathbf{B} = (\Delta \xi, \Delta \eta, h_1, h_2, h_3)$. The test case selected is the convergent-divergent channel at $Re_\tau = 395$ (case 3a), and training data for $\Delta \mathbf{B}$ is obtained by taking the differences between the RANS and NDS fields from the curved backward step (case 1) and periodic hills (case 2).



(a) $\xi$ coordinate

(b) $\eta$ coordinate

Figure 12: Mondrian forest predictions of the eigenvalue discrepancies $\Delta \xi$ and $\Delta \eta$, for the convergent-divergent channel at $Re_\tau = 395$ (case 3a). Also shown are $N_s = 50$ samples generated from MF uncertainty using the procedure described in Section 3.5.

The resulting Mondrian forest predictions for the eigenvalue discrepancies $\Delta \xi$ and $\Delta \eta$ are visualised in Figure 12. Generally, the MF predictions are in good agreement with the true discrepancies derived from comparing the NDS and RANS. However, in some regions the predictive uncertainty is high, especially for $\Delta \eta$. This could be improved by increasing the amount of representative training data, but our focus here is instead on propagating this uncertainty through the CFD solver. This is achieved by generating $N_s = 200$ sample fields for $\Delta \xi$ and $\Delta \eta$, using the procedure described in Section 3.5. The resulting samples are shown in Figure 12, where they are seen to be relatively smooth, whilst correlating closely with the MF uncertainty bounds also shown here.

In Figure 13 fifty of the samples are visualised on the Barycentric triangle (see [41]), for the points labelled (a), (b) and (c) in Figure 12. At point (c), the uncertainty in $\Delta \xi$ and $\Delta \eta$ is small and so the
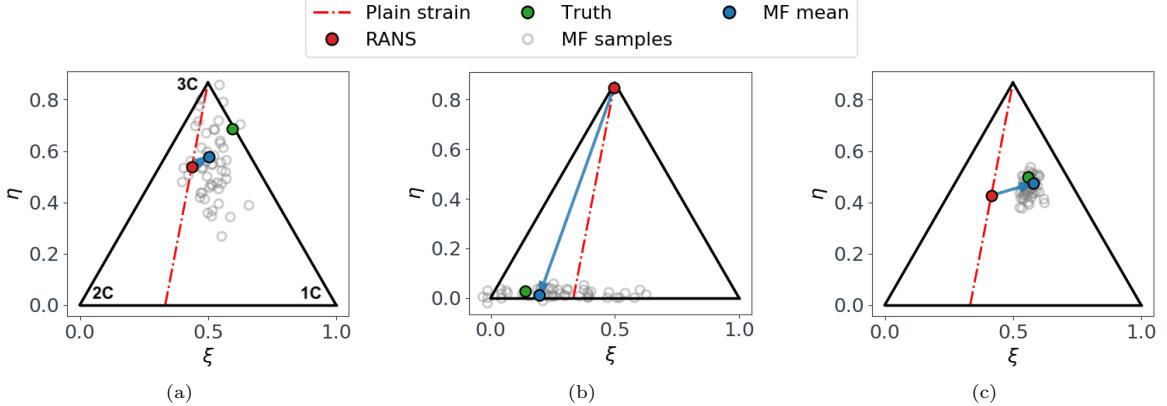
17

Figure 13: Turbulent anisotropy visualised on the Barycentric triangle [41] for the points labelled (a), (b) and (c) in Figure 12. The true anisotropy from NDS, the MF mean and LEVM predictions, and 50 MF samples are shown.

samples are clustered closely around the MF mean prediction. Whereas at (a) an (b) there is a high spread in the samples due to the high uncertainty here. When the uncertainty is high, there is a danger of samples being outside of the Barycentric triangle, leading to unrealisable Reynolds stresses. To prevent this, the modified SU2 solver (see Sec. 3.4) constrains the Cartesian coordinates $(\xi, \eta)$ to lie within the triangle:

$$
\begin{aligned}
\eta^* &= \max\left[0, \min\left(\frac{\sqrt{3}}{2}, \eta\right)\right] \\
\xi^* &= \max\left[\frac{1}{3}\eta, \min\left(1 - \frac{1}{3}\eta, \xi\right)\right]
\end{aligned}
\tag{27}
$$



Figure 14: Convergence histories for case 3a. Shown are the convergence history for the original RANS solution, and the convergence history when a randomly selected MF sample is propagated through the modified solver.

The samples can now be propagated through the modified SU2 solver, to obtain flowfield predictions for each sample. Samples are generated for $\Delta\xi$ and $\Delta\eta$ only, and for the eigenvector discrepancies $h_1, h_2, h_3$ we propagate the Mondrian forests' mean predictions. The modified solver is run for $N_s = 200$ samples, with the converged baseline RANS solution for case 2a used as the initial condition. The CFL number is reduced to 20, and the ramp parameters set to $\gamma_{max} = 0.9$ and $n_{max} = 100$. As seen in Figure 14, convergence is slower than for the original RANS simulation, but it is still acceptable.
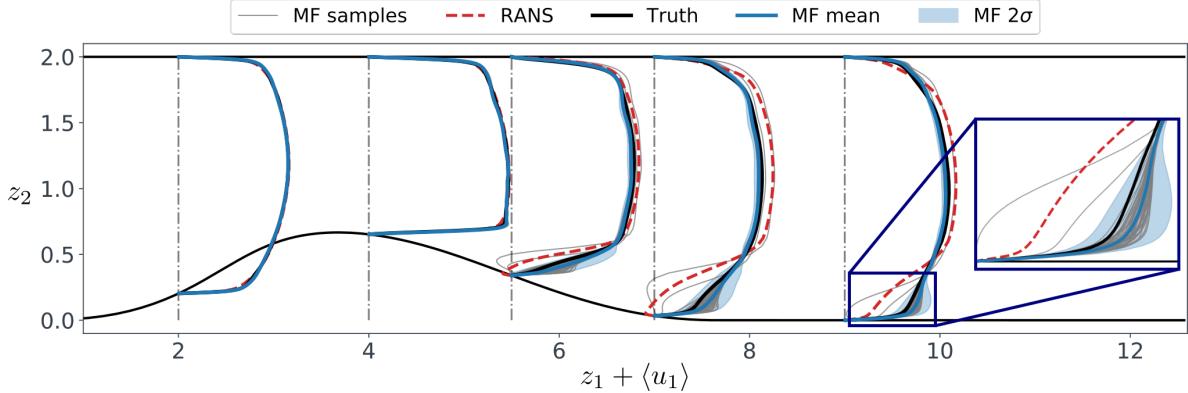
18

Figure 15: Profiles of axial velocity for the convergent-divergent channel at $Re_\tau = 395$ (case 3a), when the Mondrian forest correction for $\Delta\mathbf{B}$ is propagated through the CFD solver. The original RANS prediction, and $N_s = 50$ samples are also shown. The MF $2\sigma$ uncertainty bounds are obtained by assuming the $N_s = 200$ samples are normally distributed about the MF mean.

