

RESEARCH ARTICLE

Data-driven modelling of turbine wake interactions and flow resistance in large wind farms

Andrew Kirby^{*1} | François-Xavier Briol² | Thomas D. Dunstan³ | Takafumi Nishino¹

¹Department of Engineering Science,
University of Oxford, Oxford, UK

²Department of Statistical Science,
University College London, London, UK

³Informatics Lab, UK MetOffice, Exeter,
UK

Correspondence

*Andrew Kirby, Department of
Engineering Science, University of
Oxford, Oxford, OX1 3PJ, UK. Email:
andrew.kirby@trinity.ox.ac.uk

Abstract

Turbine wake and local blockage effects are known to alter wind farm power production in two different ways: (1) by changing the wind speed locally in front of each turbine; and (2) by changing the overall flow resistance in the farm and thus the so-called farm blockage effect. To better predict these effects with low computational costs, we develop data-driven emulators of the ‘local’ or ‘internal’ turbine thrust coefficient C_T^* as a function of turbine layout. We train the model using a multi-fidelity Gaussian Process (GP) regression with a combination of low (engineering wake model) and high-fidelity (Large-Eddy Simulations) simulations of farms with different layouts and wind directions. A large set of low-fidelity data speeds up the learning process and the high-fidelity data ensures a high accuracy. The trained multi-fidelity GP model is shown to give more accurate predictions of C_T^* compared to a standard (single-fidelity) GP regression applied only to a limited set of high-fidelity data. We also use the multi-fidelity GP model of C_T^* with the two-scale momentum theory (Nishino & Dunstan 2020, J. Fluid Mech. 894, A2) to demonstrate that the model can be used to give fast and accurate predictions of large wind farm performance under various mesoscale atmospheric conditions. This new approach could be beneficial for improving annual energy production (AEP) calculations and farm optimisation in the future.

KEYWORDS:

Class file; L^AT_EX 2_ε; Wiley NJD

1 | INTRODUCTION

The installed capacity of wind energy is projected to increase rapidly in the next decades. A major challenge in the optimisation of wind farm design is the accurate prediction of wind farm performance¹. Existing wind farm models struggle to make accurate predictions of wind farm power production. This is partly because the ‘global blockage effect’ reduces the velocity upstream of large farms and hence the energy yield². It remains unclear how global blockage should be modelled and this is the subject of a large-scale field campaign³.

Wind farms are typically modelled using engineering ‘wake’ models. These models predict the velocity deficit in the wakes behind turbines^{4,5}. To account for interactions between multiple turbines, the wake velocity deficits are superposed^{6,7}. Simple wake models can give predictions of wind farm performance with very low computational cost (10^{-3} CPU hours per simulation¹). However, wake

models do not account for the response of the atmospheric boundary layer (ABL) to the wind farm which is likely to be important for large wind farms⁸. It has been found that wake models compare poorly to Large-Eddy Simulations (LES) of large wind farms⁹.

Wind farms are also modelled in numerical weather prediction (NWP) models using farm parameterisation schemes. In these parameterisations, farms are often modelled as a momentum sink and a source of turbulent kinetic energy¹⁰. Turbine-wake interactions cannot be adequately predicted using these schemes. A new scheme was proposed¹¹ which uses a correction factor to model turbine interactions. More recently, data-driven approaches have been proposed¹² to model these effects in wind farm parameterisations.

Data-driven modelling of wind farm flows is a promising new approach¹³. Data from high-fidelity simulations with complex flow physics can be used to make predictions with low computational cost. Recent studies have applied machine learning techniques to data from a single turbine or from an existing wind farm. The data for these studies are from measurements^{14,15,16,17}, LES¹⁸ or Reynolds-Averaged Navier-Stokes (RANS) simulations^{19,20,21}. A limitation of these approaches is that they are not generalisable to different turbine layouts unless they rely on wake superposition techniques to model farm flows. Another approach is modelling the effect of turbine layout using geometric parameters¹⁷ or using the layout as a graph input to a neural network^{22,23}. However, these alternative approaches may struggle to fully capture the complex two-way interaction with the ABL as it seems impractical to prepare a data set that covers the entire range of scales involved in wind farm flows¹.

The problem of modelling wind farm flows can be split into 'internal' turbine-scale and 'external' farm-scale problems²⁴. The 'internal' problem is to determine a 'local' or 'internal' turbine thrust coefficient, C_T^* , which represents the flow resistance inside a wind farm, i.e., how the turbine thrust changes with wind speed within the farm. Nishino²⁵ proposed an analytical model for an upper limit of C_T^* by using an analogy to the classic Betz analysis. This analytical model is a function of turbine-scale induction factor but is independent of turbine layout and wind direction. Previous studies^{24,25,8} showed that C_T^* is usually lower than the limit predicted by Nishino's model and can vary significantly with turbine layout due to wake and turbine blockage effects.

The aim of this study is to develop statistical emulators of C_T^* as a function of turbine layout and wind direction. The novelty of this approach is that we are modelling the effect of turbine-wake interactions on C_T^* rather than turbine power. Both turbine-scale flows (e.g., wake effects) and farm-scale flows (e.g. farm blockage and mesoscale atmospheric response) affect turbine power within a farm. Therefore to create an emulator of turbine power, either (1) a very large set of expensive data such as finite-size wind farm LES is needed which covers a range of large-scale atmospheric conditions or (2) the model would not be generalisable to different mesoscale atmospheric responses. An emulator of C_T^* is however applicable to different atmospheric responses modelled separately, following the concept of the two-scale momentum theory^{24,8}.

In section 2 we give the definitions of key wind farm parameters in the two-scale momentum theory²⁴. Section 3 summarises the methodology of the LES and wake model simulations, followed by the machine learning approaches to develop the emulators in section 4. In section 5 we present the results from the trained emulators. These results are discussed in section 6 and concluding remarks are given in section 7.

2 | TWO-SCALE MOMENTUM THEORY

By considering the conservation of momentum for a control volume with and without a large wind farm over the land or sea surface, the following non-dimensional farm momentum (NDFM) equation can be derived²⁴,

$$C_T^* \frac{\lambda}{C_{f0}} \beta^2 + \beta^\gamma = M \quad (1)$$

where β is the farm wind-speed reduction factor defined as $\beta \equiv U_F/U_{F0}$ (with U_F defined as the average wind speed in the nominal wind farm-layer of height H_F , and U_{F0} is the farm-layer-averaged speed without the wind farm present); λ is the array density defined as $\lambda \equiv nA/S_F$ (where n is the number of turbines in the farm, A is the rotor swept area and S_F is the farm footprint area);

C_T^* is the internal turbine thrust coefficient defined as $C_T^* \equiv \sum_{i=1}^n T_i / \frac{1}{2} \rho U_F^2 n A$ (where T_i is thrust of turbine i in the farm and ρ is the air density); C_{f0} is the natural friction coefficient of the surface defined as $C_{f0} \equiv \langle \tau_{w0} \rangle / \frac{1}{2} \rho U_{F0}^2$ (where τ_{w0} is the bottom shear stress without the farm present); γ is the bottom friction exponent defined as $\gamma \equiv \log_\beta (\langle \tau_w \rangle / \tau_{w0})$ (where $\langle \tau_w \rangle$ is the bottom shear stress averaged across the farm); M is the momentum availability factor defined as,

$$M = \frac{\text{Momentum supplied by the atmosphere to the farm site **with** turbines}}{\text{Momentum supplied by the atmosphere to the farm site **without** turbines}}. \quad (2)$$

noting that this includes pressure gradient forcing, Coriolis force, net injection of streamwise momentum through top and side boundaries and time-dependent changes in streamwise velocity²⁴. The height of the farm-layer, H_F , is used to define the reference velocities U_F and U_{F0} . Equation 1 is valid so long as the same of H_F is used for both the internal and external problem. H_F is typically between $2H_{hub}$ and $3H_{hub}$ ⁸ (where H_{hub} is the turbine hub-height) and in this study we use a fixed definition of $H_F = 2.5H_{hub}$.

Patel²⁶ used an NWP model to demonstrate that, for most cases, M varied almost linearly with β (for a realistic range of β between 0.8 and 1). Therefore, M can be approximated by

$$M = 1 + \zeta(1 - \beta) \quad (3)$$

where ζ is the 'momentum response' factor or 'wind extractability' factor. Patel²⁶ found ζ to be time-dependent and vary between 5 and 25 for a typical offshore site (note that $\zeta = 0$ corresponds to the case where momentum available to the farm site is assumed to be fixed, i.e., $M = 1$).

Nishino²⁵ proposed an analytical model for C_T^* given by,

$$C_T^* = 4\alpha(1 - \alpha) = \frac{16C'_T}{(4 + C'_T)^2} \quad (4)$$

where α is the turbine-scale wind speed reduction factor defined as $\alpha \equiv U_T / U_F$ (U_T is the streamwise velocity averaged over the rotor swept area) and $C'_T \equiv T / \frac{1}{2} \rho U_T^2 A$ is a turbine resistance coefficient describing the turbine operating conditions.

For a given farm configuration at a farm site (i.e., for given set of C_T^* , λ , C_{f0} , γ and ζ) the farm wind-speed reduction factor β can be calculated using equation 1. The (farm-averaged) power coefficient C_p is defined as $C_p \equiv \sum_{i=1}^n P_i / \frac{1}{2} \rho U_{F0}^3 n A$ (P_i is power of turbine i in the farm). Using the calculated value of β , C_p can be calculated by using the expression,

$$C_p = \beta^3 C_p^* \quad (5)$$

where C_p^* is the (farm-averaged) 'local' or 'internal' turbine power coefficient defined as $C_p^* \equiv \sum_{i=1}^n P_i / \frac{1}{2} \rho U_F^3 n A$.

