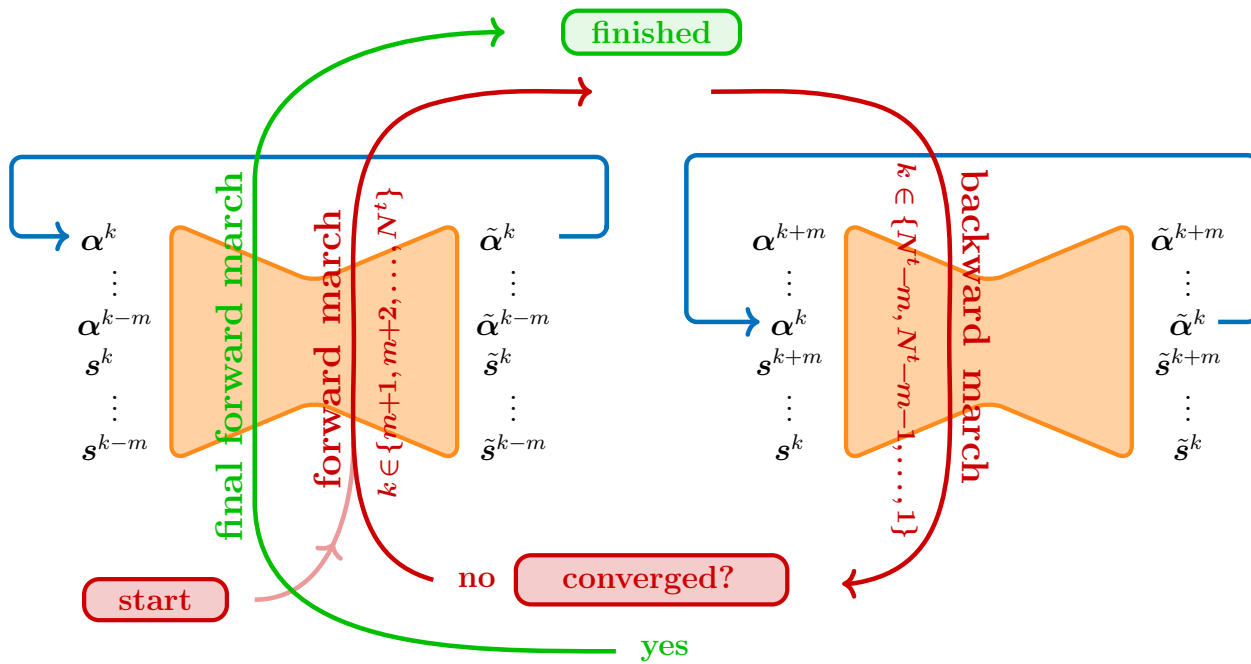


# Graphical Abstract

## Data Assimilation with Machine Learning for Dynamical Systems: Modelling Indoor Ventilation

Claire E. Heaney, Jieyi Tang, Jintao Yan, Donghu Guo, Jamesson Ipock, Sanjana Kaluvakollu, Yushen Lin, Danhui Shao, Boyang Chen, Laetitia Mottet, Prashant Kumar, Christopher C. Pain



## Highlights

### **Data Assimilation with Machine Learning for Dynamical Systems: Modelling Indoor Ventilation**

Claire E. Heaney, Jieyi Tang, Jintao Yan, Donghu Guo, Jamesson Ipock, Sanjana Kaluvakollu, Yushen Lin, Danhui Shao, Boyang Chen, Laetitia Mottet, Prashant Kumar, Christopher C. Pain

- Presentation of a method which combines data assimilation (DA) with machine learning for efficiency.
- The approach is an alternative to 4D variational DA, and uses both solution variables and observations as the input to a neural network.
- A demonstration of the data assimilation method for a dual-twin experiment.
- The assimilation of observations measured in a school classroom in north London.

# Data Assimilation with Machine Learning for Dynamical Systems: Modelling Indoor Ventilation

Claire E. Heaney<sup>a,b,\*</sup>, Jieyi Tang<sup>c,1</sup>, Jintao Yan<sup>c,1</sup>, Donghu Guo<sup>c,2</sup>, Jamesson Ipock<sup>c,2</sup>, Sanjana Kaluvakollu<sup>c,2</sup>, Yushen Lin<sup>c,d,2</sup>, Danhui Shao<sup>c,2</sup>, Boyang Chen<sup>a</sup>, Laetitia Mottet<sup>a,e</sup>, Prashant Kumar<sup>f,g</sup>, Christopher C. Pain<sup>a,b,h</sup>

<sup>a</sup>*Applied Modelling and Computation Group, Department of Earth Science and Engineering, Imperial College London, South Kensington Campus, London, SW7 2AZ, UK*

<sup>b</sup>*Centre for AI-Physics Modelling, Imperial-X, Imperial College London, White City Campus, London, W12 7SL, UK*

<sup>c</sup>*Department of Earth Science and Engineering, Imperial College London, South Kensington Campus, London, SW7 2AZ, UK*

<sup>d</sup>*Department of Electrical and Electronic Engineering, The University of Manchester, Oxford Road, Manchester, M13 9PL, UK*

<sup>e</sup>*Inria, National Institute for Research in Computer Science and Control, Bordeaux, France*

<sup>f</sup>*Global Centre for Clean Air Research (GCARE), School of Sustainability, Civil and Environmental Engineering, University of Surrey, Guildford, GU2 7XH, Surrey, UK*

<sup>g</sup>*Institute for Sustainability, University of Surrey, Guildford, GU2 7XH, Surrey, UK*

<sup>h</sup>*Data Assimilation Laboratory, Data Science Institute, Imperial College London, South Kensington Campus, London, SW7 2AZ, UK*

---

## Abstract

Data assimilation is a method of combining physical observations with prior knowledge (for instance, a computational simulation) in order to produce an improved model; that is, improved over what the physical observations or the computational simulation could offer in isolation. Recently, machine learning techniques have been deployed in order to address the significant computational burden that is associated with the procedures involved in data assimilation.

In this paper we propose an approach that uses a non-intrusive reduced-order model (NIROM) as a surrogate for a high-resolution model thereby saving computational effort. The mismatch between observations and the surrogate model is propagated forwards and backwards in time in a manner similar to 4D-variational data assimilation methods. The observations and prior are reconciled in a new way which takes full advantage of the neural network used in the NIROM and also means that there is no need to form the sensitivities explicitly when propagating the mismatch. Instead, the observations are part of the input and output of the network.

Modelling the air quality in a school classroom is the test case for our demonstration. Firstly, the data assimilation approach is shown to perform very well in a dual-twin type experiment, and secondly, the approach is used to assimilate observations collected from a classroom in Houndsfield Primary School with predictions from the NIROM.

---

\*Corresponding author

*Email address:* [c.heaney@imperial.ac.uk](mailto:c.heaney@imperial.ac.uk) (Claire E. Heaney)

<sup>1</sup>Indicates an equal contribution

<sup>2</sup>Indicates an equal contribution

## 1. Introduction

### 1.1. Background

Data assimilation (DA) aims to combine observations from sensors with predictions from a prior model (often from computer simulations) in order to obtain a more accurate prediction of the state of a system. In recent decades, DA has made a significant impact in earth science applications, for instance, its use has led to a notable increase in the accuracy of medium range weather forecasts [1]. There are two types of data assimilation in common useage: ensemble methods and variational methods [2, 3], as well as many hybrid methods [4, 5]. Ensemble approaches use statistical methods to form the prior model, whereas variational methods minimise a functional which calculates the difference between the prior and the observations [3, 4, 6]. Variational DA can be split into 3D Variational (3D-Var) DA and 4D Variational (4D-Var) DA. The former disregards the time dependency of the observations, whereas the latter takes this into account and is regarded as state-of-the-art DA for operational numerical weather prediction [6–8].

One issue for practitioners of DA is its high computational cost. For ensemble methods, this arises from the need to perform many forward simulations. For variational methods, it arises from the need to march forwards and backwards in time, until the observations and the computer predictions are reconciled. Machine learning (ML) methods are a natural tool with which one can alleviate some of the computational burden of DA, and recently there has been considerable interest in combining ML with DA. In this paper, we wish to exploit ML techniques, through a reduced-order model framework, in order to assimilate data efficiently. We propose an approach which assimilates data through time, propagating the mismatch both forwards and backwards in time in a manner similar to 4D-Var methods. We exploit the architecture and design of the ML model to develop a DA algorithm which does not require differentiation of the model (although this is not generally an issue for ML models). Instead, the observations are part of the input and output of the network.

### 1.2. Related work

Similarities between variational data assimilation and machine learning techniques have been elucidated by a number of authors [8–11]. One such example of this is that both techniques minimise a functional: in DA, the functional describes the mismatch between the observations and the computer prediction, and in ML, the functional that is minimised during training (for supervised or self-supervised learning paradigms) describes the difference between the learned output of a network and the desired output. Both DA and ML apply gradient descent methods to minimise these functionals. A second similarity shared by DA and ML is the concept of the adjoint model, used in DA to march backwards in time, and also used in ML where it is described as the backpropagation algorithm. For a general discussion of these points, see [8–11].

A number of recent articles demonstrate the promise of combining Machine Learning (ML) and DA. In some work, forward and adjoint models are replaced by fast-running ML-based surrogates, which can



36 capitalise on the ability to differentiate the adjoint using the backpropagation functionality commonly  
37 found within ML libraries [6, 7, 12, 13]. Other work attempts to represent the model error rather than the  
38 model itself with ML surrogates [14–18]. Considering ensemble DA approaches, several examples exist of  
39 using ML to facilitate DA [19–25]: both Proper Orthogonal Decomposition (POD) [19] and autoencoder  
40 networks [20, 21] have been used to identify a low dimensional space; Multilayer Perceptrons [22, 23] and  
41 Long Short-Term Memory networks have been used to make predictions in time [19, 20, 24]; and residual  
42 blocks have been used to predict measurements from geophysical logs [25]. Test cases include Lorenz  
43 systems [21, 24]; modelling the surface temperature of the sea [19]; modelling the movement of CO<sub>2</sub> in an  
44 indoor environment [20]; global weather predictions [22] and inversion of electromagnetic measurements  
45 from boreholes [25].