The resulting profiles of axial velocity are shown in Figure 15. There are slight differences between the MF mean and the *truth* from the NDS solution, which might be due to the small amount of the original RANS anisotropy tensor which was blended in for stability (recall $\gamma_{max} = 0.9$). This parameter could potentially be increased, at the expense of less favourable convergence behaviour. However, the predictions are generally in good agreement with the NDS, and are a significant improvement over the original RANS. Also shown are a selection of the samples obtained from the aforementioned sampling procedure. Despite the relatively high uncertainties in $\Delta\xi$ and $\Delta\eta$ observed in Figure 12, the uncertainty bounds obtained for the velocity are relatively small overall. This is likely to be because the large uncertainties in $\Delta\xi$ and $\Delta\eta$ are located away from the wall, where the Reynold stresses are relatively small. The uncertainty in velocity is larger in the separated flow region downstream of the divergent section. This suggests that the onset of flow separation, and the flow reattachment, are sensitive to the uncertainty in the ML predicted eigenvalues. This is in agreement with the studies of Gorlé et al. [46], who quantify the epistemic uncertainty of RANS models by perturbing the eigenvalues and eigenvectors of the anisotropy tensor.
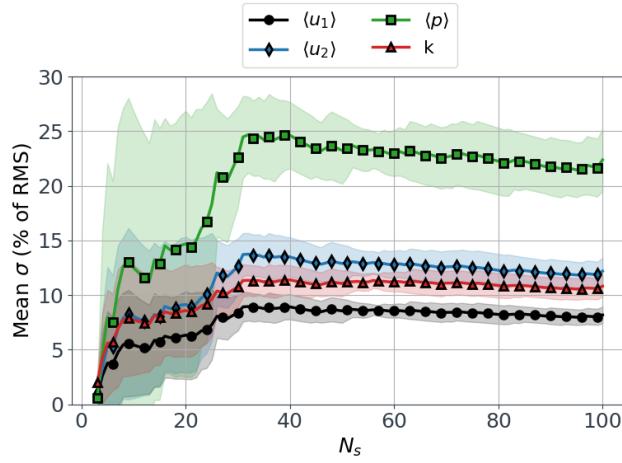


Figure 16: Spatial average of standard deviation in flow quantities across different numbers of samples. For each number of samples, ten different random selections of samples are taken from the full $N_s = 200$ set. The solid lines depict the mean of the spatial averages across the ten trials, whilst the shaded areas depict the standard deviation of the spatial averages.

A final important aspect of the proposed data-driven framework is its computational cost. Training

and prediction times for the Mondrian forests are of the order of seconds for the datasets used here (see Fig. B.17b), hence the cost of propagating the MF predictions through the CFD solver is the primary concern. With the propagation of each sample taking approximately as long to converge as the original RANS solution, there remains the question of how many samples must be propagated in order to achieve well converged uncertainty estimates for the output quantities of interest. To explore this, we monitor the spatially averaged standard deviation of $\langle u_1 \rangle$, $\langle u_2 \rangle$, $\langle p \rangle$ and $k$, whilst randomly selecting different numbers of samples from the 200 converged CFD solutions. As seen in Figure 16, satisfactory convergence is generally achieved with $N_s \geq 30$, although the uncertainty bounds for the pressure field are slightly more sensitive to the exact samples taken. If the uncertainties for the eigenvector rotations were also propagated, it is likely that more samples would be required due to the increased number of degrees of freedom.

To provide some points of comparison, in order to quantify the epistemic uncertainty of a RANS model, Gorlé et al. [46] and Xiao et al. [47] use $N_s = 5$ and $N_s = 100$ CFD evaluations respectively. Gorlé et al. [46] are able to only run $N_s = 5$ evaluations, because they consider the corners of the Barycentric triangle to be the limiting states of $(\eta, \xi)$, and therefore only perturb $(\eta, \xi)$ towards these three states. A similar approach cannot be taken in our case, since we are quantifying the uncertainty arising from augmenting the RANS model with ML predictions, and not the uncertainty of the RANS model itself. Figures 13a and 13b show that the MF's uncertainty can be strongly anisotropic, which implies uniformly perturbing towards the corners isn't appropriate. On a similar note, Figure 13 shows that the *true* turbulence anisotropy state is often not located at one of the corners. It follows that simply perturbing to a single corner state is unlikely to be a reliable way to improve RANS predictions. This provides further motivation for data-driven approaches such as the one proposed in this paper.

## 6. Conclusions

For the emerging field of data-driven turbulence modelling to enter mainstream use, methods to quantify the uncertainties arising from such approaches will be necessary. To this end, the present work explores the use of a recently proposed machine learning algorithm, Mondrian regression forests, for data-driven turbulence modelling.

In Section 4 Mondrian forests were found to offer comparable accuracy to random forests when predicting the turbulence anisotropy constant derived from near direct simulations (NDS). Interestingly, feature selection indicated that irrelevant features did not degrade the Mondrian forests' accuracy in this case. However, this may be specific to the flow cases and feature set used in this paper, and it would be prudent to perform a similar feature selection procedure on any new training data. The Mondrian forests' uncertainty estimates appeared to provide a good measure of prediction confidence, with high uncertainty in regions of the flow where predictions were far from the truth data. Adding more representative training data generally reduced both predictive errors and uncertainty. Comparing to an a priori distance measure suggested by Ling and Templeton [11], the Mondrian forest uncertainty estimates showed better correlation with predictive errors with fewer regions of significant under or over-confidence. Additionally, probability calibration curves suggest the Mondrian forest uncertainty estimates are well calibrated, which is important in order for the prediction intervals to be trusted. Furthermore, hyperparameter tuning demonstrated that the Mondrian forests achieve converged uncertainty estimates with a considerably computational cost compared to random forests with jackknife estimates.