3 | WIND FARM SIMULATIONS

In this study we model wind farms as arrays of actuator discs (or aerodynamically ideal turbines operating below the rated wind speed). This is because, in real wind farms, the effects of turbine wake interactions on the farm performance are most significant when they operate below the rated wind speed. The 'internal' thrust coefficient C_T^* is an important wind farm parameter which includes the effect of turbine interactions (including both wake and local blockage effects). In this study we will be modelling the effect of turbine layout on C_T^* for aligned turbine layouts with various wind directions and a fixed turbine resistance of $C'_T = 1.33$. We chose $C'_T = 1.33$ because it leads to a turbine induction factor of 1/4 which is close to a typical value for modern large wind turbines. As such we will be considering

$$C_T^* = f(S_x, S_y, \theta) \quad (6)$$

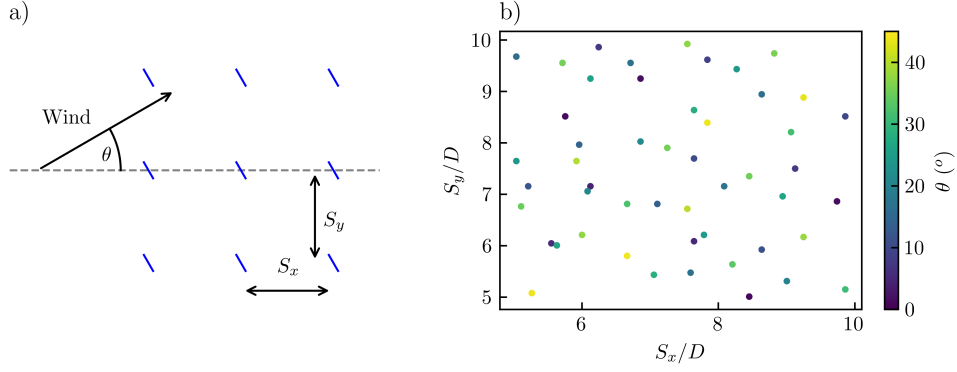


Figure 1 Design of numerical experiments: a) input parameters, b) maximin design of LES.

where S_x is the turbine spacing in the x direction, S_y is the turbine spacing in the y direction and θ is the wind direction relative to the x direction (see figure 1a). However the true function C_T^* cannot be easily evaluated so we will instead investigate C_T^* using computer codes. One computer code we will use is LES (see section 3.1) to estimate C_T^*

$$C_{T,LES}^* = f_{LES}(S_x, S_y, \theta). \quad (7)$$

We assume that the function f_{LES} is close to the true function f because of the accuracy of LES to model wind farm flows. We will also use a wake model (see section 3.2) to provide cheap approximations of C_T^* according to

$$C_{T,wake}^* = f_{wake}(S_x, S_y, \theta). \quad (8)$$

Engineering problems are often investigated using complex computer models. Evaluating the output of such computer models for a given input can be very computationally expensive. Therefore a common objective is to create a cheap statistical model of the expensive computer model; this is commonly known as emulation of computer models^{27,28}. In this study we aim to develop a statistical emulator which can cheaply emulate f_{LES} .

The emulators will only be valid for aligned layouts of wind turbines and for a given turbine resistance (here we use $C'_T = 1.33$). We consider the input parameters for a realistic range of turbine spacings¹: $S_x \in [5D, 10D]$, $S_y \in [5D, 10D]$ and $\theta \in [0^\circ, 45^\circ]$ where D is the diameter of the turbine rotor swept area. In this study D is set as 100m and the turbine hub height is also 100m. We only need to consider wind directions of $\theta \in [0^\circ, 45^\circ]$ because of symmetry in the aligned turbine layouts. If θ is negative than the turbine layout given by (S_x, S_y, θ) is exactly the same as $(S_x, S_y, -\theta)$. When $\theta > 45^\circ$, then (S_x, S_y, θ) and $(S_y, S_x, 90^\circ - \theta)$ give identical layouts.

In this study we build several emulators to predict f_{LES} . The models are trained using data from low-fidelity (wake model) and high fidelity (LES) wind farm simulations. One evaluation of $C_{T,wake}^*$ takes approximately 130 seconds on a single CPU and $C_{T,LES}^*$ requires around 400 CPU hours on a supercomputer. We use a space filling maximin design^{29,30} to select training points in the parameter space. The maximin algorithm selects points which maximises the minimum distance to other points and to the boundaries. This provides a good coverage of the domain which ensures that the emulators can give good predictions across the whole of the domain³¹. Figure 1b shows the LES training points in the parameter space.

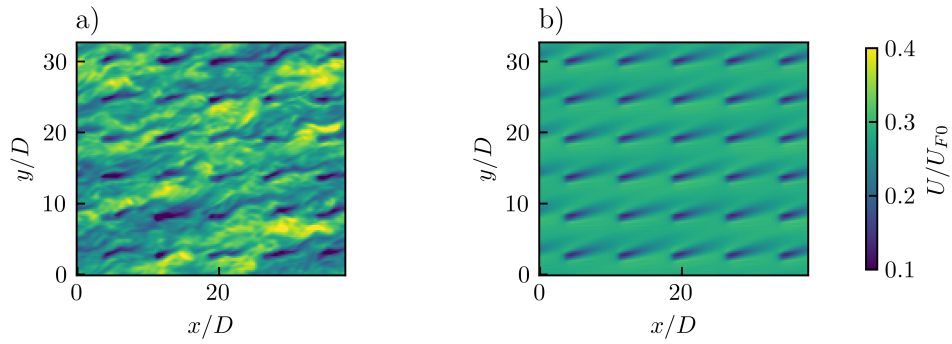


Figure 2 LES a) instantaneous and b) time-averaged flow fields over a periodic turbine array ($S_x/D = 7.59$, $S_y/D = 5.47$ and $\theta = 37.6^\circ$).

3.1 | Large-Eddy Simulations

This study uses the data from 50 high-fidelity (LES) simulations of wind farms published in a previous study⁸. Here we give a brief summary of the LES methodology. The LES models a neutrally stratified atmospheric boundary layer over a periodic array of actuator discs, which face the wind direction θ and exert uniform thrust. The resolution is 24.5m in the horizontal directions (4 points across the rotor diameter) and 7.87m in the vertical. This is a coarse horizontal resolution; however using a correction factor for the turbine thrust³² makes the $C_{T,LES}^*$ values insensitive to horizontal resolution⁸. For all simulations the vertical domain size was fixed at 1km and the horizontal extent varied with turbine layout but was at least 3.14km. The horizontal boundary conditions were periodic (essentially an infinitely-large wind farm). The bottom boundary used a no-slip condition with the value of eddy viscosity specified following the Monin-Obukhov similarity theory for a surface roughness length of $z_0 = 1 \times 10^{-4}$ m. The top boundary had a slip condition with zero vertical velocity. The flow was driven by a pressure gradient forcing which was constant and in the direction θ throughout the domain. Figure 2 shows the instantaneous and time-averaged hub height velocities from one wind farm LES. See the original paper⁸ for further details of the LES.

3.2 | Wake model simulations

Wake models are a cheap low-fidelity approach to modelling wind farm aerodynamics compared to expensive high-fidelity LES simulations¹. We use the wake model proposed by Niayafar and Porté-Agel³³ to evaluate $C_{T,wake}^*$ as a cheap approximation of C_T^* . We use the Python package PyWake³⁴ to implement the wake model. The turbine thrust coefficient C_T is needed as an input for the wake model. We use the value of C_T^* predicted by equation 4 as the value of C_T . For the turbine operating conditions used in this study ($C_T' = 1.33$) the wake model has C_T equal to 0.75 for all turbines. To model actuator discs, we consider a hypothetical turbine which has a constant C_T for all wind speeds. We calculate $C_{T,wake}^*$ for a single turbine at the back of a large farm (marked X in figure 3). The farm simulated using the wake model is 10km long in the streamwise direction and 4km long in the cross-streamwise direction. The farm size was chosen so that C_T^* no longer varied with increasing farm size. The wake growth parameter is calculated using $k^* = 0.38I + 0.004$ where I is the local streamwise turbulence intensity. The local streamwise turbulence intensity is estimated using the model proposed by Crespo and Hernández³⁵. The background turbulence intensity (TI) is set as a typical value of 10%.

The velocity incident to the turbine is calculated by averaging the velocity across the disc area. We use a 4×3 cartesian grid with Gaussian quadrature coordinates and weights on the disc to average the velocity. The disc-averaged velocity, U_T is then calculated by multiplying the averaged incident velocity by $(1 - a)$ where a is the turbine induction factor set by the value of C_T' (using the expression $a = C_T'/(4 + C_T')$). To calculate the farm-average velocity, U_F , we average the velocity across a volume around the

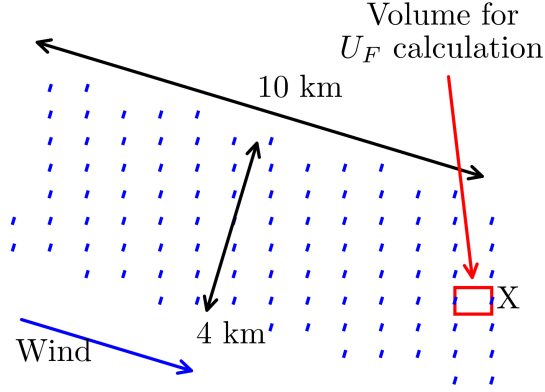


Figure 3 Example of wind farm layout for wake model simulations.

single turbine. The volume has dimensions of S_y in the y direction, S_x in the x direction and 250m in the z direction (the height of the nominal farm layer used in the previous LES study⁸). To calculate the average velocity, we discretise the volume into 200 points in the horizontal directions and 20 points in the vertical. This was sufficient for the calculation of $C_{T,wake}^*$ to not vary with further discretisation. Figure 3 shows an example of the farm layout for the wake model simulations.