46 Turning to some examples of variational methods which use ML (the focus of this paper), Chennault  
47 et al. [13] used a surrogate model (based on a multi-layer perceptron) to replace the forward and adjoint  
48 models, and incorporated this in a 4D-Var DA method applied to a Lorenz problem. They found that  
49 using information from the adjoint resulted in an improved surrogate model, especially when the amount  
50 of training data was limited. Farchi et al. [17] constructed a 4D-Var method that attempts to learn  
51 the error in a Computational Fluid Dynamics (CFD) model from noisy and sparse observations applied  
52 to a two-layer 2D quasi-geostrophic model. The DA and ML models are used alternately, producing  
53 significantly better results than using the original model. Hatfield et al. [7] used a surrogate model based  
54 on a multi-layer perceptron to replace a discretised system modelling gravity wave drag. Applied to  
55 non-orographic gravity wave drag in an atmospheric model, no statistically significant differences were  
56 found between the ML-based results and the HFM results. For low dimensional problems such as some  
57 Lorenz problems, applying a dimensionality reduction method such as POD is not necessary [24] and was  
58 not done by Hatfield et al. [7], Chennault et al. [13], Farchi et al. [17]. For high-dimensional problems,  
59 incorporating POD into the surrogate model can save a significant amount of computational effort [26–  
60 28]. Gong et al. [29] implemented a 3D-Var DA method in latent space using a combination of POD  
61 and an autoencoder to find the low dimensional space, and a nearest neighbour method to model the  
62 parameter dependence. The original high-fidelity model (HFM) had 5000 grid points, represented by  
63 30 latent variables, and over 24,000 combinations of material parameters were used. Maulik et al. [6]  
64 used POD and an LSTM to construct a surrogate model with which to predict the geopotential height  
65 of a 500 hPa pressure surface. With just over 10,000 spatial degrees of freedom represented by 5 POD  
66 modes, observations at 5000 spatial locations are assimilated. The use of surrogate models for forward  
67 and adjoint calculations results in a reduction in computational effort over classical methods by 4 orders  
68 of magnitude. Silva et al. [30] used a Generative Adversarial Network (GAN) to predict the spread of a  
69 virus through an idealised town with a compartmental model modified to include spatial variation. They  
70 were able to exploit the adjoint-like nature of the GAN to assimilate data, so that given the number  
71 of infections, they could estimate two  $R_0$  numbers, one for each of two classes of person, for example.  
72 Convergence was achieved relatively quickly, as the method was able to resolve the mismatch and update  
73 the model parameters accurately within a few forwards and backwards iterations. Regazzoni et al. [31]  
74 tackle a multiscale problem with a slightly different approach to the previously mentioned work. They use  
75 DA to identify model parameters that led to the results at the faster scale. A machine learning method  
76 is used to model the behaviour over the slower timescale.

77 For many, a method of choice for learning time-dependent behaviour is an LSTM network. However, the  
78 predictions of these networks have sometimes been shown to be unphysical [32–34], either diverging or  
79 reaching a steady state after a certain time. Within our surrogate model, along with POD, we propose  
80 the use of an adversarial autoencoder (AAE) [35] to model the time-dependent behaviour, as it includes  
81 an adversarial layer which attempts to keep predictions close to those seen in the training data. Good  
82 results for prediction have been seen in [36, 37] using a similar network, although, in those papers, the  
83 network architecture was slightly modified so that the output had a different dimension to the input (i.e.,  
84 was an encoder-decoder rather than an autoencoder).

### 85 *1.3. Contribution*

86 A new method for data assimilation is proposed in which data is assimilated whilst marching forwards and  
87 backwards in time, in a manner similar to 4D-Var DA. To ease the computational burden, a surrogate  
88 model is used, which combines POD (for dimensionality reduction) and an AAE (for predicting the  
89 evolution in time of the system). Instead of calculating the gradient of the mismatch between the prior  
90 and the observations, the AAE takes as an input both the observations and the POD coefficients (which  
91 represent the evolution of the system) and is able to provide predictions for the POD coefficients which are  
92 consistent with the observations. Whilst performing time marching, a mismatch functional is evaluated,  
93 and, if this increases, relaxation is applied. Similar to Amendola et al. [20], Peyron et al. [21], Gong et al.  
94 [29], for example, we assimilate data in the low-dimensional space identified by a dimensionality reduction  
95 method. In those articles, an autoencoder is used to find the reduced space, whereas in this paper, POD  
96 is used.

97 The test case used in this paper investigates the air flows and air quality within a naturally ventilated  
98 classroom in Houndsfield Primary School, north London. During winter, a measurement campaign was  
99 undertaken, which collected observations of CO<sub>2</sub>, relative humidity and temperature from 18 locations [38].  
100 Measurements of pollution are also available from this dataset, but are not used in this study. An air  
101 purifier is present in the classroom, which has the dual purpose of improving air quality, in particular by  
102 removing pollution particles, and also reducing the number of particles that may contain the SARS-CoV-2  
103 virus [38]. In this study we also compare CO<sub>2</sub> with a viral load tracer since the former is often used as  
104 a proxy for the virus concentration [39]. Subsequently, a CFD solver was used to obtain predictions of  
105 the indoor flows and air quality within the classroom. A ML-based surrogate model is constructed from  
106 solutions of the CFD model (snapshots) and two methods of predicting in time with the surrogate are  
107 investigated. A method for DA is then developed using the surrogate and a set of observations (generated  
108 from the CFD model). A dual-twin experiment is carried out to test how well the approach can assimilate  
109 data. Finally, observations collected from the classroom over the period of an hour are assimilated into  
110 the ML surrogate model.

111 The main contributions of this paper are (i) to present a method which combines data assimilation  
112 (4D-Var) with machine learning; (ii) to compare two methods for making predictions in time with the  
113 adversarial autoencoder; (iii) to demonstrate the data assimilation approach for a dual-twin experiment;  
114 (iv) to assimilate observations taken from a classroom in Houndsfield Primary School, north London. The  
115 novelty of this work is twofold. It comes from having the sensor values as direct inputs and outputs of  
116 the AAE, and also using the AAE for prediction and DA. As a consequence of this, no differentiation of

117 the surrogate adjoint model is needed, even though the proposed method has similarities with the 4D-Var  
 118 method.

119 The remainder of this article is arranged as follows. The methodology is presented in Section 2 and results  
 120 are shown in Section 3. Section 4 draws conclusions and discusses possible directions for future work.

## 121 2. Methodology

122 A non-intrusive reduced-order model is used to reduce the computational cost of the forward model,  
 123 thereby making data assimilation tractable. The two main elements of this method are explained in  
 124 Section 2.1 (finding a low-dimensional space which can represent well the CFD solutions) and Section 2.2  
 125 (building a model which can predict the behaviour of the system in the low-dimensional space). The  
 126 proposed method for data assimilation is presented in Section 2.3.

### 127 2.1. Dimensionality Reduction

128 Proper Orthogonal Decomposition (POD) was used to reduce the dimension of the high-fidelity results  
 129 generated by the CFD code that solved the discretisation of the Navier-Stokes Equations. For the test  
 130 case used in this paper, the seven fields of interest are velocity (one field for each of the three components),  
 131 relative humidity, temperature, carbon dioxide levels and a field which represents the potential viral load  
 132 in the air. The latter was based on carbon dioxide levels but included a half-life to take account of how  
 133 long a virus might remain active in the air. We chose the value of 35 minutes for this, which is within the  
 134 possible range of values for SARS-CoV-2 [40]. Each field was normalised or scaled so that the values at  
 135 each node lie within the range  $[-1,1]$  in order to take account of the different units of the seven solution  
 136 fields and the different physical process being represented by each field. Therefore, the snapshots matrix  
 137 has dimension  $7N$  by  $M$ , where  $N$  is the number of nodes, the number of fields is 7, and  $M$  is the number  
 138 of snapshots used in building the surrogate model. Using Singular Value Decomposition, the snapshots  
 139 matrix  $\Phi$  is decomposed in the following manner

$$140 \quad [\phi^1 \ \phi^2 \ \dots \ \phi^M] =: \Phi = \mathbf{U}\Sigma\mathbf{V}^T \quad (1)$$

141 where  $\phi^k \in \mathbb{R}^{7N}$  is a column vector containing the seven normalised solution fields at time level  $k$ ,  $\mathbf{U}$  and  
 142  $\mathbf{V}$  are orthogonal matrices containing the left and right singular vectors respectively, and the singular  
 143 values  $\{\sigma_i\}_{i=1}^M$  are found on the diagonal of the matrix  $\Sigma$ . An empirical measure of how much information  
 144 is captured by using  $N^{\text{POD}}$  of the possible  $M$  basis functions is given by

$$145 \quad \frac{\sum_{i=1}^{N^{\text{POD}}} \sigma_i^2}{\sum_{i=1}^M \sigma_i^2}. \quad (2)$$

146 Once  $N^{\text{POD}}$  has been chosen, the first  $N^{\text{POD}}$  columns of  $\mathbf{U}$  are taken as the basis functions, which  
 147 are stored in the matrix  $\mathbf{R} \in \mathbb{R}^{7N \times N^{\text{POD}}}$ . Choosing the basis functions to be the left singular vectors  
 148 minimises the least square error of the projection error of the snapshots. Of the total number of snapshots,  
 149 denoted by  $N^s$ ,  $M$  snapshots are used to find the low-dimensional space and are also used to determine  
 150 the mappings related to the normalisation of data. The remaining  $N^s - M$  snapshots are used as unseen

151 or test data. All the snapshots are projected onto the low-dimensional space by the basis functions as  
 152 follows

$$153 \quad \boldsymbol{\alpha}^k = \mathbf{R}^T \boldsymbol{\phi}^k \quad \forall k, \quad (3)$$

154 where  $\boldsymbol{\alpha}^k$  is a vector containing the POD coefficients associated with time level  $k$ . Once the POD  
 155 coefficients for the snapshots in the training dataset are obtained, each POD coefficient is scaled or  
 156 normalised so that its values over time lie in the range  $[0,1]$ . This process is known as normalisation in  
 157 machine learning terminology and is considered good practice when training neural networks. For more  
 158 details on the theory of POD and the optimality of the basis functions see [41].