To explore the suitability of Mondrian forests in a data-driven turbulence modelling framework, in Section 5 they were trained to predict the turbulent anisotropy tensor discrepancies. Using a modified RANS solver similar to that described by Kaandorp and Dwight [15], the Mondrian forest predicted mean discrepancy vector was successfully propagated forward to obtain converged flowfields. Depending on the use case propagating only the mean prediction might be sufficient, and the Mondran forests' uncertainties can still be used to decide whether it is worth proceeding with this, or whether more suitable training data is required. However, in other cases it might be desirable to propagate the uncertainties themselves. As demonstrated in Section 5, this allows for the quantification of uncertainty in quantities of interest such as velocity and pressure due to the uncertainty in the Mondrian forest's predictions.

20

The above quantification of uncertainties is an important step if such approaches are to enter mainstream use, but future work must also address other sources of uncertainty in the ML predictions. For example, the point-wise locality assumption made in the present framework, and any differences between the RANS and NDS mean velocity and pressure fields. The first error source could be mitigated by introducing additional flow features to account for non-locality. The second error source could perhaps be reduced by modifying the present framework so that input-output training data is derived purely from NDS data. Improved accuracy could also be achieved by improving the Mondrian forest implementation itself.

Nevertheless, based on the aforementioned findings, for data-driven turbulence modelling Mondrian forests appear to offer a promising alternative to random forests, and other probabilistic methods such as BART trees or Bayesian neural networks. A final potential area of future work is online learning. Mondrian forests and trees can be updated in an online fashion as new data becomes available, and the resulting models should be identical to their batch trained counterparts. This capability can save time, as the model must not be completely retrained when new data is obtained. Additionally, the Mondrian trees are only altered in the region of feature space populated by the new online data. As an example, this would allow for new higher Reynolds number training data to be added, without affecting the model's predictions at a lower Reynolds number. Such capability might be important in an industrial setting where repeatability is crucial.

### Acknowledgement

### Appendix A. Mondrian forests and random forests

*Appendix A.1. Decision trees*

Following the notation of [26], we describe a decision tree by the tuple $(\mathsf{T}, \boldsymbol{\delta}, \boldsymbol{\xi})$, where $\mathsf{T}$ is the tree, $\boldsymbol{\delta}_j \in \{1, \ldots, d\}$ is the *split dimension* and $\boldsymbol{\xi}_j \in \mathbb{R}$ is the *split location* for the $j^{th}$ node in the tree. A decision tree trained on the training data $(\boldsymbol{X}_N, \mathbf{y}_N)$ is a hierarchical partitioning of the input data. At each node in the tree, the data is split in a binary fashion

$$
\begin{aligned}
B_{left(j)} &:= \left\{ \mathbf{x} \in B_j : x_{\delta_j} \le \zeta_j \right\} \\
B_{right(j)} &:= \left\{ \mathbf{x} \in B_j : x_{\delta_j} > \zeta_j \right\}
\end{aligned}
\tag{A.1}
$$

where the $j^{th}$ block of data is

$$
B_j = (l_{j1}, u_{j1}] \times \cdots \times (l_{jd}, u_{jd}],
\tag{A.2}
$$

with $l_{jd}$ and $u_{jd}$ the lower and upper bounds of the rectangular block $B_j$ along dimension $d$. As an example, for the decision tree fitted to the input data $\boldsymbol{x} \in [0,1]^2$ in Figure 2a, the root node $\varrho$ splits with $\delta_1 = 2$ and $\xi_1 = 0.6$, leading to its right child node $j = right(\varrho)$ having the data block $B_j = (0,1] \times (0.6, 1]$.

There are many induction algorithms available to learn a decision tree structure from training data, such as the popular CART algorithm [31] for classification and regression. Generally, these algorithms learn the tree structure $\mathsf{T}$ and leaf node parameters $\boldsymbol{\theta}$ by greedily optimising an appropriate criterion, such as mean squared error (MSE) for regression. The parameter $\boldsymbol{\theta}_j$ parametrises the conditional distribution $p(y|\mathbf{x} \in B_j)$, and for regression is simply the mean of the $K_j$ number of responses $y_K$ residing in the leaf node's block $B_j$:

$$
\boldsymbol{\theta}_j = \frac{1}{|K_j|} \sum_{q \in N_j} y_q
\tag{A.3}
$$

For a given test data point $\tilde{\mathbf{x}}$, a prediction is made by walking through the decision tree to identify its corresponding leaf node $j$ and then returning the parameter $\boldsymbol{\theta}_j$.

*Appendix A.2. Random forests*

Let $\boldsymbol{\varphi}$ be a random forest consisting of $L$ decision trees $\mathsf{T}_1, \ldots, \mathsf{T}_L$. For the data point $\tilde{\mathbf{x}}$, the random forest's prediction is an average of each tree's prediction

$$\boldsymbol{\varphi}(\tilde{\mathbf{x}}) = \frac{1}{L} \sum_{i=1}^{L} w_i \mathsf{T}_i(\tilde{\mathbf{x}}) \tag{A.4}$$

where $w_i$ is the $i$-th weighting term[8]. The individual trees in a random forest are randomised with bootstrap aggregation (*bagging*), where each tree is trained on a slightly different subset of the training data, and additionally the set of candidate splits within each node are randomly sub-sampled. This randomisation may slightly increase bias, but it significantly decreases over-fitting (prediction variance). The random forest regressor algorithm from `github.com/scikit-learn/scikit-learn` is used in this paper. Unless otherwise stated, all hyper-parameters are kept at their defaults (as of version 0.22.1).