4 | MACHINE LEARNING METHODOLOGY

4.1 | Gaussian Process regression

We will use Gaussian process (GP) regression³⁶ to build statistical emulators of f_{LES} . A Gaussian process is a stochastic process $g \sim \mathcal{GP}(m, k)$ described by a mean function $m(v) = \mathbb{E}[g(v)]$ and a covariance function $k(v, v') = \mathbb{E}[(g(v) - m(v))(g(v') - m(v'))]$. In our case $v = (S_x, S_y, \theta)$. We will use such a stochastic process as a model of f_{LES} , the true mapping from v to $C_{T,LES}^*$. Each realisation from this process will therefore be a function which could plausibly represent this mapping. The mean function represents the expected output value at an input $v = (S_x, S_y, \theta)$. The covariance function gives the covariance between output values at v and v' . Examples of covariance functions include squared exponential, rational quadratic and periodic functions³⁶. Different covariance functions will give differently shaped GPs. For example the squared exponential covariance function will give very smooth GPs whereas the periodic function will give GPs with a periodic structure. Other types of structure, for example symmetry, can also be encoded in the covariance function. Therefore the expected shape (for example smoothness) of the expected relationship and any properties (for example discontinuities or symmetries) need to be considered when choosing a covariance function for GP regression.

Let $V = (v_1, \dots, v_n)^T$ be a collection of design points then $m_V = (m(v_1), \dots, m(v_n))^T$ is the mean vector and $k_{VV} = (k(v_i, v_j))$ is the covariance matrix. We will start by positing a GP model with mean m and covariance k (called the ‘prior GP’), then condition this GP on LES observations; the outcome is a new GP (called the ‘posterior GP’). This gives the posterior distribution $g|V, C_{T,LES}^* \sim \mathcal{GP}(\bar{m}_{\sigma^2}, \bar{k}_{\sigma^2})$. \bar{m}_{σ^2} is the posterior mean function given by $\bar{m}_{\sigma^2}(v) = m(v) + k_{vV}(k_{VV} + \sigma^2 I_{n \times n})^{-1}(C_{T,LES}^* - m_V)$ where $k_{vV} = (k(v, v_1), \dots, k(v, v_n))$ and $I_{n \times n}$ is the identity matrix of size n . The posterior mean function \bar{m}_{σ^2} is used to make predictions at $v = (S_x, S_y, \theta)$. The posterior covariance function \bar{k}_{σ^2} quantifies the uncertainty in our prediction at $v = (S_x, S_y, \theta)$. The posterior covariance function is given by $\bar{k}_{\sigma^2}(v, v') = k(v, v') - k_{vV}(k_{VV} + \sigma^2 I_{n \times n})^{-1}k_{Vv'}$.

Often in GP regression a zero prior mean is used. However, using an informative prior mean can improve the accuracy of the trained model. By using a prior mean, many of the trends in f_{LES} can be incorporated into our model prior to making expensive

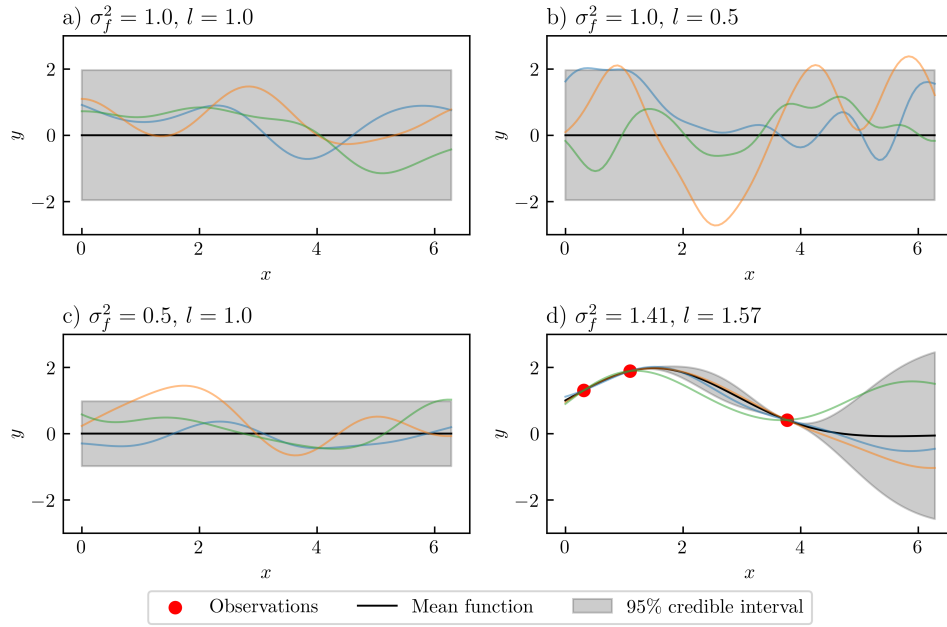


Figure 4 Demonstration of basic GP regression: a) shows the prior mean and covariance function prior to fitting with 3 GPs drawn from the distribution shown in colour; b) shows the effect of decreasing the lengthscale hyperparameter; c) the effect of variance hyperparameter; and d) the posterior mean and covariance functions.

evaluations of $C_{T,LES}^*$. Therefore, after training our model will likely better describe the true relationship between S_x, S_y, θ and f_{LES} . In this study, we will use both $C_{T,wake}^*$ and the analytical model of C_T^* as the prior mean for the standard GP regression. For the wake model prior mean we also vary the specified ambient TI input parameter.

We expect f_{LES} to be a smooth function of input variables S_x, S_y and θ , and to vary more rapidly with θ than S_x or S_y . Therefore we will use an anisotropic squared-exponential covariance function,

$$k(v, v') = \sigma_f^2 \exp\left(-\frac{(S_x - S'_x)^2}{2l_1^2}\right) \exp\left(-\frac{(S_y - S'_y)^2}{2l_2^2}\right) \exp\left(-\frac{(\theta - \theta')^2}{2l_3^2}\right) \quad (9)$$

where $\sigma_f^2 > 0$ is the signal variance hyperparameter and $l_i > 0$ is the lengthscale hyperparameter for each dimension. This is also called an ARD (automatic relevance detection) kernel. If we consider $v = v'$ then we can see that σ_f^2 determines the variance of $g(v)$. Therefore σ_f^2 determines the prior uncertainty the model has about the value of $g(v)$. As the lengthscale hyperparameter l_i gets smaller then $k(v, v')$ decreases (for $v \neq v'$). Equally if l_i increases then $k(v, v')$ will also increase. A GP with a small l_i will therefore vary more rapidly across the parameter space in the i th dimension.

Due to numerical issues associated with the matrix inversion/linear system solve operations in the formulae for the posterior GP, it is common to add a nugget $\sigma^2 > 0$ to the kernel matrix. The hyperparameters σ_f^2 and l_i are selected automatically during the fitting process by maximising the log marginal likelihood³⁶. This approach selects the model which maximises the fit to the data.

Figure 4 shows the impact of the hyperparameters in an example GP regression setting (using the squared exponential covariance function). The mean function and 95% credible interval (± 1.96 times the standard deviation) prior to fitting are shown in figure 4a with 3 GPs drawn from the distribution (coloured lines). The effect of decreasing the lengthscale hyperparameter l_i is shown in figure 4b. The prior mean and 95% credible interval are unchanged however the example GPs drawn vary more rapidly because of the shorter lengthscale. Figure 4c shows the same setup as figure 4a but with a smaller value of σ_f^2 . The example GPs still vary slowly but the magnitude of the variations is now smaller. Figure 4d shows the GPs conditioned on observations with hyperparameters selected by maximising the log marginal likelihood.

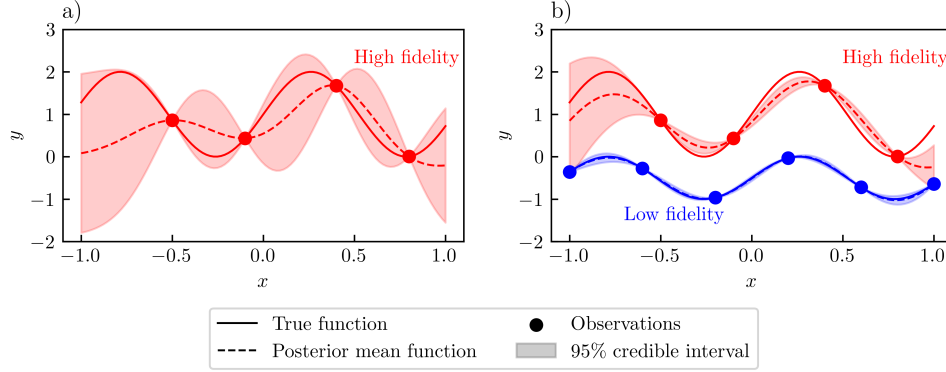


Figure 5 Demonstration of a) basic GP regression and b) multi-fidelity GP regression. In this example $f(x) = 1 + \sin(6x)$ for the high-fidelity data and $f(x) = -0.5 + 0.5\sin(6x)$ for the low-fidelity data.