## 159 2.2. Prediction

160 An adversarial autoencoder [35] is chosen to learn the behaviour of the system through time. An autoen-  
 161 coder is a neural network which attempts to learn the identity map through two networks, an encoder  
 162 and a decoder. In addition to this, the adversarial autoencoder contains a discriminator network, the  
 163 purpose of which is to encourage the latent space to follow a (given) prior distribution ( $P_{\text{prior}}$ ). This  
 164 means that the latent space is less likely to have gaps and will perform better for unseen data [35]. See  
 165 Figure 1 for a schematic diagram of an adversarial autoencoder, including the encoder,  $\mathcal{E}^{\text{nc}}$ , the decoder,  
 166  $\mathcal{D}^{\text{ec}}$ , and discriminator,  $\mathcal{D}^{\text{is}}$ . The discriminator is used during training, but is not needed for prediction  
 167 (i.e., inference), for which the map learned by the AAE can be written as:

$$168 \quad \tilde{\boldsymbol{x}} = f^{\text{AAE}}(\boldsymbol{x}) \equiv \mathcal{D}^{\text{ec}}(\mathcal{E}^{\text{nc}}(\boldsymbol{x})), \quad (4)$$

169 where  $\boldsymbol{x}$  and  $\tilde{\boldsymbol{x}}$  represent the input and output of the network respectively. During training, there are  
 170 three steps undertaken for each subset of training data (mini-batch). First, the reconstruction error is  
 171 minimised (the difference between the input of the encoder and the output of the decoder); second, the  
 172 discriminator is trained with samples generated by the encoder (labelled as false) and samples from the  
 173 prior distribution (labelled as true); and finally, the encoder is trained to fool the discriminator. This  
 174 training procedure can be summarised as

$$175 \quad \min_{\mathcal{E}^{\text{nc}}, \mathcal{D}^{\text{ec}}} \mathbb{E}(\|\boldsymbol{x} - \tilde{\boldsymbol{x}}\|_2^2) + \min_{\mathcal{E}^{\text{nc}}} \max_{\mathcal{D}^{\text{is}}} (\mathbb{E}_{\boldsymbol{z} \sim P_{\text{prior}}}[\log \mathcal{D}^{\text{is}}(\boldsymbol{z})] + \mathbb{E}_{\boldsymbol{x} \sim P_{\text{data}}}[\log(1 - \mathcal{D}^{\text{is}}(\mathcal{E}^{\text{nc}}(\boldsymbol{x})))]), \quad (5)$$

176 where  $\mathbb{E}$  is the expectation,  $\boldsymbol{z} \sim P_{\text{prior}}$  is a sample from the desired distribution and  $\boldsymbol{x} \sim P_{\text{data}}$  is  
 177 a sample from the training data (POD coefficients and sensor values). The first term represents the  
 178 reconstruction error (involving the encoder and decoder), and the second term represents the adversarial  
 179 training (involving the encoder and the discriminator). The input of the AAE is made up of a sequence of  
 180  $m + 1$  POD coefficients (representing the CFD model) and sensor values (representing the observations).  
 181 To generate the sensor values used in training, the CFD model is evaluated at the sensor locations. When  
 182 predicting, only the first  $m$  sensor values need to be used along with the POD coefficients at  $m$  time levels,  
 183 and the algorithm will predict the POD coefficients and the sensor values at the subsequent time level.  
 184 When assimilating data, actual observations at the sensor locations can be used, therefore  $m + 1$  sensor  
 185 values and POD coefficients at  $m$  time levels will be used to predict the POD coefficients at the subsequent  
 186 time level.

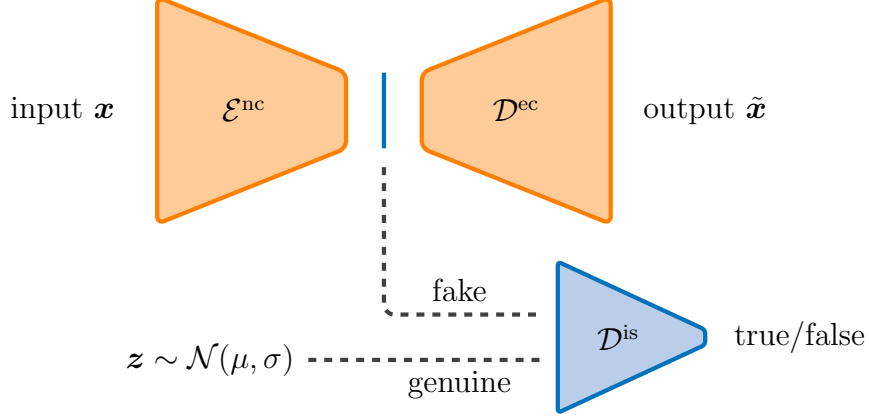


Figure 1: Schematic diagram of an adversarial autoencoder. The encoder is represented by  $\mathcal{E}^{\text{nc}}$ ; the blue line represents the adversarial layer (and is the output of the encoder); the decoder network  $\mathcal{D}^{\text{ec}}$  maps the values in the adversarial layer to the output; the input to the discriminator  $\mathcal{D}^{\text{is}}$  is either a (genuine) sample from the prior distribution (here, a Gaussian distribution  $\mathcal{N}(\mu, \sigma)$  with mean  $\mu$  and variance  $\sigma$ ) or a (fake) sample from the output of the encoder.

187 Autoencoders are typically used for compression. In order to use this type of network to predict in time,  
 188 two strategies were tested, both using a time-marching approach. Both methods can be applied to the  
 189 same trained AAE. When being used for prediction, the map learned by the AAE can be written as

$$190 \quad \tilde{\mathbf{X}} = f^{\text{AAE}}(\mathbf{X}), \text{ where } \mathbf{X} = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{x}^{k-1} \\ \vdots \\ \mathbf{x}^{k-m} \end{pmatrix}, \tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}^k \\ \tilde{\mathbf{x}}^{k-1} \\ \vdots \\ \tilde{\mathbf{x}}^{k-m} \end{pmatrix}, \text{ and } \mathbf{X}, \tilde{\mathbf{X}} \in \mathbb{R}^{(m+1) \times F}. \quad (6)$$

191  $F$  represents the number of features (here, the number of POD coefficients plus the number of sensor  
 192 values observed at each sensor multiplied by the number of sensors),  $\mathbf{x}^k$  represents the POD coefficients  
 193 and sensor values at time level  $k$  and  $\tilde{\mathbf{x}}^k$  represents the corresponding part of the output. Here,  $\mathbf{x}^k$  and  $\tilde{\mathbf{x}}^k$   
 194 are row vectors. The input and output of the AAE are 2D arrays,  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  respectively, as convolutional  
 195 layers are used within the encoder and decoder. Once the AAE has been trained to reproduce these values  
 196 as closely as it can (whilst simultaneously attempting to make the latent space Gaussian), we can make  
 197 predictions in time, using either of the following two methods.

198 The first method (“constrained prediction”) requires  $m$  time levels ( $k-1, k-2, \dots, k-m$ ), from what is  
 199 sometimes referred to as a “look-back window”, in order to make a prediction for the  $k$ th time level. Both  
 200 known and unknown values make up the input of the AAE. As  $\mathbf{x}^k$  is unknown initially, we approximate  
 201 it with  $\mathbf{x}^{k-1}$ . The values in Equation (6) are passed through the (trained) AAE resulting in an updated  
 202 approximation for  $\mathbf{x}^k$ , written as  $\tilde{\mathbf{x}}^k$ . Disregarding the other outputs of the AAE (for the known or  
 203 previous time levels), we iterate (by repeatedly passing the updated values for  $\mathbf{x}^k$  and the unchanged  
 204 known values through the AAE) until the values of  $\mathbf{x}^k$  do not change within a given tolerance. We then  
 205 take values at time levels ( $k, k-1, \dots, k-m+1$ ) to be the known values and, from these, make a  
 206 prediction for time level  $k+1$ , and so on. This procedure is described in Algorithm 1 and Figure 2.

---

**Algorithm 1** Constrained Prediction. By constraining the values of the sensors and POD coefficients at the previous (known) time levels, we can approximate the values of the sensors and POD coefficients at the future time level. The quantity represented by  $\|\cdot\|_2$  is the  $L_2$  norm, which, is related to the mean square error.

---

1: Given the maximum number of iterations ( $N^{its}$ ); the trained AAE ( $f^{AAE}$ ); a tolerance for the stopping criterion ( $\varepsilon$ ); and  $m$  previous solutions (POD coefficients and sensor values),  $\mathbf{x}^{k-1}, \mathbf{x}^{k-2}, \dots, \mathbf{x}^{k-m}$ , this algorithm will make predictions for the POD coefficients and sensor values at future time levels.

2: **for** time level  $k \in$  desired range of time levels **do**

3:     *!! approximate the solution at time level k*

4:      $\mathbf{x}^{k,0} = \mathbf{x}^{k-1}$

5:     **for** iteration  $i = 1, 2, \dots, N^{its}$  **do**

6:

$$\begin{pmatrix} \tilde{\mathbf{x}}^{k,i-1} \\ \tilde{\mathbf{x}}^{k-1} \\ \vdots \\ \tilde{\mathbf{x}}^{k-m} \end{pmatrix} = f^{AAE} \begin{pmatrix} \mathbf{x}^{k,i-1} \\ \mathbf{x}^{k-1} \\ \vdots \\ \mathbf{x}^{k-m} \end{pmatrix} \quad (7)$$

7:      $\mathbf{x}^{k,i} = \tilde{\mathbf{x}}^{k,i-1}$

8:     **if** ( $\|\tilde{\mathbf{x}}^{k,i-1} - \mathbf{x}^{k,i-1}\|_2^2 < \varepsilon$  or  $i == N^{its}$ ) **then**

9:          $\mathbf{x}^k = \mathbf{x}^{k,i}$

10:         exit iteration loop

11:     **end if**

12:     **end for**

13: **end for**

---

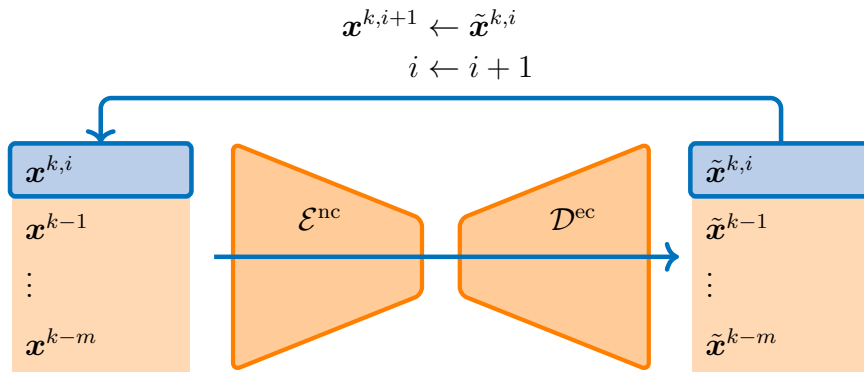


Figure 2: Schematic diagram illustrating the method of constrained prediction. The input includes known values ( $\mathbf{x}^{k-1}, \mathbf{x}^{k-2}, \dots, \mathbf{x}^{k-m}$ ) (shaded in orange), and unknown values  $\mathbf{x}^{k,i}$  (shaded with blue). The approximation of the solution,  $\mathbf{x}^{k,i}$ , is fed into the AAE along with the known values and the output of the AAE,  $\tilde{\mathbf{x}}^{k,i}$ , is taken to be the updated solution. This is fed back into the input of the AAE with the known values and the iteration process continues until the solution does not change to within a given tolerance.