*Appendix A.2.1. Jackknife variance estimates*

Wager et al. [17] propose the use of jackknife re-sampling to provide confidence intervals for random forest predictions. The *infinitesimal jackknife*, provides a variance measure, i.e., $V_{IJ} = \mathrm{var}(\boldsymbol{\varphi}(\tilde{\mathbf{x}}))$ for a random forest's predicted response to the input $\tilde{\mathbf{x}}$. For a random forest $\boldsymbol{\varphi}$ trained on the data $(\boldsymbol{X}_N, \mathbf{y}_N)$, the infinitesimal jackknife variance is given by

$$V_{IJ} = \sum_{i=1}^{N} \left[ \frac{1}{L} \sum_{j=1}^{L} (|\kappa_i|_j - 1) \left( \mathsf{T}_j(\tilde{\mathbf{x}}) - \overline{\mathsf{T}}(\tilde{\mathbf{x}}) \right) \right] \tag{A.5}$$

where $|\kappa_i|_j$ denotes the number of times $\kappa_i$ appears in the $j^{th}$ bootstrap sample, $\mathsf{T}_j(\tilde{\mathbf{x}})$ is the predicted response of the $j^{th}$ tree, and $\overline{\mathsf{T}}(\tilde{\mathbf{x}})$ is the mean of $\mathsf{T}_j(\tilde{\mathbf{x}})$ over $j = 1, \ldots, L$. The variance in (A.5) can be biased upwards when $L$ is small, therefore a bias corrected version is suggested

$$V_{IJ-U} = V_{IJ} - \frac{N}{L^2} \sum_{j=1}^{L} \left( \mathsf{T}_j(\tilde{\mathbf{x}}) - \overline{\mathsf{T}}(\tilde{\mathbf{x}}) \right)^2. \tag{A.6}$$

A *Jackknife-after-bootstrap* variance estimate is also defined, and Wager et al. [17] suggest that the arithmetic mean of this and $V_{IJ-U}$ provide more unbiased variance estimates. However, since the infinitesimal jackknife $V_{IJ-U}$ tends to overestimate the variance, this is used alone here as a conservative estimate. The standard deviation $\sigma = \sqrt{V_{IJ-U}}$ provides *confidence intervals*. They measure how far the random forest prediction $\boldsymbol{\varphi}(\tilde{\mathbf{x}})$ – obtained by building $M$ trees on a sample – is from its *expected value*, which is the average random forest prediction obtained by building $M$ trees across different samples. They do not provide *prediction intervals*, meaning they do not tell us how far the predicted value $\boldsymbol{\varphi}(\tilde{\mathbf{x}})$ is from the true value $\tilde{y}$.

*Appendix A.3. Mondrian forests*

*Appendix A.3.1. Mondrian trees*

Mondrian trees are restrictions of Mondrian processes to a finite set of points. Mondrian processes are families $\mathcal{M}_t : t \in [0, \infty)$ of random hierarchical binary partitions of $\mathbb{R}^d$, where $\mathcal{M}_t$ is a refinement of $\mathcal{M}_s$ wherever $t > s$.[9] A Mondrian tree $\mathcal{T}$ is a tuple $(\mathsf{T}, \boldsymbol{\delta}, \boldsymbol{\xi}, \boldsymbol{\zeta})$, where $(\mathsf{T}, \boldsymbol{\delta}, \boldsymbol{\xi})$ is a decision tree. The additional parameter $\boldsymbol{\zeta} = \{\zeta\}_{j \in \mathsf{T}}$ specifies the split time $\zeta$ for each node $j$. Split times increase with depth of the tree, and are involved in the setting of a hierarchical prior (amongst other things).

---

[8] All $w_i$'s are usually set to unity for random forests.

[9] $t$ is referred to as a *time*, but this should not be confused with a physical time related to the data, or to discrete time in an online learning setting.

*Appendix A.3.2. Hierarchical prior and predictive posterior*

Mondrian trees are probabilistic models, which determine $p_\mathcal{T}(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y}))$. Lakshminarayanan et al. [26] assume the responses in each leaf node are Gaussian distributed, so that every node $j \in \mathsf{T}$ has a mean parameter $\mu_j$. A hierarchical Gaussian prior is then used for $\boldsymbol{\mu} = \{\mu_j : j \in \mathsf{T}\}$, such that

$$\mu_\varrho|\mu_H \sim \mathcal{N}\left(\mu_H, \phi_\varrho\right), \quad \mu_j|\mu_{\text{parent}(j)} \sim \mathcal{N}\left(\mu_{\text{parent}(j)}, \phi_j\right), \tag{A.7}$$

where $\varrho$ denotes the root node, and $\phi_j = \gamma_1\sigma(\gamma_2\zeta) - \gamma_1\sigma(\gamma_2\zeta parent(j))$. The sigmoid function $\sigma(t) = (1 + e^{-t})^{-1}$ encodes the prior assumption that children are expected to be more similar to their parent nodes as tree depth increases. The hyperparameters $\mu_H, \gamma_1, \gamma_2$ are set according to Appendix B in Ref. [26].

As discussed in Appendix A.1, for a typical decision tree the predicted response $\tilde{y}$ for a test point $\tilde{\mathbf{x}}$ is simply the average of the responses in $B^x_{\text{leaf}(\tilde{\mathbf{x}})}$. With Mondrian trees, the test point can *branch off* the existing tree at any point along the path from the root node to leaf($\tilde{\mathbf{x}}$). Hence, the predictive posterior over $y$ is a weighted mixture of Gaussians along the path from the root node to leaf($\tilde{\mathbf{x}}$)

$$p_\mathcal{T}\left(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y})\right) = \sum_{j \in \text{path}(\text{leaf}(\tilde{\mathbf{x}}))} w_j \mathcal{N}\left(m_j, v_j\right), \tag{A.8}$$

where the weight $w_j$ describes the probability of branching off just before reaching the $j^{th}$ node, and $m_j$ and $v_j$ are the predictive mean and variance at the $j^{th}$ node. As the test point $\tilde{\mathbf{x}}$ moves further away from the training data at a given node $j$, i.e. as $\tilde{\mathbf{x}}$ moves away from the block $B^x_j$ in Figure 2b, the probability $w_j$ increases. This causes Mondrian forests to exhibit higher uncertainty as $\tilde{\mathbf{x}}$ moves further away from the training data $(\boldsymbol{X}, \mathbf{y})$. Additionally, the Mondrian forest's predicted response $\tilde{y}$ approaches the prior as $\tilde{\mathbf{x}}$ moves further away from the training data. To obtain the predictive mean and variance at each node, we require $p_\mathcal{T}(\boldsymbol{\mu}|(\boldsymbol{X}, \mathbf{y}))$, the posterior over $\boldsymbol{\mu}$. This is computed using Gaussian belief propagation. Lakshminarayanan et al. [26] note that, since a hierarchical tree structure is used, posterior inference can be performed with a computational cost of $O(N)$ for $N$ training samples. This is compared to a Gaussian process whose computational cost is typically $O(N^3)$.