4.2 | Non-linear multi-fidelity Gaussian Process regression

In many applications there are several computational models available. These models can have varying accuracies and computational costs. The models which are more computationally expensive typically give more accurate predictions. The GP regression framework can be extended to combine information from low and high-fidelity models³⁷. This type of modelling uses the low-fidelity observations to speed up the learning process and the high-fidelity observations to ensure accuracy. In our scenario we will combine evaluations of from a low-fidelity ($C_{T,wake}^*$) and a high-fidelity ($C_{T,LES}^*$) model. Note that for the multi-fidelity models in this study we set the ambient TI to 10% for the wake model and use a zero prior mean. We will keep the number of high-fidelity training points fixed at 50 and we will vary the number of low-fidelity training points used.

We combine information from our high and low-fidelity models using a nonlinear information fusion algorithm³⁸. The framework is based on the autoregressive multi-fidelity scheme given by:

$$g_{high}(v) = \rho(g_{low}(v)) + \delta(v) \quad (10)$$

where $g_{low}(v)$ is a model with a GP denoted f_{wake} and $g_{high}(v)$ is a model with a GP denoted f_{LES} . ρ is a model with a GP which maps the low-fidelity output to the high-fidelity output and $\delta(v)$ is a model with a GP which is a bias term. The non-linear multi-fidelity framework can learn non-linear space-dependent correlations between models of different accuracies. To reduce the computational cost and complexity of implementation the autoregressive scheme given by equation 10 is simplified. Firstly, the GP prior $g_{low}(v)$ is replaced by the GP posterior $g_{low,*}(v)$ and secondly the GPs ρ and δ are assumed to be independent. Equation 10 can then be summarised as

$$g_{high}(v) = h_{high}(v, g_{low,*}(v)) \quad (11)$$

where h_{high} is a model with a GP which has both v and $g_{low,*}(v)$ as inputs. More details of h_{high} and the implementation of the multi-fidelity framework are given in Perdikaris *et. al.*³⁸.

Figure 5 shows an example of how a multi-fidelity GP can outperform a standard GP regression. We implement the non-linear multi-fidelity framework using the ‘emukit’ package³⁹. We first maximise the log marginal likelihood whilst keeping the Gaussian noise variance fixed at a low value of 1×10^{-6} . The fitting process is then repeated whilst allowing the Gaussian noise variance to be optimised too. This is to prevent a high noise local optima from being selected.

5 | RESULTS

In this study, we build various statistical emulators of f_{LES} using different techniques and compare the performance. A summary of the techniques is shown in the list below:

- 1 Standard Gaussian Process regression (see section 4.1)
 - a **GP-analytical-prior**: Gaussian Process using analytical model (equation 4) prior mean
 - b **GP-wake-TI10-prior**: Gaussian Process using wake model (section 3.2) with ambient TI=10% prior mean
 - c **GP-wake-TI1-prior**: Gaussian Process using wake model with ambient TI=1% prior mean
 - d **GP-wake-TI5-prior**: Gaussian Process using wake model with ambient TI=5% prior mean
 - e **GP-wake-TI15-prior**: Gaussian Process using wake model with ambient TI=15% prior mean
- 2 Non-linear multi-fidelity Gaussian Process regression (see section 4.2)
 - a **MF-GP-nlow500**: multi-fidelity Gaussian Process using 500 low-fidelity training points
 - b **MF-GP-nlow250**: multi-fidelity Gaussian Process using 250 low-fidelity training points
 - c **MF-GP-nlow1000**: multi-fidelity Gaussian Process using 1000 low-fidelity training points

The code used to produce the results in this section is available open-access at the following GitHub repository: https://github.com/AndrewKirby2/ctstar_statistical_model.

5.1 | Performance of standard GP regression

We first assessed the accuracy of the standard GP models (section 4.1) by performing leave-one-out cross-validation (LOOCV). This is a method of estimating the accuracy of a statistical model when making predictions on data not used to train the model. We trained our model on 49 of the 50 training points and then calculated the prediction accuracy for the single high-fidelity data point which is excluded from the training set. This is then repeated for all data points in turn, and we took the average accuracy as an estimate of the model test accuracy. The standard GP models were implemented using the 'GPY' package⁴⁰.

The standard GP gave accurate predictions of f_{LES} with average errors of less than 2%. Table 1 shows the accuracy of the standard GP models compared to the analytical and wake models. We calculated the errors by using the expression $|\overline{m}_{\sigma^2} - C_{T,LES}^*|/0.75$ where \overline{m}_{σ^2} is the posterior mean function of the emulator. The reference value for C_T^* of 0.75 was chosen because this is the prediction from the analytical model. Both GP models give similar maximum errors of approximately 6%. Using the wake model as a prior mean gave a lower mean absolute error of 1.26%. The GP models reduced the average prediction error and significantly reduced the maximum error compared to the wake model and analytical model of C_T^* .

Table 1 Accuracy of models for C_T^* prediction.

Model	MAE (%)	Maximum error (%)
GP-analytical-prior	1.87	6.09
GP-wake-TI10-prior	1.26	6.11
Analytical model	5.26	22.0
Wake model (TI=10%)	4.60	9.28

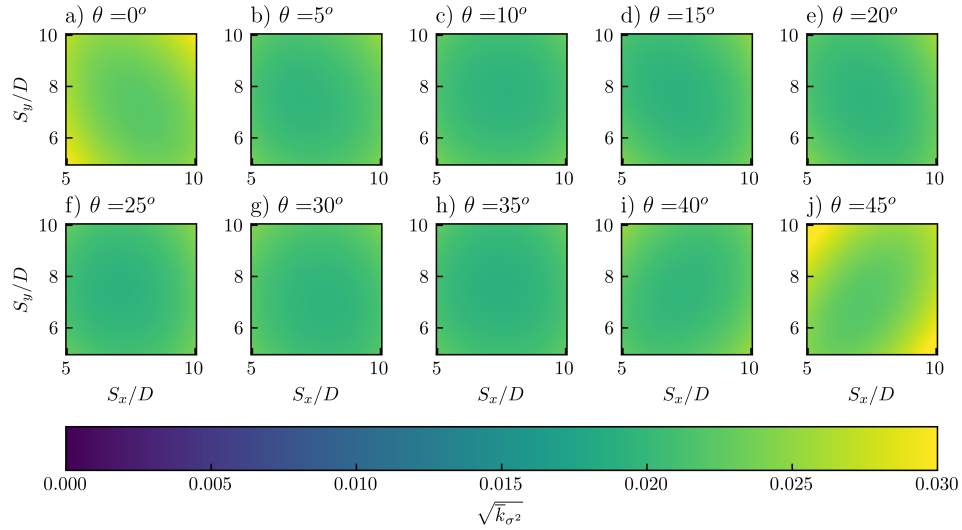


Figure 6 Posterior variance function of GP-wake-TI10-prior model.

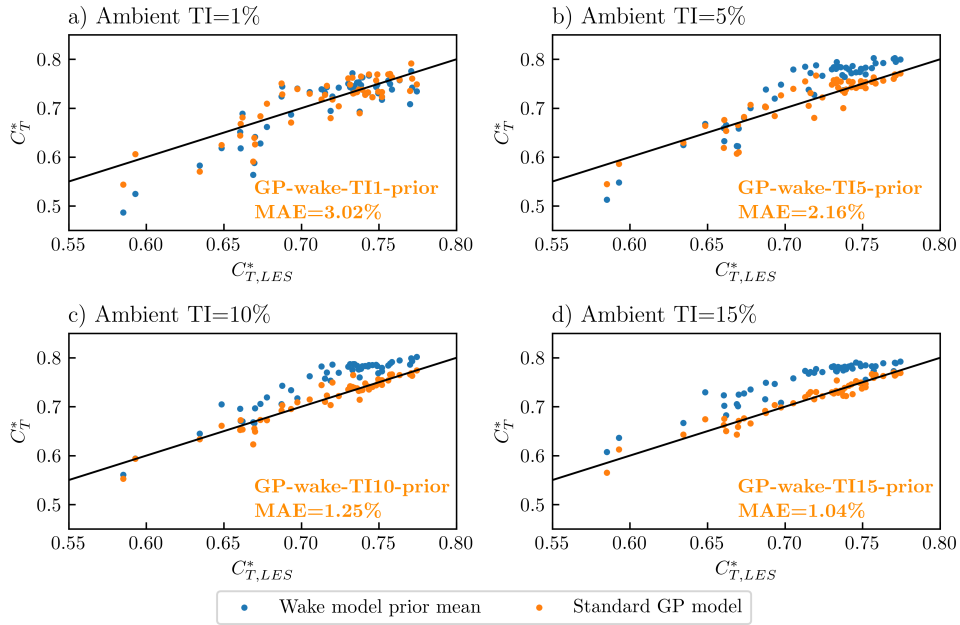


Figure 7 Sensitivity of fitted GP models to the ambient TI chosen for wake model prior means.

The model **GP-wake-TI10-prior** has a high degree of confidence when making predictions in regions of the parameter space. Figure 6 shows the square root of the posterior covariance function \bar{k}_{σ^2} , which quantifies the uncertainty of the emulator. The uncertainty is uniform throughout the parameter space with regions of slightly higher uncertainty at $\theta = 0^\circ$ and 45° .

We also assessed the sensitivity of the model accuracy to the ambient TI used in the wake model prior mean. Figure 7 shows the impact of ambient TI on the wake model prior mean and the fitted GP model. Increasing the ambient TI increased the value of $C_{T,wake}^*$. This is because of the enhanced wake recovery behind wind turbines. Increasing the ambient TI in the wake model results in $C_{T,wake}^*$ overpredicting $C_{T,LES}^*$. The MAE from the LOOCV procedure for each fitted GP is shown in the bottom right corner.