207 The second method (“prediction by backpropagation”) uses the trained decoder to make predictions in

208 a manner similar to that described in Silva et al. [30]. Again, the method requires  $m$  previous time  
209 levels to make a prediction for the the solution at a future time level. A set of randomly chosen latent  
210 variables ( $\tilde{\mathbf{z}}^k$ ) is decoded and the difference between the output corresponding to the known values  
211 ( $\tilde{\mathbf{x}}^{k-1}, \tilde{\mathbf{x}}^{k-2}, \dots, \tilde{\mathbf{x}}^{k-m}$ ) and the known values themselves ( $\mathbf{x}^{k-1}, \mathbf{x}^{k-2}, \dots, \mathbf{x}^{k-m}$ ) is minimised with  
212 respect to the latent variables. This effectively finds the latent variables which result in an output close  
213 to the values at the  $m$  known time levels. Once this has converged, the  $(m + 1)$ th value in the output is  
214 taken as the prediction of the POD coefficients and sensor values at the future time level, time level  $k$ .  
215 The process is repeated for successive time levels as described in Algorithm 2 and Figure 3.

216 We note that, in the explanation of these two time-marching schemes, as we are predicting only and not  
217 assimilating data, we predict the future value of both the POD coefficients and the sensor values. When  
218 assimilating data, the future value of the sensor values will be known, so  $\mathbf{s}^k$  can be treated as the other  
219 known variables. We also note that both prediction methods require  $m$  known values in order to make  
220 a prediction for the subsequent time level. Initially, these  $m$  values come from the high-fidelity model,  
221 however, after predicting for  $m$  time levels, predictions are based on previous predictions from the AAE  
222 and no more values are needed from the high-fidelity model.

---

**Algorithm 2** Prediction by backpropagation. By minimising the difference between the output of the decoder and the solutions (POD coefficients and sensor values) at the known time levels, a set of the latent variables which produce the known values can be obtained. The decoder can then be used to obtain a prediction for the solution at time level  $k$ .

---

- 1: Given the trained decoder ( $\mathcal{D}^{\text{ec}}$ ) and  $m$  previous solutions (POD coefficients and sensor values),  $\mathbf{x}^{k-1}, \mathbf{x}^{k-2}, \dots, \mathbf{x}^{k-m}$ , this algorithm will make predictions for the POD coefficients and sensor values at future time levels.
- 2: **for** time level  $k \in$  desired range of time levels **do**
- 3:     *!! set all  $d$  latent variables to random values at time level  $k$*
- 4:      $\tilde{\mathbf{z}}^k = \mathcal{N}^d(0, 1)$
- 5:     *!! evaluate the decoder*
- 6:

$$\begin{pmatrix} \tilde{\mathbf{x}}^k \\ \tilde{\mathbf{x}}^{k-1} \\ \vdots \\ \tilde{\mathbf{x}}^{k-m} \end{pmatrix} = \mathcal{D}^{\text{ec}}(\tilde{\mathbf{z}}^k) \quad (8)$$

- 7:     *!! minimise the difference between the output of the decoder and the desired output with respect to the latent variables — note that this is only done for the known values*
- 8:

$$\mathbf{z}^k = \arg \min_{\forall \mathbf{z}^k} \left( \begin{pmatrix} \tilde{\mathbf{x}}^{k-1} \\ \tilde{\mathbf{x}}^{k-2} \\ \vdots \\ \tilde{\mathbf{x}}^{k-m} \end{pmatrix} - \begin{pmatrix} \mathbf{x}^{k-1} \\ \mathbf{x}^{k-2} \\ \vdots \\ \mathbf{x}^{k-m} \end{pmatrix} \right) \quad (9)$$

- 9:     *!! evaluate the decoder for an input of  $\mathbf{z}^k$*
  - 10:      $\mathcal{D}^{\text{ec}}(\mathbf{z}^k)$
  - 11:     *!! Discard the AAE's approximation of the solution at the previous time levels, and keep the approximation for time level  $k$ .*
  - 12:      $\mathbf{x}^k = \tilde{\mathbf{x}}^k$
  - 13: **end for**
- 

### 223 2.3. Data Assimilation

224 In this section, we describe the method used to perform data assimilation in the reduced space (i.e. the  
 225 low-dimensional space found by POD). Results for time series predictions presented in Section 3.3 will  
 226 show that both the constrained prediction and prediction by backpropagation methods outlined in the  
 227 previous section give similar results, so the first method (constrained prediction) is used here for its  
 228 simplicity and fast convergence.

229 In order to perform data assimilation with the AAE as trained in the previous section, let us assume  
 230 that we have observations for all the sensor values with which we trained the AAE (that is CO<sub>2</sub>, relative  
 231 humidity, velocities, temperature). When predicting in time with the AAE (and not performing DA),



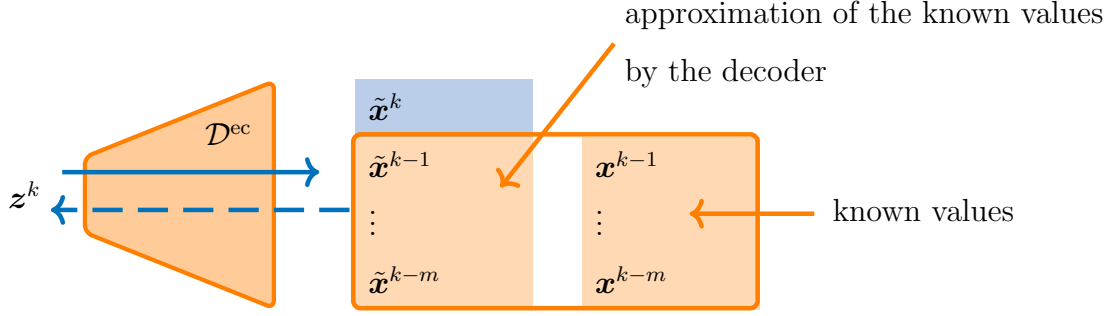


Figure 3: A schematic diagram of the “prediction by backpropagation” method. The difference between known solutions and solutions predicted by the decoder is minimised by backpropagating the difference through the decoder network. This finds a set of latent variables that corresponds to the known values. The remaining values in the output of the decoder are taken to be the solutions (POD coefficients and sensors values) at time level  $k$  ( $\tilde{\mathbf{x}}^k$ ).

232 the known values were the previous values of the POD coefficients and sensor values at  $m$  time levels,  
 233 and the unknown values were the POD coefficients and sensor values at the future time level, level  $k$  (see  
 234 Algorithm 1). Now, as we are assimilating data as well as predicting in time, the known values include  
 235 the future value of the observations. For the two different cases of prediction and data assimilation, the  
 known and unknown values are shown in Table 1. Algorithm 1 is rewritten in order to assimilate data with

	known values	unknown values
prediction (only)	$\{\boldsymbol{\alpha}^{k'}\}_{k'=k-m}^{k-1}$ and $\{\mathbf{s}^{k'}\}_{k'=k-m}^{k-1}$	$\boldsymbol{\alpha}^k$ and $\mathbf{s}^k$
data assimilation	$\{\boldsymbol{\alpha}^{k'}\}_{k'=k-m}^{k-1}$ and $\{\mathbf{s}^{k'}\}_{k'=k-m}^k$	$\boldsymbol{\alpha}^k$

Table 1: The known and unknown values used for prediction and for data assimilation.

236  
 237 the AAE in Algorithm 3. So, given an initial set of POD coefficients from time levels  $(1, 2, \dots, m)$  and  
 238 a set of observations throughout time,  $\{\mathbf{s}\}_{k=1}^{N^t}$ , we approximate the POD coefficients at time level  $m + 1$   
 239 using Algorithm 3. The first time we pass the input through the AAE the first time, we approximate  $\boldsymbol{\alpha}^k$   
 240 by  $\boldsymbol{\alpha}^{k-1}$ . Once we have a converged solution for the POD coefficients at the future time level,  $\boldsymbol{\alpha}^k$ , we  
 241 march forwards in time until the final time level is reached,  $N^t$ , say. Now we march backwards in time,  
 242 treating as known the observations  $\{\mathbf{s}^{k'}\}_{k'=N^t-m}^{N^t}$  and POD coefficients  $\{\boldsymbol{\alpha}^{k'}\}_{k'=N^t-m+1}^{N^t}$ , and as unknown,  
 243 the solution  $\boldsymbol{\alpha}^{N^t-m}$ . Once we have a converged solution at time level  $N^t - m$ , we continue marching back  
 244 in time until we reach the first time level. We march forwards and backwards in time until the method  
 245 converges, see Figure 4. There are two additions we make to this general procedure. First, rather than  
 246 minimising the least squares mismatch over time, we rely on the AAE to produce solutions consistent  
 247 with the known or fixed values, thereby producing solutions which we expect to reconcile the mismatch  
 248 between the observations and the predictions. We do, however, make use of a functional based on the

249 least squares mismatch over time, defined as follows:

$$250 \quad \mathcal{M} = \mathcal{M}^{\text{fwd}} + \mathcal{M}^{\text{bwd}} \quad (10)$$

$$251 \quad \mathcal{M}^{\text{fwd}} = \sum_{k=m+1}^{N^t} \left( \hat{\mathbf{s}}^k - \tilde{\phi}^k(\mathbf{x}_p) \right)^T \mathbf{W} \left( \hat{\mathbf{s}}^k - \tilde{\phi}^k(\mathbf{x}_p) \right) \quad (11)$$

$$252 \quad \mathcal{M}^{\text{bwd}} = \sum_{k=1}^{N^t-m} \left( \hat{\mathbf{s}}^k - \tilde{\phi}^k(\mathbf{x}_p) \right)^T \mathbf{W} \left( \hat{\mathbf{s}}^k - \tilde{\phi}^k(\mathbf{x}_p) \right) \quad (12)$$

253 where  $\hat{\mathbf{s}}^k$  are the sensor observations at time level  $k$ ,  $\mathbf{x}_p$  are the sensor locations,  $N^t$  is the number of time  
 254 levels over which data is assimilated,  $\mathbf{W}$  is the covariance matrix, here, taken as the identity matrix and  
 255  $\tilde{\phi}^k(\mathbf{x}_p)$  is the current approximation of the sensor values at time level  $k$ . This is calculated by using the  
 256 inverse of Equation (3),

$$257 \quad \tilde{\phi}^k = \mathbf{R}\tilde{\alpha}^k \quad \forall k, \quad (13)$$

258 and evaluating  $\tilde{\phi}^k$  at  $\mathbf{x}_p$  by interpolating from the nearest nodal values. This can be done by pre-  
 259 multiplying  $\tilde{\phi}^k$  by an interpolation matrix which will be sparse, containing non-zero values associated with  
 260 the nearest nodes. The functional is evaluated after every backward march. After the second backward  
 261 march (when there are at least two functional evaluations available), if the functional has increased in  
 262 value (and therefore the match to the observed data has become less good), we apply relaxation over the  
 263 next iteration, as follows:

$$264 \quad \alpha^k = \alpha^{k,\text{old}} + \gamma \left( \tilde{\alpha}^k - \alpha^{k,\text{old}} \right) \quad (14)$$

265 where  $\gamma \in [0, 1]$  controls the amount of relaxation applied,  $\alpha^{k,\text{old}}$  is the solution approximated during the  
 266 previous backwards march,  $\tilde{\alpha}^k$  is the solution predicted by the AAE, and  $\alpha^k$  is the approximation to the  
 267 solution at time level  $k$  (after relaxation has been applied). Initially, the value of  $\gamma$  is set to be one, and  
 268 when the functional increases the value of  $\gamma$  is halved.