*Appendix A.3.3. Forests of Mondrian trees*

Lakshminarayanan et al. [25] propose combining Mondrian trees to form a Mondrian forest as a way to reduce over-fitting behaviour. The prediction of a Mondrian forest is then the average prediction from the $L$ number of Mondrian trees

$$p(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y})) = \frac{1}{L}\sum_j^L p_j(\tilde{y}|\tilde{\mathbf{x}}, (\boldsymbol{X}, \mathbf{y})). \tag{A.9}$$

Again, this predictive posterior over $\tilde{y}$ is a mixture of Gaussians, and it is straightforward to calculate the predictive mean and variance from this. For further details see Section 3 in [26].

Following the approach of other practitioners, such as Mourtada et al. [48], bootstrap aggregation is not used when training Mondrian forests in this paper. Due to the randomisation involved in the construction of Mondrian trees, each tree is still likely to be different. The Mondrian forest regressor algorithm from github.com/scikit-garden/scikit-garden is used in this paper. Lakshminarayanan et al. [26] control the depth of their Mondrian trees by stopping the splitting of nodes which have less than a given number of data points. To provide a fairer comparison with the random forest algorithm, in this paper we set this hyperparameter to two and additionally limit the depth of Mondrian trees with the $\mathcal{D}_{max}$ hyperparameter.

# Appendix B. Feature selection and hyperparameter tuning

This section documents preliminary studies done to select training features and suitable hyperparameters. In the studies here, random forests and Mondrian forests are trained on all of the flow cases listed in Table 1, with the turbulent anisotropy constant $C_{aniso}$ defined in (22) the response variable to be predicted.

*Appendix B.1. Measuring predictive accuracy*

To quantify predictive accuracy two metrics are used. The coefficient of determination, or $R^2$ score, given by

$$R^2 = 1 - \frac{\sum_{n=1}^{N}(y_n - \varphi(\boldsymbol{x}_n))^2}{\sum_{n=1}^{N}(y_n - \bar{y})^2} \tag{B.1}$$

measures how well the true data $y$ is replicated by the model's predictions, $\varphi(\boldsymbol{x}_n)$. The $R^2$ score spuriously increases with the dimension $d$ (recall $\boldsymbol{x} \in \mathbb{R}^D$), therefore score is adjusted to correct for this:

$$\bar{R}^2 = 1 - (1 - R^2)\frac{N - 1}{N - D - 1}. \tag{B.2}$$

The mean absolute error (MAE) provides a measure of error in the model's predictions

$$MAE = \frac{1}{N} \sum_{n=1}^{N} |y_n - \varphi(\boldsymbol{x}_n)|, \tag{B.3}$$

and is given as a percentage of the true mean response $\bar{y}$ in this paper. MAE is chosen over root-mean-square error (RMSE) due to RMSE's sensitivity to sample size and exaggeration of small numbers of large errors.

When MAE is calculated based on the training response $\mathbf{y}$ and training data predictions $\varphi(\boldsymbol{x})$ it is referred to as the *training error*, meanwhile when it is based on the test data $y^*$ and $\varphi(\boldsymbol{x}^*)$ it is referred to as the *test error*. The out-of-sample error, or test error, is important as it assesses how well the machine learning model generalises to data it has not seen during training. As is standard practice amongst machine learning practitioners, cross validation is used to provide fair error measures. The commonly used $k$-fold cross validation involves randomly partitioning $\boldsymbol{D}_{1:N}$ into $k$ folds. Training is then performed on $k - 1$ folds, with the left-out fold used as test data. This is performed $k$ times so that all folds are used as test data, and the averaged errors over the $k$ folds are taken. Roberts et al. [49] notes that performing such a strategy on spatial data results in serious underestimation of predictive errors due to the dependence between neighbouring observations. To avoid this problem leave-one-group-out (LOGO) cross validation is used in this paper, where each flow case in Table 1 is used as a fold.

*Appendix B.2. Feature selection*

Generally, increasing the dimension $d$ of the input feature space increases the computational cost of training and predictions, and as noted by Lakshminarayanan et al. [26], irrelevant features can be detrimental to Mondrian forest predictions. Therefore, backward elimination is used to remove less important features. At each iteration the least important feature $x_d$ in $\boldsymbol{x} = (x_1, \ldots, x_D)^T$ is removed, and the forest is retrained on the reduced data. Feature importance is measured using permutation importance, proposed by Breiman [30]. The permutation importance $PI(x_d)$ quantifies the change in predictive error due to permuting the feature $x_d$ in the input data.

Backward elimination is performed on Random forests and Mondrian forests, with features iteratively removed one at a time until $d = 2$. For both forests the maximum depth is set at $\mathcal{D}_{max} = 40$ and the number of trees at $M = 160$. The mean $\bar{R}^2$ scores from LOGO cross-validation across all ten flow cases are plotted against $d$ in Figure B.17a. The key points from this plot are as follows:

- The random forest (RF) generally has a $1 - 2\%$ higher $\bar{R}^2$ score than the Mondrian forest (MF), but the MF still returns competitive accuracy.

- Many features can be removed without affecting the predictive accuracy of the MF or RF. This is perhaps unsurprising, since many of the 47 invariants are zero for two-dimensional flows. For $d < 12$ the $\bar{R}^2$ scores begin to decrease as features containing important information are removed.