The fitted GPs became more accurate when the wake model ambient TI was increased. Increasing the ambient TI for the wake model causes the wakes to recover faster. The wakes become shorter in the streamwise direction and wider in the spanwise direction. As such, $C_{T,wake}^*$ becomes less sensitive to the turbine layout. When an ambient TI of 1% and 5% is used for the wake model, $C_{T,wake}^*$ is more sensitive to turbine layout than $C_{T,LES}^*$ (figures 7a and 7b). When the ambient TI is increased to 10% and above, the relationship between $C_{T,wake}^*$ and $C_{T,LES}^*$ becomes simpler (figures 7c and 7d). This seems to explain why the fitted GPs become more accurate.

5.2 | Performance of non-linear multi-fidelity GP regression

We then assessed the accuracy of the multi-fidelity GP models (section 4.2). All models used the 50 high-fidelity ($C_{T,LES}^*$) training points and a varying number of low-fidelity ($C_{T,wake}^*$) training points (using an ambient TI of 10% for $C_{T,wake}^*$). The results from LOOCV are shown in table 2. For the LOOCV we train our model on 49 out of the 50 high-fidelity data points and all low-fidelity data points. Then we average the error in predicting the high-fidelity data point left of the training set and repeat this in turn for data points. Increasing the number of low-fidelity training points from 250 to 500 reduced the mean and maximum error. However, increasing this to 1000 low-fidelity training points did not increase accuracy and increased the fitting and prediction time. This is because the number of high-fidelity training points is fixed. There is a threshold where the model of the relationship between f_{LES} and f_{wake} , denoted ρ , limits the final accuracy of the emulator of f_{LES} .

The posterior mean \bar{m}_{σ^2} of $g_{low}(v)$ is an emulator of f_{wake} and $g_{high}(v)$ is an emulator of f_{LES} . Figure 8 gives the predictions from the posterior mean of $g_{high}(v)$ (for **MF-GP-nlow500**). The lowest \bar{m}_{σ^2} values were for a wind direction of $\theta = 0^\circ$. \bar{m}_{σ^2}

Table 2 Performance of the multi-fidelity Gaussian Process models.

Model	MAE (%)	Maximum error (%)	Training time (s)	Prediction time (s)
MF-GP-nlow250	1.46	7.12	6.15	0.00157
MF-GP-nlow500	0.828	3.75	9.73	0.00167
MF-GP-nlow1000	0.866	3.55	26.8	0.00236

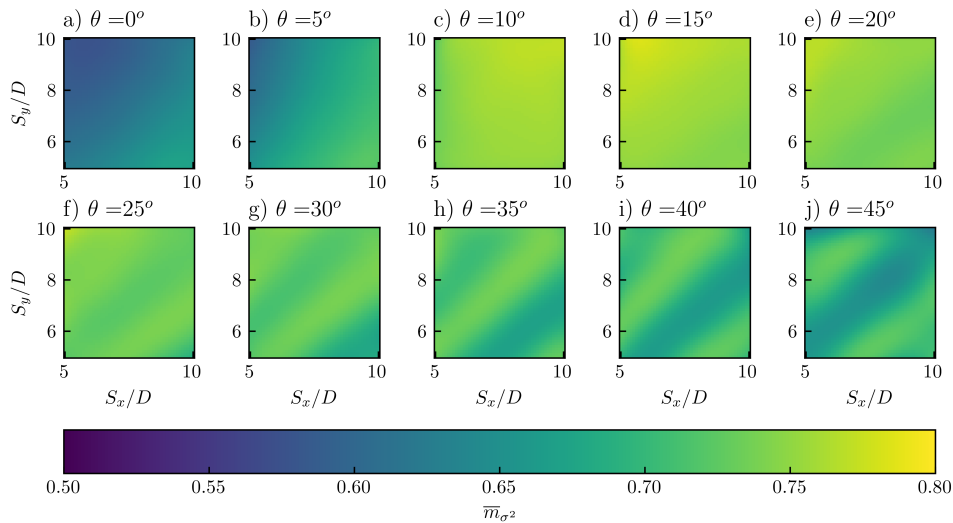


Figure 8 Posterior mean function for $g_{high}(v)$ of **MF-GP-nlow500**.

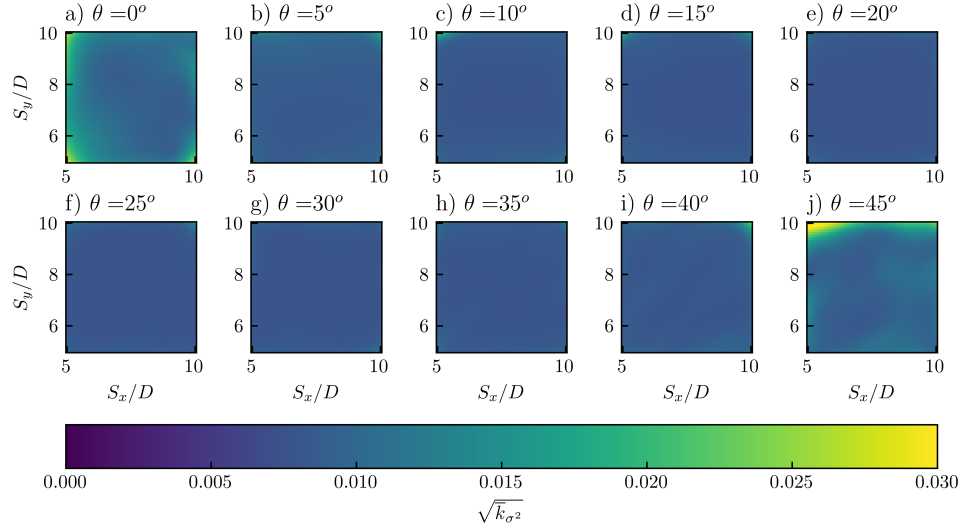


Figure 9 Posterior variance function for $g_{high}(v)$ of **MF-GP-nlow500**.

increased rapidly with θ reaching a maximum of slightly over 0.75 at $\theta = 10^\circ$. For large values of θ (above $\theta = 25^\circ$) there were local minima in \overline{m}_{σ^2} which appear in figure 8 as diagonal strips of low \overline{m}_{σ^2} values. The main diagonal strip occurs along the line of $S_y = S_x \tan(\theta)$. There are two smaller strips either side of with positions given by $S_y = 2 \tan(\theta)$ and $S_y = 0.5 \tan(\theta)$ (this is discussed further in section 6).

The uncertainty the model **MF-GP-nlow500** has in predicting f_{LES} is shown in figure 9. The model uncertainty is uniform throughout the parameter space with slightly higher values at $\theta = 0^\circ$ and 45° . Compared to the posterior variance of **GP-wake-TI10-prior** (shown in figure 6) the uncertainty is lower. By incorporating information from $C_{T,wake}^*$, the multi-fidelity GP model has more confidence about predicting f_{LES} .

The prediction errors from the LOOCV (for **MF-GP-nlow500**) are shown in figure 10. The box plot of prediction errors in figure 10a shows that this model had no significant bias whereas both the wake and analytical models systemically overestimated $C_{T,LES}^*$. Figures 10b-d show that for the statistical model there appears to be no part of the parameter space which had larger errors.

The multi-fidelity approach used in this study builds a statistical model of both the low-fidelity (f_{wake}) and high-fidelity (f_{LES}) model. We can use the posterior means of $g_{low}(v)$ and $g_{high}(v)$ to see the differences between the wake model and LES. The posterior mean for both models are shown in figure 11. For the wake model the change in \overline{m}_{σ^2} with θ is greater than for the LES (especially between $\theta = 0^\circ$ and 10°). For larger values of θ , there is a larger difference in \overline{m}_{σ^2} between waked and unwaked layouts for the low-fidelity model compared to the high-fidelity one. This suggests that the wake model is more sensitive to changes in wind directions than the LES.

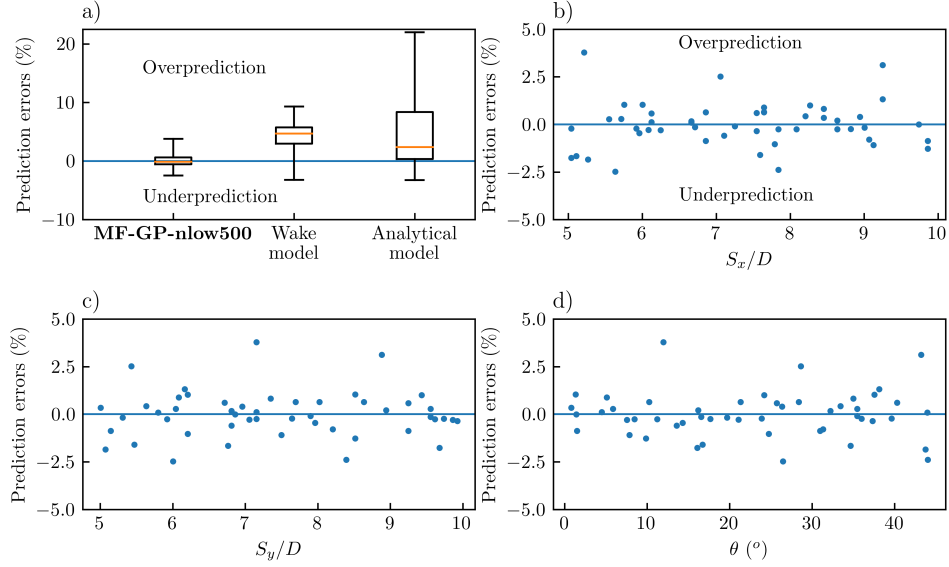


Figure 10 Comparison of LOOCV prediction errors (%) for different models a) and LOOCV prediction error (%) of MF-GP-nlow500 against input parameters b) S_x/D , c) S_y/D and d) $\theta(^{\circ})$. Note that for the box plot in a) the orange line is the median LOOCV error and the box is the interquartile range of LOOCV error.