269 The second addition made to the basic procedure of marching forwards and backwards in time is prioritisation.  
 270 By introducing a parameter  $\beta_{j\ell} \in [0, 2]$  for each observed value  $\ell$  at each sensor  $j$ , we can control  
 271 the importance of each observation. For example, if the observation for field  $\ell$  at the  $j$ th sensor at time  
 272 level  $k$  is represented by  $\hat{\mathbf{s}}_{j\ell}^k$ , then we prioritise as follows:

$$273 \quad \mathbf{s}_{j\ell}^k = \beta_{j\ell} \hat{\mathbf{s}}_{j\ell}^k + (1 - \beta_{j\ell}) \tilde{\mathbf{s}}_{j\ell}^k, \quad (15)$$

274 where the tilde sign ( $\tilde{\phantom{x}}$ ) represents the output of the AAE and  $\mathbf{s}_{j\ell}^k$  represents the approximation to this  
 275 sensor value which will be used as in input to the AAE at the next iteration. The observation itself,  
 276  $\hat{\mathbf{s}}_{j\ell}^k$ , is fixed. If this observation is to be ignored, then  $\beta_{j\ell}$  can be set to zero. If the AAE is to be  
 277 relied upon to enforce the observation, then  $\beta_{j\ell}$  can be set to one. If more priority is to be given to  
 278 this observation then  $\beta_{j\ell} > 1$ . This is a very flexible way of ignoring some observations and prioritising  
 279 other observations without having to retrain the neural network. See Algorithm 4 for a description of  
 280 the DA algorithm proposed in this paper.  $\beta_{j\ell} = 1$  might be seen as the default prioritisation. Figure 5  
 281 illustrates the DA process. The blue loop represents the inner iteration which, for the forward march,  
 282 obtains the POD coefficients at the future time level,  $\alpha^k$  given the previous  $m$  values of POD coeffi-  
 283 cients,  $\{\alpha^{k-1}, \alpha^{k-2}, \dots, \alpha^{k-m}\}$  and observations  $\{\mathbf{s}^k, \mathbf{s}^{k-1}, \dots, \mathbf{s}^{k-m}\}$ . The forward march iterates from  
 284  $k = m + 1$  to  $N^t$  with an increment of 1, and once complete, the backward march begins, where the

285 solution at time level  $k$  is obtained from future solutions  $\{\boldsymbol{\alpha}^{k+m}, \boldsymbol{\alpha}^{k+m-1}, \dots, \boldsymbol{\alpha}^{k+1}\}$  and observations  
 286  $\{\mathbf{s}^{k+m}, \mathbf{s}^{k+m-1}, \dots, \mathbf{s}^k\}$ . The backward march iterates from  $k = N^t - m$  to 1 with a decrement of 1. One  
 287 forward and one backward march makes up an outer iteration which is indicated in red. When the outer  
 288 iterations converge, one final forward march is performed before the algorithm finishes.

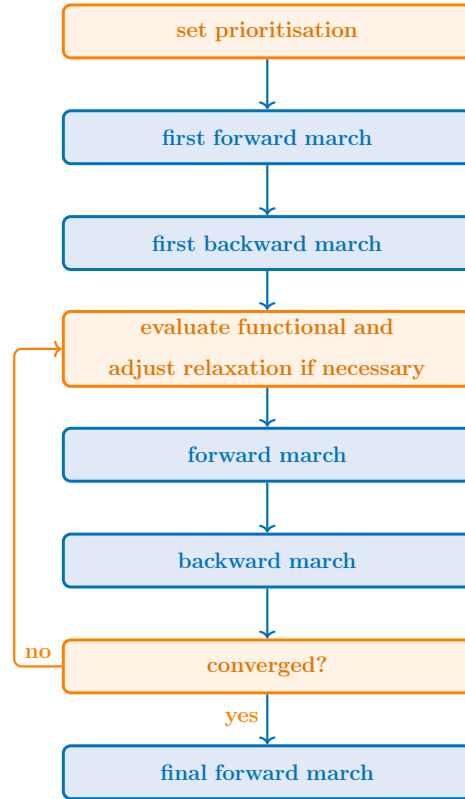


Figure 4: Forward and backward marching that occurs when performing the data assimilation.

---

**Algorithm 3** Constrained Prediction for use with DA. By constraining the values of the sensors and POD coefficients at the previous (known) time levels, and the sensor values at the future time level, we can approximate the POD coefficients at the future time level.

---

```

1: This algorithm will make improved predictions for the POD coefficients for each time level in the
   set of desired values  $\mathcal{K}$ , given the maximum number of iterations ( $N^{its}$ ); the trained AAE ( $f^{AAE}$ ); a
   tolerance for the stopping criterion ( $\varepsilon$ );  $m$  previous solutions (POD coefficients),  $\{\alpha^{k'}\}_{k'=k-m}^{k-1}$ ;  $m+1$ 
   sensor values,  $\{s^{k'}\}_{k'=k-m}^k$ ; the range of time levels  $\mathcal{K}$ ; and the relaxation parameter  $\gamma$ .
2: for time level  $k \in \mathcal{K}$  do
3:   !! approximate the solution at time level k
4:   if forward then
5:      $\alpha^{k,0} = \alpha^{k-1}$ 
6:   else
7:      $\alpha^{k,0} = \alpha^{k+1}$ 
8:   end if
9:   for iteration  $i = 1, 2, \dots, N^{its}$  do
10:    form  $\mathbf{X}$  and pass through the autoencoder
11:     $\tilde{\mathbf{X}} = f^{AAE}(\mathbf{X})$ 
12:    !! update the approximation of the POD coefficients using values in  $\tilde{\mathbf{X}}$ 
13:     $\alpha^{k,i} = \tilde{\alpha}^{k,i-1}$ 
14:    if ( $\|\tilde{\alpha}^{k,i-1} - \alpha^{k,i-1}\|_2^2 < \varepsilon$  or  $i == N^{its}$ ) then
15:       $\alpha^k = \alpha^{k,i}$ 
16:      exit iteration loop
17:    end if
18:  end for
19: end for

```

---

---

**Algorithm 4** This algorithm details the data assimilation algorithm proposed here, which marches forwards and backwards through time, adjusting the relaxation parameter,  $\gamma$ , according to the value of the mismatch functional.

---

```

1: set prioritisation coefficients  $\beta_{j\ell}$  for all the sensor values
2: set  $\mathcal{M}^0 = 10^{10}$ 
3: set  $\gamma = 1$ 
4: for iteration  $i = 1, \dots, N^{its}$  of outer loop do
5:   !! forward march with  $\mathcal{K} = \{m + 1, m + 2, \dots, N^t\}$  using Algorithm 3
6:   for  $k = m + 1, m + 2, \dots, N^t$  do
7:     predict  $\alpha^k$  given  $\{\alpha^{k'}\}_{k'=k-m}^{k-1}$  and  $\{\mathbf{s}^{k'}\}_{k'=k-m}^k$ .
8:   end for
9:   !! backward march with  $\mathcal{K} = \{N^t - m, N^t - m - 1, \dots, 1\}$  using Algorithm 3
10:  for  $k = N^t - m, N^t - m - 1, \dots, 1$  do
11:    predict  $\alpha^k$  given  $\{\alpha^{k'}\}_{k'=k-m}^{k-1}$  and  $\{\mathbf{s}^{k'}\}_{k'=k-m}^k$ .
12:  end for
13:  !! relaxation
14:  evaluate the functional,  $\mathcal{M}$ , at iteration  $i$  using Equation (10)
15:  if ( $\mathcal{M}^i > \mathcal{M}^{i-1}$ ) then
16:     $\gamma \leftarrow 0.5\gamma$ 
17:  end if
18:  !! check convergence
19: end for
20: !! final forward march with  $\mathcal{K} = \{m + 1, m + 2, \dots, N^t\}$  using Algorithm 3
21: for  $k = m + 1, \dots, N^t$  do
22:   predict  $\alpha^k$  given  $\{\alpha^{k'}\}_{k'=k-m}^{k-1}$  and  $\{\mathbf{s}^{k'}\}_{k'=k-m}^k$ 
23: end for