- The problem of irrelevant features degrading MF accuracy is not observed here. This may be because most of the non-zero features are relevant in the flow cases considered.
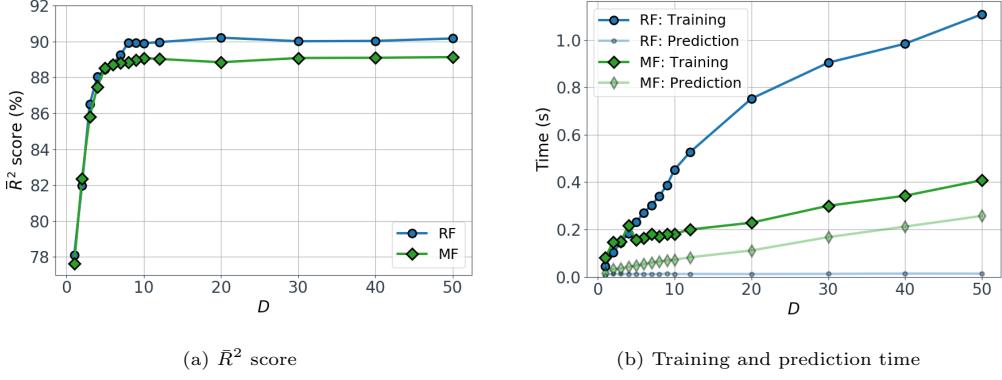
24

(a) $\bar{R}^2$ score



(b) Training and prediction time

Figure B.17: Effect of number of features. Features are removed using backward elimination based on permutation importance. $\bar{R}^2$ scores are the mean scores from LOGO cross-validation. Results are averaged over three runs. Training and prediction times are measured on a personal laptop with a $8^{th}$ generation Intel® Core™ i5 processor.

Since the feature space can be reduced to $d = 12$ with no significant loss of accuracy, this is done for the remainder of this paper. When backward elimination is halted at $d = 12$ with the Mondrian forest, the following feature set is obtained:

$$\boldsymbol{x} = \left\{ \widehat{\boldsymbol{S}}^2, \widehat{\boldsymbol{\Omega}}^2, \widehat{\mathbf{A}}_p^2, \widehat{\mathbf{A}}_k^2, \widehat{\boldsymbol{\Omega}}\widehat{\mathbf{A}}_k, \widehat{\boldsymbol{\Omega}}\widehat{\mathbf{A}}_k \widehat{\boldsymbol{S}}^2, \widehat{\boldsymbol{\Omega}}^2 \widehat{\mathbf{A}}_k \widehat{\boldsymbol{S}}, \widehat{\mathbf{A}}_k^2 \widehat{\mathbf{A}}_p \widehat{\boldsymbol{S}}, \widehat{\mathbf{A}}_k^2 \widehat{\boldsymbol{S}} \widehat{\mathbf{A}}_p \widehat{\boldsymbol{S}}^2, Re_d, \hat{k}, \hat{\omega} \right\}^T, \tag{B.4}$$

and it is used for all subsequent RF and MF computations in this paper. As shown in Figure B.17b, the reduced feature space allows for reduced training/prediction times, with no appreciable loss in predictive accuracy.

### Appendix B.3. Hyper-parameter tuning

To determine suitable settings for the number of trees $M$ and maximum depth of trees $D_{max}$ an exhaustive grid search is performed. LOGO cross validation is performed for every combination of the hyperparameters in the sets $M = \{2, 4, 6, 8, 10, 20, 40, 80, 100\}$ and $\mathcal{D}_{max} = \{1, 5, 10, 20, 40, 60, 80, 100\}$, and the resulting $\bar{R}^2$ scores are shown in Figures B.18a and B.18b. Both RF's and MF's display similar trends here. The high training and test errors for low values of $\mathcal{D}_{max}$ in Figure B.18a indicate that the shallow trees suffer from under-fitting to the $C_{aniso}$ relationships in the data. Whereas the moderate difference between training and test errors at higher $\mathcal{D}_{max}$ values is possibly due to over-fitting and a lack of generalisation. Increasing $M$ (Fig. B.18b) yields only a slight reduction in MAE, suggesting that if interpretability of the ML model is important single trees ($M = 1$) could be used with only a slight increase in predictive errors.

A second grid search is performed over a larger range of $M$ and $\mathcal{D}_{max}$ to determine the sensitivity of the prediction uncertainty $\sigma(C_{aniso})$ to the hyperparameters. For the RF's $\sigma(C_{aniso})$ is estimated using the jackknife (Sec. Appendix A.2.1), while for the MF's $\sigma(C_{aniso})$ is naturally returned when predicting $C_{aniso}$ (see Sec. Appendix A.3.3). Figures B.18c and B.18d show that the MF's achieve converged $\sigma$ estimates with a relatively low number of trees ($M > 160$). On the other hand, convergence isn't achieved with the RF's even with $M > 1000$ trees. Increasing $M$ further (i.e. towards $N$) is impractical due to the significant computational cost of computing the infinitesimal jackknife for large values of $M$.

Based on the aforementioned hyperparameter tuning, $M = 160$ and $\mathcal{D}_{max} = 40$ is chosen for all subsequent MF's to give low MAE's and converged $\sigma$ estimates.
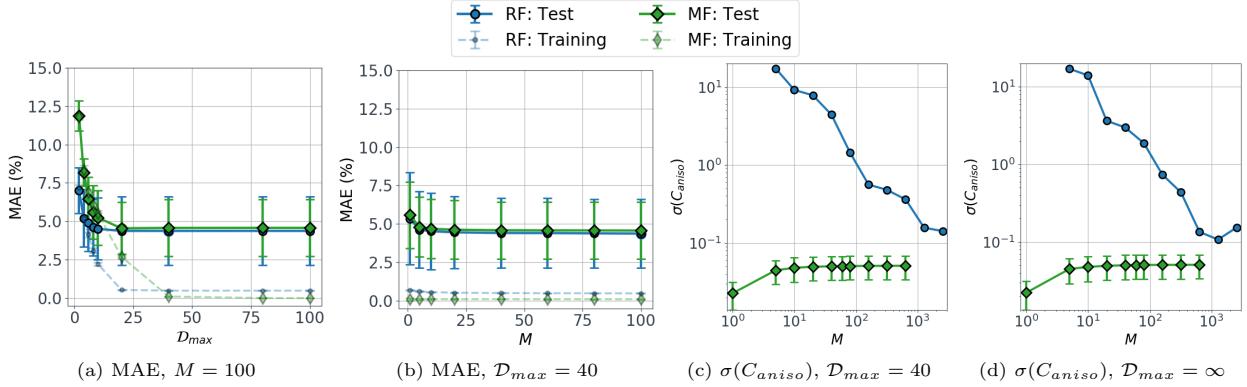
Figure B.18: Effect of hyper-parameters on MAE and predicted $\sigma(C_{aniso})$. Due to the computational cost of jackknifing, the RF results in c) and d) are for case 3b only. All other results are averaged across all ten folds of the LOGO cross-validation. Error bars represent standard deviation of MAE and predicted $\sigma(C_{aniso})$ across individual folds in cross-validation.