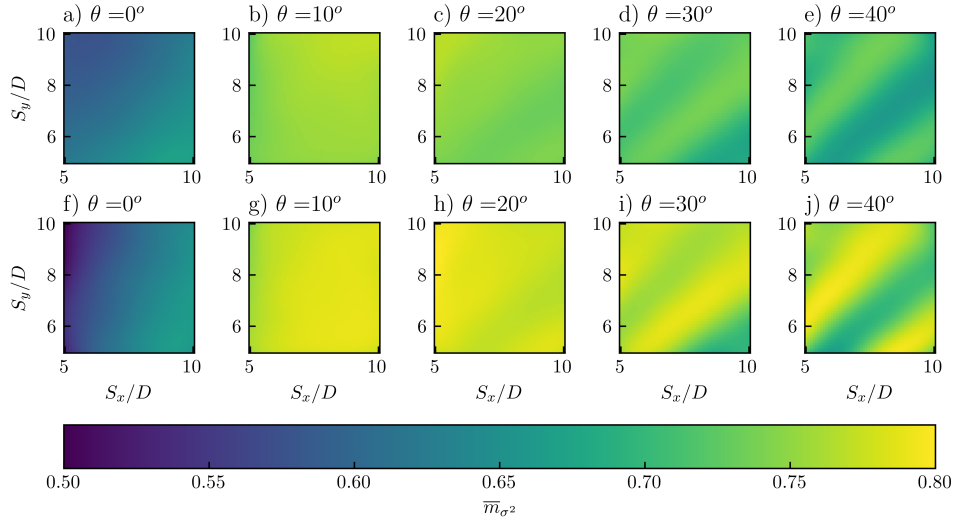


Figure 11 Posterior mean function of MF-GP-nlow500 for different values of θ for a) to e) $g_{high}(v)$ and f) to j) $g_{low}(v)$.

5.3 | Prediction of wind farm performance

We use the predicted values of $C_{T,LES}^*$ from the emulators to predict the power output of wind farms under various mesoscale atmospheric conditions, following the concept of the two-scale momentum theory. We predict the (farm-averaged) turbine power coefficient C_p using $C_{T,LES}^*$ predictions from MF-GP-nlow500. We call this prediction of farm performance $C_{p,model}$. Firstly, we use the $C_{T,LES}^*$ prediction from the LOOCV procedure as C_T^* in equation 1 to calculate β for a given value of wind extractability ζ . We substitute this value of β into the expression $C_p = \beta^3 C_T^{*\frac{3}{2}} C_T'^{-\frac{1}{2}}$ (which is only valid for actuator discs) to calculate $C_{p,model}$. We compare the value of $C_{p,model}$ with the turbine power coefficient recorded in the LES, $C_{p,LES}$. The effect of the coarse LES

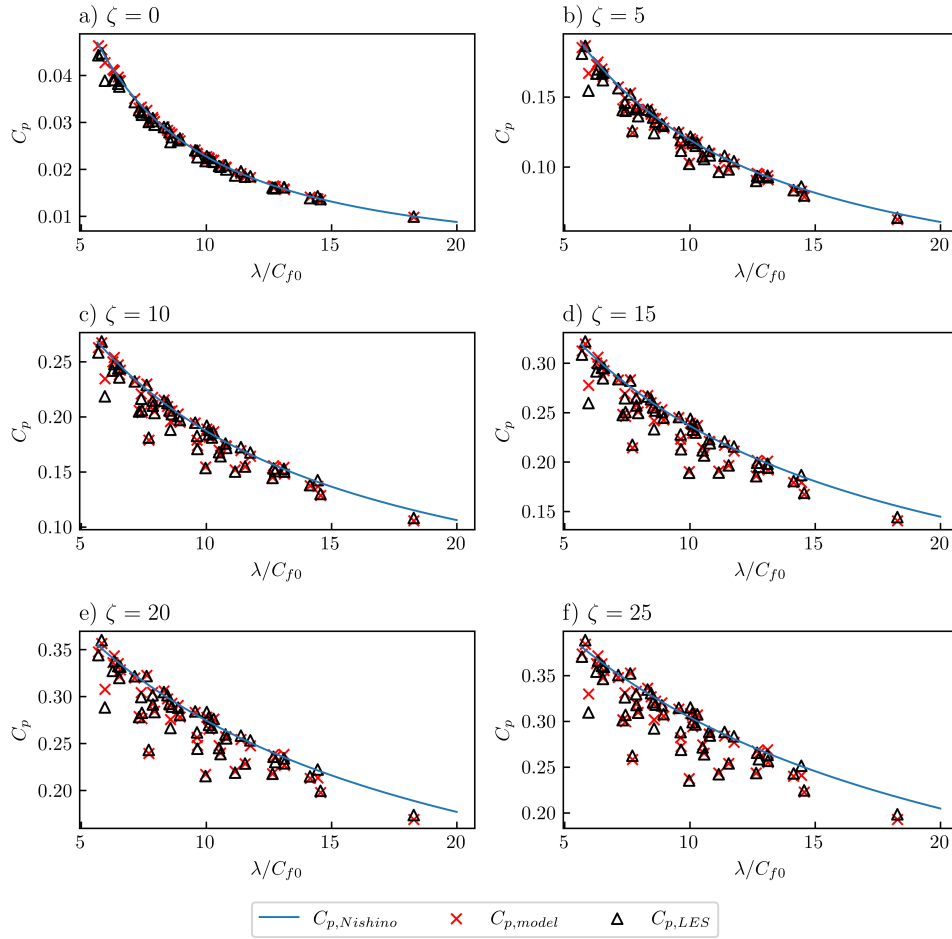


Figure 12 Comparison of C_p predictions with LES results for a realistic range of ζ values.

resolution on turbine thrust (and hence also ABL response and C_p) has already been corrected⁸. The LES was performed with periodic horizontal boundary conditions and a fixed momentum supply, i.e., $\zeta = 0$. However, the $C_{p,LES}$ has also been adjusted for a given ζ by scaling the velocity fields assuming Reynolds number independence⁸.

Similarly, the analytical model of C_T^* can be used to give a theoretical prediction of wind farm performance called $C_{p,Nishino}$ ⁸, which is given by

$$C_{p,Nishino} = \frac{64C_T'}{(4 + C_T')^3} \left[\frac{-\zeta + \sqrt{\zeta^2 + 4 \left(\frac{16C_T'}{(4 + C_T')^2} \frac{\lambda}{C_{f0}} + 1 \right) (1 + \zeta)}}{2 \left(\frac{16C_T'}{(4 + C_T')^2} \frac{\lambda}{C_{f0}} + 1 \right)} \right]^3. \quad (12)$$

We will compare the accuracy of both $C_{p,model}$ and $C_{p,Nishino}$ in predicting $C_{p,LES}$.

Both $C_{p,model}$ and $C_{p,LES}$ are shown in figure 12 for a realistic range of wind extractability factors, along with the results from $C_{p,Nishino}$ (equation 12). $C_{p,Nishino}$ provides an approximate upper limit of farm-averaged C_p as it predicts very well the effects of array density and large-scale atmospheric response. The statistical model accurately predicts the effect of turbine layout on farm performance which becomes more important with larger ζ values. As ζ increases, there is a larger difference between $C_{p,LES}$ and $C_{p,Nishino}$. Also, $C_{p,model}$ becomes slightly less accurate when ζ increases.

Table 3 shows the average prediction errors of $C_{p,model}$ and $C_{p,Nishino}$. We quantified the mean absolute error using two different reference powers. Using $C_{p,LES}$ as the reference power, $C_{p,Nishino}$ had an error of around 5% and the error increases

Table 3 Comparison of models for C_p prediction.

$\frac{1}{50} \sum_{i=1}^{50} C_{p,i} - C_{p,LES} /C_{p,LES}$			$\frac{1}{50} \sum_{i=1}^{50} C_{p,i} - C_{p,LES} /C_{p,Betz}$		
ζ	$C_{p,Nishino}$	$C_{p,model}$	ζ	$C_{p,Nishino}$	$C_{p,model}$
0	2.82%	2.15%	0	0.142%	0.108%
5	4.38%	1.48%	5	0.954%	0.338%
10	5.16%	1.35%	10	1.67%	0.459%
15	5.66%	1.30%	15	2.24%	0.542%
20	6.02%	1.26%	20	2.72%	0.601%
25	6.30%	1.24%	25	3.11%	0.648%

with ζ . The mean absolute error of $C_{p,model}$ was typically less than 1.5% and this decreased slightly as ζ increases (due to the reference power $C_{p,LES}$ increasing). We also use the power of an isolated ideal turbine, $C_{p,Betz}$, as a reference power. $C_{p,Betz}$ is calculated using the actuator disc theory with the expression $C_{p,Betz} = 64C'_T/(4 + C'_T)^3$ (note that in this study $C'_T = 1.33$ and hence $C_{p,Betz} = 0.563$). In this case the mean absolute error increased with ζ for both $C_{p,model}$ and $C_{p,Nishino}$. However, the average prediction error of $C_{p,model}$ remained below 0.65%.