```

---

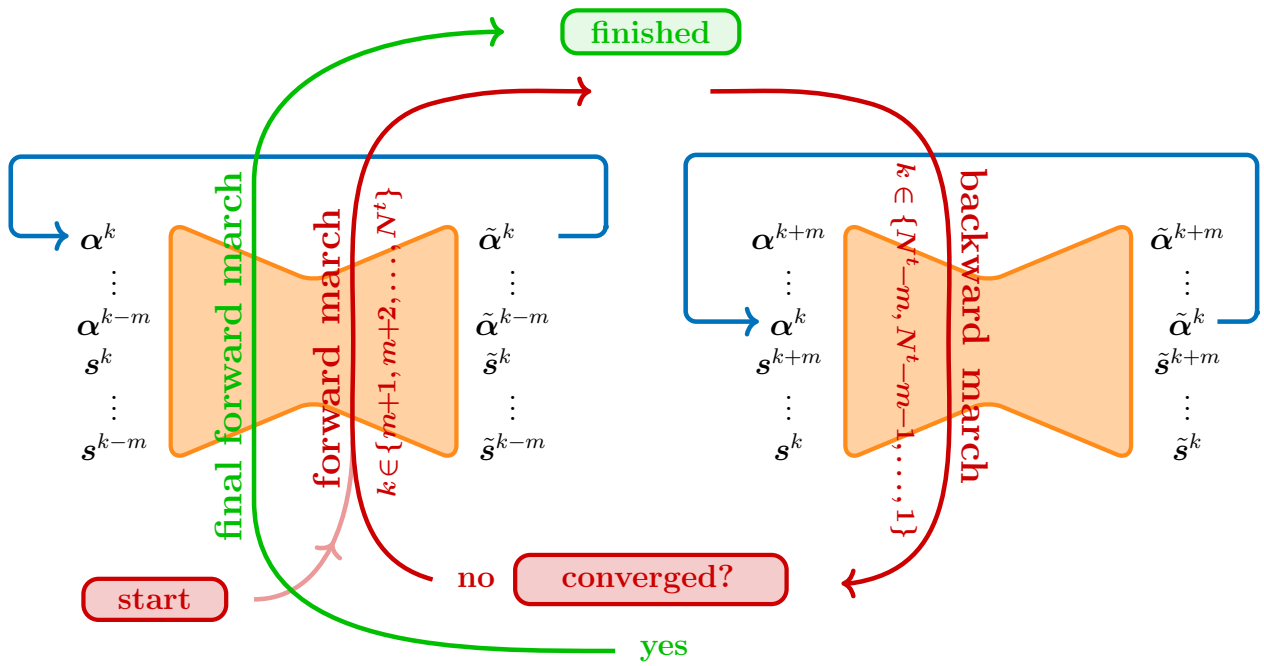


Figure 5: A schematic diagram of the data assimilation process.

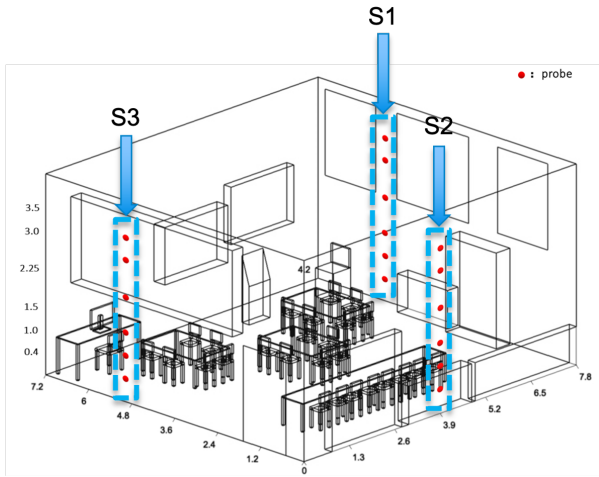
### 289 3. Results

290 First we present a description of the test case, explaining where the observations were taken and how  
291 the test case was modelled numerically. In Section 3.2 we give details on how the dimension of the CFD  
292 results was reduced; in Section 3.3 we present predictions by the surrogate model of the air flows and air  
293 quality in the room. These predictions are based on the constrained method of prediction. In Section 3.4  
294 we describe how the data assimilation approach was tested; first, on a dual-twin type experiment and  
295 second, with observations from a classroom.

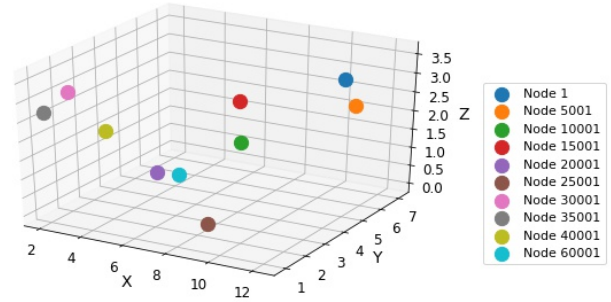
#### 296 3.1. Test Case

297 A classroom in Houndsfield Primary School, north London, is used to demonstrate the proposed meth-  
298 ods. As part of a measurement campaign to investigate indoor air quality, observations of CO<sub>2</sub>, relative  
299 humidity and temperature were collected at 18 locations over a two week period [38]. From this, we select  
300 observations collected over an hour whilst the classroom was occupied by 26 students and one teacher.  
301 The measurements have been time-averaged over 1 min. Heating was supplied by a radiator and a Dyson  
302 fan, the latter was also used to clean the air and remove the virus-laden particles from the air. Figure 6  
303 shows a schematic diagram of the classroom, including desks, chairs, windows, the door to an internal  
304 corridor, the fan, the radiator (on the same wall as the sensors labelled S1) and the sensor locations. The  
305 classroom measures 7.2 m by 7.8 m with a height of 4.2 m. The door to the internal corridor was open  
306 (seen in Figure 6a on the same wall as sensors labelled S3) and the top-hinged windows on the right were  
307 also open (seen in Figure 6a opposite the door to the internal corridor and near the sensors labelled S1.  
308 Another view of this can be seen in Figure 7, which includes the teacher and students.

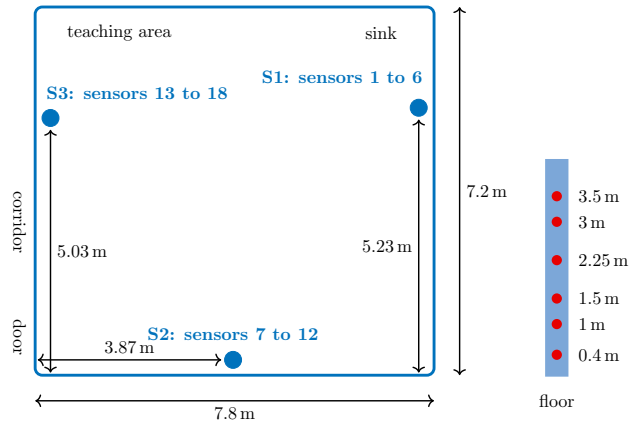
309 In addition to collecting observations, a Large Eddy Simulation (LES) was performed to generate high-  
310 fidelity numerical results of the air flows and air quality in the classroom by using the open source CFD  
311 code Fluidity [42, 43]. Based on the finite element method, this code solves the Navier-Stokes equations for  
312 unsteady, incompressible and viscous flow with an advanced mesh capability, which includes unstructured,  
313 anisotropic and adaptive meshing [44]. Fluidity is able to predict the evolution of components in the air  
314 (i.e., CO<sub>2</sub>, water vapour, temperature and virus concentration) by the transportation of scalar fields as  
315 tracers. The principle behind the LES model is based on fully resolving large-scale turbulent structures  
316 and modelling the less energetic small-scale structures through a subgrid-scale model. With this method,  
317 Fluidity has been able to resolve turbulent features of urban flows and the dispersion of passive tracers  
318 down to length scales of 10 cm and time scales of less than a second [45–47]. Thermal stratification  
319 is modelled by the Boussinesq approximation which assumes that the predominant influence of density  
320 variation is on the buoyancy forces. The geometry of the classroom is created using the Indoor Geometry  
321 Generator (IGG) [48], which is a software package that allows easy definition of objects such as furniture,  
322 doors, and windows, as well as people sitting or standing. Figure 7 shows the layout of the room used for  
323 the CFD simulations. The dimensions and locations of desks, students, fan and radiator are identical to  
324 those of the classroom. To replicate the indoor-outdoor exchange of air, there are two additional buffer  
325 regions attached to the left and right of the classroom. Natural ventilation is simulated by imposing  
326 inflow and outflow boundary conditions on certain walls of the buffer regions, which allows air to go in  
327 or out of the room through the door which leads to the internal corridor (left) and through the top-  
328 hinged windows (right). In Figure 7, the inlets are shaded yellow and the outlets are shaded green. Inlet



(a) Layout of the classroom and its contents



(b) Node locations for some time history plots



(c) Sensor locations

Figure 6: Classroom geometry, sensor and node locations used for time history plots. (a) and (c) Location of the 18 sensors used to collect  $\text{CO}_2$ , relative humidity and temperature data. (Sensors 1 to 6 share the same  $x$  and  $y$  coordinates, with the sensor index increasing with height. Similarly for sensors 7 to 12 and 13 to 18.) See Kumar et al. [38] for more details.



329 velocities are specified on the nearest facing planes of the buffer regions of 0.5 m/s (left) and 1.0 m/s  
 330 (right). Open boundary conditions are specified on the far walls of the buffer regions (shaded green), by  
 331 setting non-hydrostatic pressure to zero. On all other external boundaries and internal walls, floors and  
 332 ceilings, we use a natural boundary condition of zero shear stress with a normal velocity of zero. For  
 333 any remaining surfaces, such as the surface of desks, objects and people, a no-slip boundary condition is  
 334 applied. Regarding boundary conditions for temperature, zero heat flux was applied on the walls, floor,  
 335 furniture and ceiling. Within the classroom, a Dirichlet temperature boundary condition was applied  
 336 to the radiator enforcing a temperature of 32 °C. In order to account for the body heat of individuals,  
 337 a Robin or heat transfer boundary condition was applied imposing a temperature of 34 °C and using a  
 338 surface heat transfer coefficient of 10 W m<sup>-2</sup> K<sup>-1</sup>. A Dirichlet temperature boundary condition was also  
 339 applied at the largest monitor with a temperature of 40 °C, which can be seen in Figure 6a on the desk  
 340 near the array of sensors labelled S3. Dirichlet boundary conditions are also adopted at the ‘mouth’  
 341 and ‘nose’ area of people for other passive scalar fields (i.e., CO<sub>2</sub>, relative humidity, water vapour and  
 342 virus concentration), in order to reproduce the inhalation and exhalation process. Initial conditions are  
 343 given in Table 2 and boundary conditions for CO<sub>2</sub>, relative humidity and viral load are consistent with  
 344 these values. The simulation is started impulsively, so the value of velocity is adjusted until it satisfies  
 345 continuity during the first time step. An adaptive time stepping scheme is applied in this simulation,  
 346 which means that the time step is chosen so that a Courant number of 10 is obtained. The anisotropic  
 347 unstructured mesh is adapted so that the smallest element edge length is 0.1 m and the largest element  
 348 edge length is 2.5 m.

	left buffer	school classroom	right buffer
CO <sub>2</sub> (ppm)	400	430	400
temperature (K)	279	289	279
relative humidity (%)	40	40	40
viral load (copies)	0	0	0
velocity	<i>impulsively started</i>		

Table 2: The initial conditions used in the CFD simulation. The values are constant throughout each region.