# References

## References

[1] J. Wild, High-Performance High-Lift Design for Laminar Wings, in: J.-P. S. D. Knörzer, C. Warsop, C. Diaconescu (Ed.), Proc. Seventh Eur. Aeronaut. Days, Aviation in Europe - Innovation for Growth, London, UK, 2015, pp. 305–310. doi:10.2777/62810.

[2] P. G. Tucker, J. R. DeBonis, Aerodynamics, computers and the environment, Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences 372 (2014) 20130331. doi:10.1098/rsta.2013.0331.

[3] A. J. Brand, J. Peinke, J. Mann, Turbulence and wind turbines, Journal of Physics: Conference Series 318 (2011) 072005. doi:10.1088/1742-6596/318/7/072005.

[4] L. Raynal, F. Augier, F. Bazer-Bachi, Y. Haroun, C. Pereira Da Fonte, CFD Applied to Process Development in the Oil and Gas Industry - A Review, 2016. doi:10.2516/ogst/2015019.

[5] A. D. Scillitoe, P. G. Tucker, P. Adami, Large Eddy Simulation of Boundary Layer Transition Mechanisms in a Gas-Turbine Compressor Cascade, J. Turbomach. 141 (2019) 1–10. doi:10.1115/1.4042023.

[6] D. Mehta, A. H. van Zuijlen, B. Koren, J. G. Holierhoek, H. Bijl, Large Eddy Simulation of wind farm aerodynamics: A review, J. Wind Eng. Ind. Aerodyn. 133 (2014) 1–17. doi:10.1016/j.jweia.2014.07.002.

[7] B. Blocken, LES over RANS in building simulation for outdoor and indoor applications: A foregone conclusion?, Build. Simul. 11 (2018) 821–870. doi:10.1007/s12273-018-0459-3.

[8] J. C. Hunt, A. M. Savill, Guidelines and criteria for the use of turbulence models in complex flows, 2005. doi:10.1017/CBO9780511543227.008.

[9] P. G. Tucker, Trends in turbomachinery turbulence treatments, Prog. Aerosp. Sci. 63 (2013) 1–32. doi:10.1016/j.paerosci.2013.06.001.

[10] K. Duraisamy, G. Iaccarino, H. Xiao, Turbulence Modeling in the Age of Data, Annu. Rev. Fluid Mech. 51 (2019) 357–377. doi:10.1146/annurev-fluid-010518-040547.

[11] J. Ling, J. Templeton, Evaluation of machine learning algorithms for prediction of regions of high Reynolds averaged Navier Stokes uncertainty, Phys. Fluids 27 (2015) 85103. doi:10.1063/1.4927765.

[12] J. Ling, A. Kurzawski, J. Templeton, Reynolds averaged turbulence modelling using deep neural networks with embedded invariance, J. Fluid Mech 807 (2016) 155–166. doi:10.1017/jfm.2016.615.

[13] A. P. Singh, R. Matai, A. Mishra, K. Duraisamy, P. A. Durbin, Data-driven augmentation of turbulence models for adverse pressure gradient flows, in: 23rd AIAA Computational Fluid Dynamics Conference, American Institute of Aeronautics and Astronautics, 2017, pp. 1–20. doi:10.2514/6.2017-3626.

[14] J. L. Wu, H. Xiao, E. Paterson, Physics-informed machine learning approach for augmenting turbulence models: A comprehensive framework, Phys. Rev. Fluids 7 (2018) 1–42. doi:10.1103/PhysRevFluids.3.074602.

[15] M. L. Kaandorp, R. P. Dwight, Data-driven modelling of the reynolds stress tensor using random forests with invariance, Computers & Fluids 202 (2020) 104497. doi:10.1016/j.compfluid.2020.104497.

[16] D. H. Wolpert, The Lack of a Priori Distinctions between Learning Algorithms, Neural Comput. 8 (1996) 1341–1390. doi:10.1162/neco.1996.8.7.1341.

[17] S. Wager, T. Hastie, B. Efron, Confidence intervals for random forests: The jackknife and the infinitesimal jackknife, J. Mach. Learn. Res. 15 (2014) 1625–1651.

[18] J. L. Wu, J. X. Wang, H. Xiao, J. Ling, A Priori Assessment of Prediction Confidence for Data-Driven Turbulence Modeling, Flow, Turbul. Combust. 99 (2017) 25–46. doi:10.1007/s10494-017-9807-0.

[19] E. J. Parish, K. Duraisamy, A paradigm for data-driven predictive modeling using field inversion and machine learning, J. Comput. Phys. 305 (2016) 758–774. doi:10.1016/j.jcp.2015.11.012.

[20] N. Geneva, N. Zabaras, Quantifying model form uncertainty in Reynolds-averaged turbulence models with Bayesian deep neural networks, J. Comput. Phys. 383 (2019) 125–147. doi:10.1016/j.jcp.2019.01.021.

[21] C. A. Blauw, R. P. Dwight, Bayesian Additive Regression Trees for data-driven RANS turbulence modelling, Ph.D. thesis, Delft University of Technology, 2019.

[22] J. Hensman, N. Fusi, N. D. Lawrence, Gaussian processes for big data, in: Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, UAI'13, AUAI Press, Arlington, Virginia, USA, 2013, p. 282–290.