6 | DISCUSSION

Data-driven modelling of the internal turbine thrust coefficient C_T^* is a novel approach to modelling turbine-wake interactions. Data-driven models of wind farm performance typically focus on predicting the power output, which, however, depends on flow physics across a wide range of scales. Current data-driven approaches are either not generalisable to different atmospheric responses, or would require a very large set of expensive training data, such as finite-size wind farm LES data. Data-driven models of C_T^* captures the effects of turbine-wake interactions, whilst also being applicable to different atmospheric responses (following the concept of the two-scale momentum theory).

The statistical emulator of C_T^* developed in this study was able to predict the farm power C_p of Kirby et. al.⁸ with an average error of less than 0.65%. The high accuracy and very low computational cost of this approach shows the potential of this approach for modelling turbine-wake interactions. It has several advantages over traditional approaches using the superposition of wake models. Information from turbulence-resolving LES is included which ensures a high accuracy. It will also be more advantageous as wind farms become larger because wake models struggle to capture the complex multi-scale flows physics which are important for large farms. The statistical model of C_T^* may therefore allow fast and accurate predictions of wind farm performance.

All emulators developed in this study gave substantially better predictions of $C_{T,LES}^*$ compared to the analytical and wake models. Both the mean and maximum prediction errors were reduced by the emulators. The standard GP regression approach had a mean prediction error of 1.26% and maximum error of approximately 6%. The accuracy depends on the size of the LES data set and could be further decreased with a larger training set. The multi-fidelity GP approach gave more accurate predictions of $C_{T,LES}^*$ compared to the standard GP regression. This is because non-linear information fusion algorithm has incorporated information from many low-fidelity data points to improve the emulator of the high-fidelity (LES) model. This approach has the advantage that, unlike the standard GP regression approach, it is not necessary to evaluate the prior mean before making a prediction. Therefore, to predict C_T^* it is only necessary to evaluate the posterior mean of the high-fidelity emulator for a specific turbine layout.

The shape of the posterior mean in figure 8 gives insights into the physics of turbine-wake interactions. This is because $C_{T,LES}^*$ is low when a layout has a high degree of turbine-wake interactions. For the turbine operating conditions used, $C_{T,LES}^*$ is close to 0.75 when a layout has a small degree of wake interactions. Figure 8a shows $C_{T,LES}^*$ when the wind direction is perfectly aligned

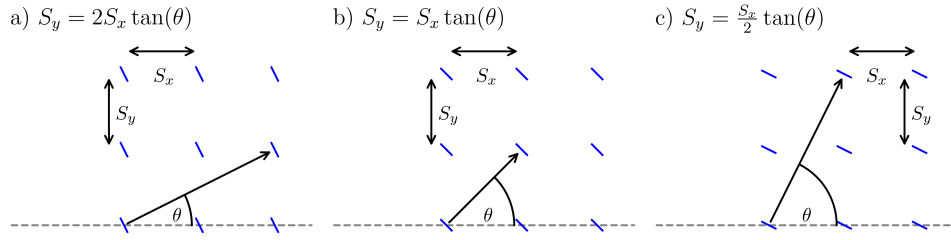


Figure 13 Alignment of turbines for different combinations of S_x , S_y and θ .

with the rows of turbines ($\theta = 0$). This gives wind farms with a high degree of wake interactions which results in low $C_{T,LES}^*$ values. For $\theta = 0^\circ$, increasing S_x/D increases C_T^* because there is a larger streamwise distance between turbines for the wakes to recover. When the cross-streamwise spacing (S_y/D) is increased the degree of wake interactions increases, i.e., $C_{T,LES}^*$ decreases. This is because there is a lower array density which results in a lower turbulence intensity within the farm and hence slower wake recovery. Yang⁴¹ found that increasing the cross-streamwise spacing in infinitely-large wind farms increased the power of individual turbines and concluded that this was due to reduced wake interactions. However, the increase in turbine power found by Yang⁴¹ may be also explained by to a faster farm-averaged wind speed caused by a reduced array density rather than reduced wake interactions.

When the wind direction θ increases, $C_{T,LES}^*$ increases to a maximum of just over 0.75 at $\theta = 10^\circ$ (figure 8c). This result agrees qualitatively with another study⁴² in which it was found that the maximum farm power was produced by an intermediate wind direction. When θ increases above 20° regions of low $C_{T,LES}^*$ appear diagonally (see figures 8f-j). The regions of low $C_{T,LES}^*$ are centred on the surfaces given by $S_y = 2S_x \tan(\theta)$, $S_y = S_x \tan(\theta)$ and $S_y = 0.5S_x \tan(\theta)$. These regions correspond to turbines being aligned along different axes throughout the farm (see figure 13). There are longer streamwise distance between turbines for these arrangements (compared to $\theta = 0^\circ$) and so the $C_{T,LES}^*$ values are higher than for $\theta = 0^\circ$.

The accuracy of the statistical emulators could be further improved in future studies. Both the standard and multi-fidelity GP models can be improved by adding more evaluations of $C_{T,LES}^*$. From table 2, the accuracy of the multi-fidelity GP models did not improve once we used more than 500 $C_{T,wake}^*$ evaluations. This shows that the error in predicting $C_{T,LES}^*$ for MF-GP-nlow500 is not due to the model of f_{wake} . Instead the error arises from the learnt relationship between f_{wake} and f_{LES} .

The statistical emulators developed are not applicable to all wind farms because of the limited nature of our data set. A limitation of the developed model is that it is only applicable to farms with perfectly aligned layouts. It should also be noted that our model was trained on data from simulations of a neutrally stratified boundary layer. Therefore a larger LES data set with an extended parameter space would be required to account for the effect of atmospheric stability on wake interactions and the resulting C_T^* . Another limitation of our model is that it assumes all turbines have the same resistance coefficient C_T' . It is likely that this condition can be strictly satisfied only in the fully developed region of a large farm where the wind speed does not change in the streamwise or cross-streamwise directions.

Although we considered only actuator discs in this study for demonstration, the proposed approach using a data-driven model of C_T^* can be applied to power prediction of real turbines as well in future studies. In this study, we calculate $C_{p,model}$ using the expression $C_{p,model} = \beta^3 C_T^{*\frac{3}{2}} C_T'^{-\frac{1}{2}}$. This assumes that the relationship between C_p^* and C_T' is given by $C_p^* = C_T^{*\frac{3}{2}} C_T'^{-\frac{1}{2}}$, which is only valid for actuator discs. For real turbines, the relationship between C_p^* and C_T' can be calculated using BEM theory⁴³ according to the turbine design and operating conditions (noting that the turbine induction factor can still be estimated as $a = C_T'/(4 + C_T')$). $C_{p,model}$ can then be calculated using equation 5 with β found using equation 1. However, for a data-driven model of C_T^* to be applicable to real turbines, it will be necessary to model the impact of a variable C_T' rather than assuming a fixed C_T' value as in this study.

7 | CONCLUSIONS

In this study we proposed a new data-driven approach to modelling turbine wake interactions and resulting flow resistance in large wind farms. We developed statistical emulators of the farm-internal turbine thrust coefficient $C_{T,LES}^*$ as a function of turbine layout and wind direction. C_T^* represents the flow resistance within a wind farm and reflects the characteristics of the turbine-scale flows including wake and turbine blockage effects. We developed several emulators using both standard GP regression and multi-fidelity GP regression. The standard GP was trained using data from 50 infinitely-large wind farm LES (and using a low-fidelity wake model as a prior mean). The multi-fidelity GP was trained using data from both LES and wake model simulations. We estimated the test accuracy of the model by performing leave-one-out cross-validation and assessed the error in predicting $C_{T,LES}^*$. All emulators had a mean test error of less than 2% for predicting $C_{T,LES}^*$. The multi-fidelity GP gave the best performance with a mean prediction error of 0.849% and maximum prediction error of 3.78% with no bias for under or over-prediction. This is low compared to the mean error of the wake model (4.60%) and analytical C_T^* model (5.26%) which both had a bias for overpredicting $C_{T,LES}^*$.

We used an emulator of $C_{T,LES}^*$ to make predictions of wind farm performance under various mesoscale atmospheric conditions (characterised by the wind extractability factor ζ) using the two-scale momentum theory²⁴. Our predictions of farm power production had an average error of less than 1.5% under realistic wind extractability scenarios compared to the LES. When the error in power prediction is expressed relative to the power of an isolated ideal turbine the average prediction error is less than 0.7%. We also used a previously proposed analytical model of C_T^* ²⁵ to predict farm power output with an average error of less than 3.5% (with the power of an isolated turbine as the reference power). The analytical model correctly predicts the trends in farm performance with array density under different scenarios of large-scale atmospheric response, although it tends to overpredict the power where turbine-wake interactions are important. Using statistical emulators of C_T^* is a new approach to modelling turbine-wake interactions and flow resistance within large wind farms. The approach can be extended in future studies by increasing the size of the training data set, for example, to account for the effects of C_T' and atmospheric stability conditions on C_T^* . The very low computational cost and high accuracy of the model could be beneficial for future wind farm optimisation.

ACKNOWLEDGMENTS

The first author (AK) acknowledges the NERC-Oxford Doctoral Training Partnership in Environmental Research (NE/S007474/1) for funding and training.

Author contributions

T.N. derived the theory. A.K. and T.D.D. performed the simulations. F-X.B. provided assistance and guidance for the machine learning methodology. A.K. wrote the paper with corrections from T.N., F-X.B and T.D.D.

Financial disclosure

None reported.

Conflict of interest

The authors report no conflict of interest.