349 Solutions of the high-fidelity model were generated every 5 s over a period of one hour. Although the  
 350 solutions were generated using mesh adaptivity, in order to apply POD, all the solutions were interpolated  
 351 onto the unstructured mesh that was generated at the end of the simulation, which had 192,060 nodes.  
 352 Using solutions on different-sized meshes would be possible, but would require the use of space-filling  
 353 curves [49], spatially varying kernels [50] or graph neural networks [51]. Here we choose to interpolate all  
 354 the solutions onto a single (unstructured) mesh, as a fast and straightforward option whilst we investigate  
 355 the capability of the newly proposed DA method. The solutions or snapshots were separated into two  
 356 groups: a training or seen dataset corresponding to a time window of [5, 2880] seconds (576 time levels);  
 357 and a test or unseen dataset corresponding to a window of [2885, 3600] seconds (144 time levels). Within  
 358 these time windows, solutions are available every 5 s.

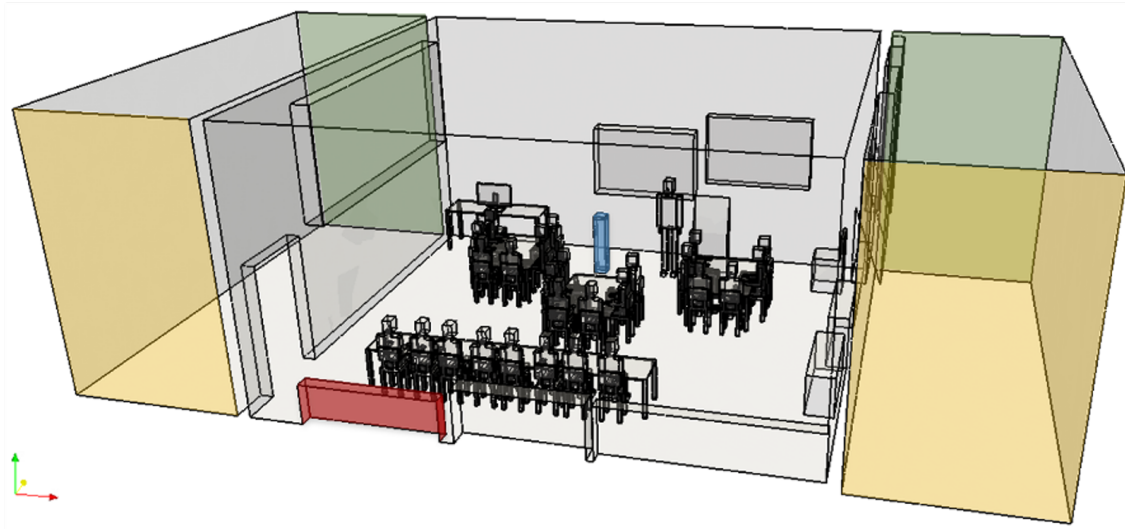


Figure 7: Geometry of the classroom at Houndsfield Primary School used in the CFD simulations showing the furniture, 26 students and the teacher. The blue rectangular box represents a Dyson fan; the red rectangular box represents a radiator. The two yellow surfaces denote two inlets with imposed velocities of 0.5 m/s (left) and 1.0 m/s (right), where incoming air has the temperature of 279 K. The two green surfaces denote two outlets with open pressure conditions.

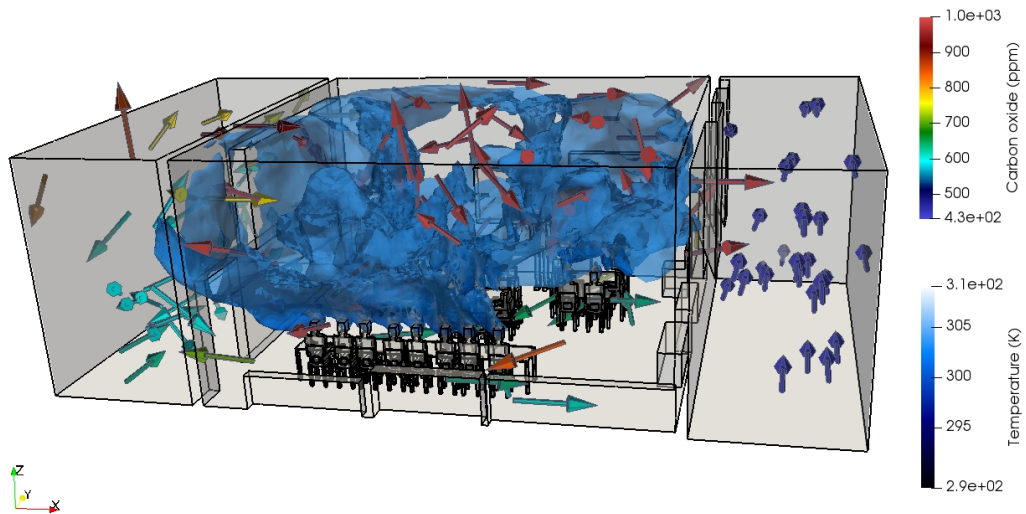


Figure 8: Thermal  $\text{CO}_2$  plume within the naturally ventilated classroom simulated after one hour. The iso-surface represents a level of 800 ppm of  $\text{CO}_2$  and is coloured by temperature; arrows represent air flows and are coloured by  $\text{CO}_2$  level.

359 *3.2. Dimensionality Reduction*

360 POD is applied to the training dataset to reduce the dimension of the problem. Seven solution fields (CO<sub>2</sub>,  
 361 velocities, temperature, relative humidity and viral load) defined on an unstructured mesh of 192,060 nodes  
 362 are approximated by 150 POD basis functions and coefficients, which represents capturing 90.2% of the  
 363 information contained by the snapshots in the training dataset. This corresponds to a compression ratio  
 364 of 1,344,420 (7 fields, each of 192,060 nodes) to 150, or almost 9000 to 1. Before applying POD, the seven  
 365 fields are normalised so their values lie in the range [-1,1]. After applying POD, further normalisation is  
 366 applied to the POD coefficients, so their values lie in the range [0,1], as they will be used to train a neural  
 367 network to predict the evolution of the POD coefficients. This is standard machine learning practice.  
 368 Once calculated, these normalisation functions are applied, as needed, to any unseen data. Figure 9a  
 369 shows the variance captured by the singular values (or principal components), and Figure 9b shows how  
 370 the information captured by the POD basis functions decays with increasing numbers of basis functions.  
 371 Although 150 POD coefficients capture over 90% of the information, at this stage the singular values  
 372 are decaying slowly. Retaining more coefficients would improve the accuracy, however, this would make  
 373 training the networks more difficult, so using 150 coefficients was seen as a good compromise between  
 374 accuracy and tractability of training the networks. The POD coefficients associated with the training  
 375 dataset,

$$\alpha^k = \mathbf{R}^T \phi^k \quad \forall k \in \text{training dataset}, \quad (16)$$

377 will be used to train the neural network that will predict in time and assimilate data.

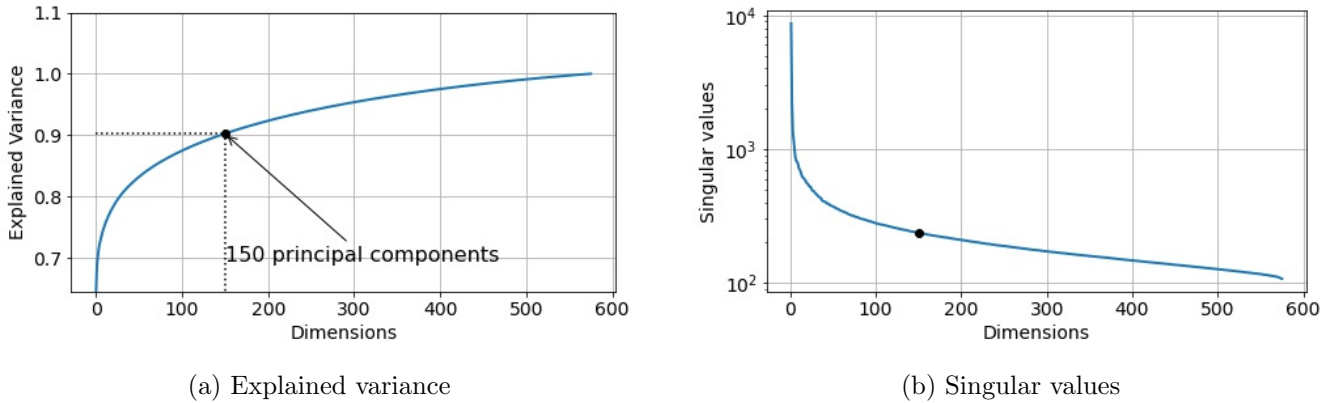


Figure 9: POD is applied to the snapshots in the training dataset in order to reduce the dimension from 1,344,420 to 150 (a compression ratio of almost 9000). Retaining 150 POD basis functions is equivalent to capturing 90.2% of the information held in the snapshots in the training dataset. The so-called explained variance is plotted on the left, and the decay of the singular values is plotted on the right.

378 *3.3. Prediction*

For prediction, an adversarial autoencoder with convolutional layers is used. For training the network, the data is assembled in the form of 2D arrays, using (normalised) POD coefficients and normalised sensor values at a number of successive time levels. The dimension of the input to the network is the number of consecutive time levels used by the number of POD coefficients + (number of sensors × number of fields observed at each sensor):

$$9 \text{ by } 150 + (18 \times 6), \text{ that is } 9 \text{ by } 258. \quad (17)$$

379 See Table 3 for a diagram of the structure of the 2D data. For flexibility, we allow for six fields of  
 380 sensor data. In the dual-twin experiment we use data from the high-fidelity model as the sensor data (by  
 381 interpolating the high-fidelity model at the sensor locations). When assimilating the observations from  
 382 the classroom, we have observations for CO<sub>2</sub>, temperature and humidity, and we set the other values  
 by interpolating the high-fidelity model at the sensor locations. The adversarial autoencoder network

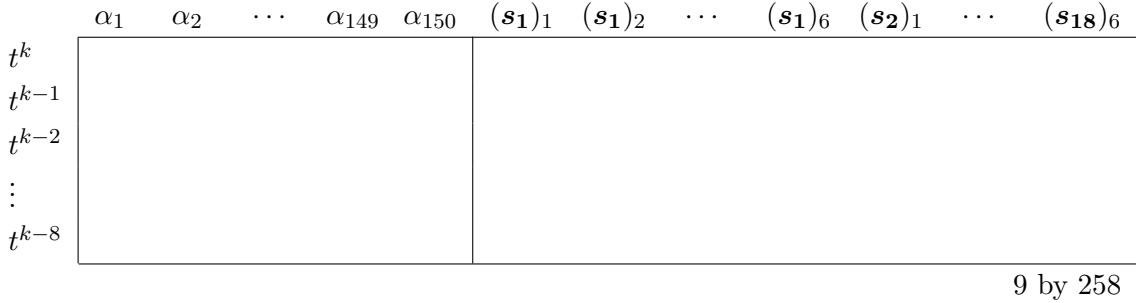


Table 3: The 2D structure of the data used for training, with 9 rows (samples) and 258 columns (features). The  $i$ th POD coefficient is represented by  $\alpha_i$ . Values associated with the  $j$ th sensor are represented by  $\mathbf{s}_j$ . Each sensor will have 6 values, one for each of the six fields  $\{(\mathbf{s}_j)_\ell\}_{\ell=1}^6$ .