[23] B. Lakshminarayanan, A. Pritzel, C. Blundell, Simple and scalable predictive uncertainty estimation using deep ensembles, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17, Curran Associates Inc., Red Hook, NY, USA, 2017, p. 6405–6416.

[24] B. Lakshminarayanan, D. Roy, Y. W. Teh, Particle gibbs for bayesian additive regression trees, in: Artificial Intelligence and Statistics, 2015, pp. 553–561.

[25] B. Lakshminarayanan, D. M. Roy, Y. W. Teh, Mondrian Forests: Efficient Online Random Forests, in: Proc. 19th Int. Conf. Artif. Intell. Stat., Cadiz, Spain, 2016, pp. 1–15.

[26] B. Lakshminarayanan, D. M. Roy, Y. W. Teh, Mondrian forests for large-scale regression when uncertainty matters, Proc. 19th Int. Conf. Artif. Intell. Stat. AISTATS 2016 51 (2016) 1478–1487.

[27] H. Xiao, P. Cinnella, Quantification of Model Uncertainty in RANS Simulations: A Review, Prog. Aerosp. Sci. 108 (2019) 1–31. doi:10.1016/j.paerosci.2018.10.001.

[28] P. Tucker, Unsteady Computational Fluid Dynamics in Aeronautics, Springer Netherlands, 2014. doi:10.1007/978-94-007-7049-2.

[29] J. Wu, H. Xiao, R. Sun, Q. Wang, Reynolds-averaged Navier-Stokes equations with explicit data-driven Reynolds stress closure can be ill-conditioned, J. Fluid Mech. 869 (2019) 553–586. doi:10.1017/jfm.2019.205.

[30] L. Breiman, Random Forests, Mach. Learn. 45 (2001) 5–32. doi:10.3390/rs10060911.

[31] L. Breiman, J. Friedman, C. J. Stone, R. A. Olshen, Classification and Regression Trees, 1st ed., Chapman & Hall, 1984.

[32] T. D. Economon, Simulation and adjoint-based design for variable density incompressible flows with heat transfer, in: 2018 Multidisciplinary Analysis and Optimization Conference, American Institute of Aeronautics and Astronautics, 2018, pp. 1–24. doi:10.2514/6.2018-3111.

[33] F. R. Menter, M. Kuntz, R. Langtry, Ten Years of Industrial Experience with the SST Turbulence Model, Turbul. Heat Mass Transf. 4 4 (2003) 625–632.

[34] Y. Bentaleb, S. Lardeau, M. A. Leschziner, Large-eddy simulation of turbulent boundary layer separation from a rounded step, J. Turbul. (2012). doi:10.1080/14685248.2011.637923.

[35] J. Fröhlich, C. P. Mellen, W. Rodi, L. Temmerman, M. A. Leschziner, Highly resolved large-eddy simulation of separated flow in a channel with streamwise periodic constrictions, J. Fluid Mech. 526 (2005) 19–66. doi:10.1017/S0022112004002812.

[36] J. P. Laval, M. Marquillie, Direct Numerical Simulations of Converging–Diverging Channel Flow, Prog. Wall Turbul. Underst. Model. ERCOFTAC Ser. 14 (2011) 203–209.

[37] L. A. Schiavo, A. B. Jesus, J. L. Azevedo, W. R. Wolf, Large Eddy Simulations of convergent–divergent channel flows at moderate Reynolds numbers, Int. J. Heat Fluid Flow 56 (2015) 137–151. doi:10.1016/J.IJHEATFLUIDFLOW.2015.07.006.

[38] R. Vinuesa, P. S. Negi, M. Atzori, A. Hanifi, D. S. Henningson, P. Schlatter, Turbulent boundary layers around wing sections up to Rec=1,000,000, 2018. doi:10.1016/j.ijheatfluidflow.2018.04.017.

[39] J. Tyacke, P. Tucker, LES of heat transfer in electronics, Applied Mathematical Modelling 36 (2012) 3112–3133. doi:10.1016/j.apm.2011.09.072.

[40] S. B. Pope, Turbulent Flows, volume 1, 2000. doi:10.1088/1468-5248/1/1/702.

[41] S. Banerjee, R. Krahl, F. Durst, C. Zenger, Presentation of anisotropy properties of turbulence, invariants versus eigenvalue approaches, Journal of Turbulence 8 (2007) 1–32. doi:10.1080/14685240701506896.

[42] T.-H. Shih, J. Zhu, J. L. Lumley, A new reynolds stress algebraic equation model, Computer Methods in Applied Mechanics and Engineering 125 (1995) 287–302. doi:10.1016/0045-7825(95)00796-4.

[43] A. Hellsten, S. Wallin, Explicit algebraic reynolds stress and non-linear eddy-viscosity models, International Journal of Computational Fluid Dynamics 23 (2009) 349–361. doi:10.1080/10618560902776828.

[44] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems, volume 30, Curran Associates, Inc., 2017, pp. 4765–4774.

[45] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, S.-I. Lee, From local explanations to global understanding with explainable AI for trees, Nature Machine Intelligence 2 (2020) 56–67. doi:10.1038/s42256-019-0138-9.

[46] C. Gorlé, S. Zeoli, M. Emory, J. Larsson, G. Iaccarino, Epistemic uncertainty quantification for Reynolds-averaged Navier-Stokes modeling of separated flows over streamlined surfaces, Phys. Fluids 31 (2019). doi:10.1063/1.5086341.

[47] H. Xiao, J. X. Wang, R. G. Ghanem, A random matrix approach for quantifying model-form uncertainties in turbulence modeling, Comput. Methods Appl. Mech. Eng. (2017). doi:10.1016/j.cma.2016.10.025.

[48] J. Mourtada, S. Gaïffas, E. Scornet, AMF: Aggregated Mondrian Forests for Online Learning, arXiv:1906.10529v1 [stat.ML] (2019) 1–32.

[49] D. R. Roberts, V. Bahn, S. Ciuti, M. S. Boyce, J. Elith, G. Guillera-Arroita, S. Hauenstein, J. J. Lahoz-Monfort, B. Schröder, W. Thuiller, D. I. Warton, B. A. Wintle, F. Hartig, C. F. Dormann, Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure, 2017. doi:10.1111/ecog.02881.