Data availability statement

The data and code that support the findings of this study are openly available at https://github.com/AndrewKirby2/ctstar_statistical_model. This includes the results from the wind farm LES and wake model simulations. The repository also includes the code for the results presented in sections 5.1, 5.2 and 5.3.

Author ORCID

A. Kirby, <https://orcid.org/0000-0001-8389-1619>; F-X. Briol <https://orcid.org/0000-0002-0181-2559>; T. Nishino, <https://orcid.org/0000-0001-6306-7702>.

References

1. Porté-Agel F, Bastankhah M, Shamsoddin S. Wind-Turbine and Wind-Farm Flows: A Review. *Boundary-Layer Meteorology* 2020; 174: 1-59. doi: 10.1007/s10546-019-00473-0
2. Bleeg J, Purcell M, Ruisi R, Traiger E. Wind farm blockage and the consequences of neglecting its impact on energy production. *Energies* 2018; 11: 1609. doi: 10.3390/en11061609
3. Carbon Trust . Global Blockage Effect in Offshore Wind (GloBE) [accessed 07/11/2022]. <https://www.carbontrust.com/our-projects/large-scale-rd-projects-offshore-wind/global-blockage-effect-in-offshore-wind-globe>; 2022.
4. Jensen NO. A note on wind generator interaction. *Risø-M-2411 Risø National Laboratory Roskilde* 1983.
5. Bastankhah M, Porté-Agel F. A new analytical model for wind-turbine wakes. *Renewable Energy* 2014; 70: 116-123. doi: 10.1016/j.renene.2014.01.002
6. Katic I, Hojstrup J, Jensen NO. A simple model for cluster efficiency. *Proceedings of the European wind energy association conference and exhibition, Rome, Italy* 1986: 407-409.
7. Zong H, Porté-Agel F. A momentum-conserving wake superposition method for wind farm power prediction. *Journal of Fluid Mechanics* 2020; 889: A8. doi: 10.1017/jfm.2020.77
8. Kirby A, Nishino T, Dunstan TD. Two-scale interaction of wake and blockage effects in large wind farms. *Journal of Fluid Mechanics* 2022; 953: A39. doi: 10.1017/jfm.2022.979
9. Stevens RJAM, Gayme DF, Meneveau C. Effects of turbine spacing on the power output of extended wind-farms. *Wind Energy* 2016; 19: 359-370. doi: 10.1002/we.1835
10. Fitch AC, Olson JB, Lundquist JK, et al. Local and mesoscale impacts of wind farms as parameterized in a mesoscale NWP model. *Monthly Weather Review* 2012; 140. doi: 10.1175/MWR-D-11-00352.1
11. Abkar M, Porté-Agel F. A new wind-farm parameterization for large-scale atmospheric models. *Journal of Renewable and Sustainable Energy* 2015; 7. doi: 10.1063/1.4907600
12. Pan Y, Archer CL. A Hybrid Wind-Farm Parametrization for Mesoscale and Climate Models. *Boundary-Layer Meteorology* 2018; 168: 469-495. doi: 10.1007/s10546-018-0351-9

13. Zehtabiyani-Rezaie N, Iosifidis A, Abkar M. Data-driven fluid mechanics of wind farms: A review. *Journal of Renewable and Sustainable Energy* 2022; 14: 32703. doi: 10.1063/5.0091980
14. Renganathan SA, Maulik R, Letizia S, Iungo GV. Data-driven wind turbine wake modeling via probabilistic machine learning. *Neural Computing and Applications* 2022; 34: 6171-6186. doi: 10.1007/s00521-021-06799-6
15. Optis M, Perr-Sauer J. The importance of atmospheric turbulence and stability in machine-learning models of wind farm power production. *Renewable and Sustainable Energy Reviews* 2019; 112: 27-41. doi: 10.1016/j.rser.2019.05.031
16. Japar F, Mathew S, Narayanaswamy B, Lim CM, Hazra J. Estimating the wake losses in large wind farms: A machine learning approach. *ISGT 2014* 2014: 1-5. doi: 10.1109/ISGT.2014.6816427
17. Yan C, Pan Y, Archer CL. A general method to estimate wind farm power using artificial neural networks. *Wind Energy* 2019; 22: 1421-1432. doi: 10.1002/we.2379
18. Zhang J, Zhao X. Wind farm wake modeling based on deep convolutional conditional generative adversarial network. *Energy* 2022; 238: 121747. doi: <https://doi.org/10.1016/j.energy.2021.121747>
19. Wilson B, Wakes S, Mayo M. Surrogate modeling a computational fluid dynamics-based wind turbine wake simulation using machine learning. *2017 IEEE Symposium Series on Computational Intelligence (SSCI) 2017*: 1-8. doi: 10.1109/SSCI.2017.8280844
20. Ti Z, Deng XW, Yang H. Wake modeling of wind turbines using machine learning. *Applied Energy* 2020; 257: 114025. doi: <https://doi.org/10.1016/j.apenergy.2019.114025>
21. Ti Z, Deng XW, Zhang M. Artificial Neural Networks based wake model for power prediction of wind farm. *Renewable energy* 2021; 172: 618-631. doi: <https://doi.org/10.1016/j.renene.2021.03.030>
22. Park J, Park J. Physics-induced graph neural network: An application to wind-farm power estimation. *Energy* 2019; 187. doi: 10.1016/j.energy.2019.115883
23. Bleeg J. A Graph Neural Network Surrogate Model for the Prediction of Turbine Interaction Loss. *Journal of Physics: Conference Series* 2020; 1618. doi: 10.1088/1742-6596/1618/6/062054
24. Nishino T, Dunstan TD. Two-scale momentum theory for time-dependent modelling of large wind farms. *Journal of Fluid Mechanics* 2020; 894: A2. doi: 10.1017/jfm.2020.252
25. Nishino T. Two-scale momentum theory for very large wind farms. *Journal of Physics: Conference Series* 2016; 753: 032054. doi: 10.1088/1742-6596/753/3/032054
26. Patel K, Dunstan TD, Nishino T. Time-dependent upper limits to the performance of large wind farms due to mesoscale atmospheric response. *Energies* 2021; 14: 6437. doi: 10.3390/en14196437
27. Sacks J, Welch WJ, Mitchell TJ, Wynn HP. Design and analysis of computer experiments. *Statistical Science* 1989; 4: 409-423. doi: 10.1214/ss/1177012413
28. Currin C, Mitchell T, Morris M, Ylvisaker D. Bayesian prediction of deterministic functions, with applications to the design and analysis of computer experiments. *Journal of the American Statistical Association* 1991; 86: 953-963. doi: 10.1080/01621459.1991.10475138

29. Johnson ME, Moore LM, Ylvisaker D. Minimax and maximin distance designs. *Journal of Statistical Planning and Inference* 1990; 26: 131-148. doi: 10.1016/0378-3758(90)90122-B
30. Santner TJ, Williams BJ, Notz W. *The design and analysis of computer experiments*. second ed. 2018.
31. Wynne G, Briol FX, Girolami M. Convergence guarantees for gaussian process means with misspecified likelihoods and smoothness. *Journal of Machine Learning Research* 2021; 22.
32. Shapiro CR, Gayme DF, Meneveau C. Filtered actuator disks: Theory and application to wind turbine models in large eddy simulation. *Wind Energy* 2019; 22: 1414-1420. doi: 10.1002/we.2376
33. Niayifar A, Porté-Agel F. Analytical modeling of wind farms: A new approach for power prediction. *Energies* 2016; 9. doi: 10.3390/en9090741
34. Pedersen MM, Laan v. dP, Friis-Møller M, Rinker J, Réthoré PE. DTUWindEnergy/PyWake: PyWake. 2021. doi: 10.5281/zenodo.2562662
35. Crespo A, Hernández J. Turbulence characteristics in wind-turbine wakes. *Journal of Wind Engineering and Industrial Aerodynamics* 1996; 61: 71-85. doi: 10.1016/0167-6105(95)00033-X
36. Rasmussen CE, Williams CKI. *Gaussian Processes for Machine Learning*. the MIT Press . 2018
37. Peherstorfer B, Willcox K, Gunzburger M. Survey of multifidelity methods in uncertainty propagation, inference, and optimization. *SIAM Review* 2018; 60. doi: 10.1137/16M1082469
38. Perdikaris P, Raissi M, Damianou A, Lawrence ND, Karniadakis GE. Nonlinear information fusion algorithms for data-efficient multi-fidelity modelling. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 2017; 473. doi: 10.1098/rspa.2016.0751
39. Paleyes A, Pullin M, Mahsereci M, Lawrence N, González J. Emulation of physical processes with Emukit. 2019.
40. GPy . GPy: A Gaussian process framework in python. <http://github.com/SheffieldML/GPy>; since 2012.
41. Yang X, Kang S, Sotiropoulos F. Computational study and modeling of turbine spacing effects infinite aligned wind farms. *Physics of Fluids* 2012; 24: 11510. doi: 10.1063/1.4767727
42. Stevens RJAM, Gayme DF, Meneveau C. Large eddy simulation studies of the effects of alignment and wind farm length. *Journal of Renewable and Sustainable Energy* 2014; 6: 023105. doi: 10.1063/1.4869568
43. Nishino T, Hunter W. Tuning turbine rotor design for very large wind farms. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 2018; 474(2220): 1–20. doi: 10.1098/rspa.2018.0237

How to cite this article: Kirby A., Briol F-X., Dunstan T.D., and Nishino T. (2022), Data-driven modelling of wind turbine wake interactions in large wind farms, *Wind Energy*, xxxx.