383 was implemented in TensorFlow [52] and a schematic diagram of the network can be seen in Figure 1.  
 384 After optimisation, the best values and settings for the hyperparameters were found and are shown in  
 385 Tables A.4 and A.5.  
 386

387 Figures 10a to 10g demonstrate how well the non-intrusive reduced-order model, referred to here as  
 388 PredAAE, predicts for the seven fields of interest (CO<sub>2</sub>, velocities, temperature, relative humidity and  
 389 viral load). The orange curves represent time series results from the high-fidelity model and the blue curves  
 390 are the time series predictions made by the PredAAE using the constrained prediction method. The time  
 391 series were collected at the location of three of the nodes of the unstructured mesh (see Figure 6c). The  
 392 training data is taken from the time interval [5,2880] seconds. The test data is taken from the remaining  
 393 time [2885, 3600] seconds. PredAAE predicts well over the training data for velocities and temperature,  
 394 but for CO<sub>2</sub>, humidity and viral load, the predictions start to stray from the high-fidelity model results.  
 395 This difference can be seen when the predictions are still within the training data time interval. A similar  
 396 collection of plots can be seen in Figure 11 for when the PredAAE method uses the backpropagation  
 397 method. The results shown here are broadly similar to those for the constrained method presented in  
 398 Figure 10, although for relative humidity and viral load, there is some improvement in the accuracy of the  
 399 predictions. We emphasise here, that (for both constrained and backpropagation methods) the reduced-  
 400 order model uses solutions at 8 previous time levels to approximate the solution at a future time level.  
 401 The model therefore requires 8 time levels as initial conditions, but once the model has made predictions  
 402 for 8 time levels, the predictions are based solely on outputs of the model and do not require any further  
 403 CFD solutions.

404 A comparison of the convergence of the two prediction methods can be seen in Figure 12, plotting the  
 405 mean square error between predictions at successive iterations. This shows that the constrained method  
 406 converges very quickly (within a few iterations), whereas the backpropagation method takes longer to  
 407 converge but converges to a lower error.

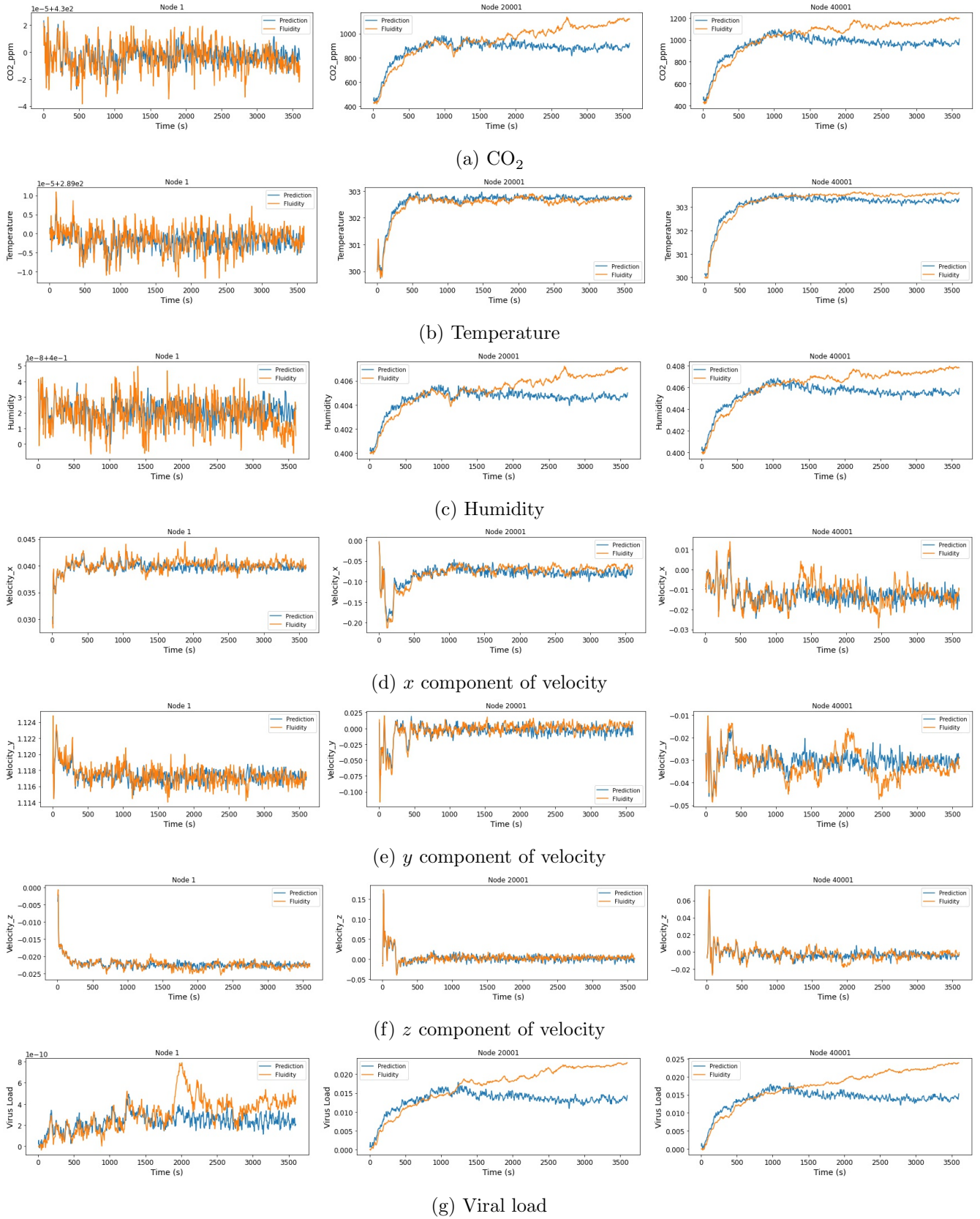


Figure 10: Time series plots of the seven fields at three locations corresponding to nodes 1, 20001 and 40001 (see Figure 6c). Orange represents the high-fidelity model results and blue represents the PredAAE prediction with the constrained method.

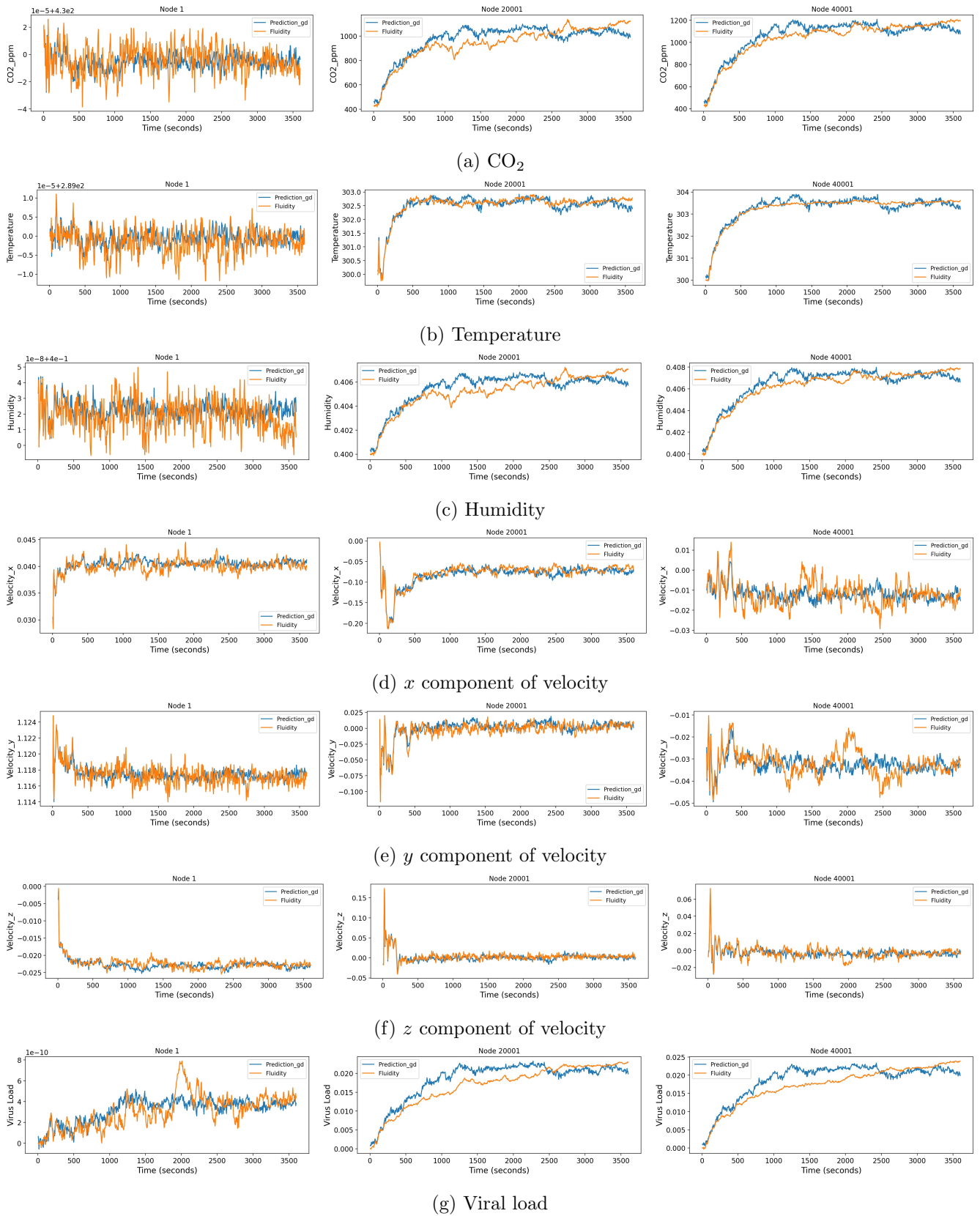
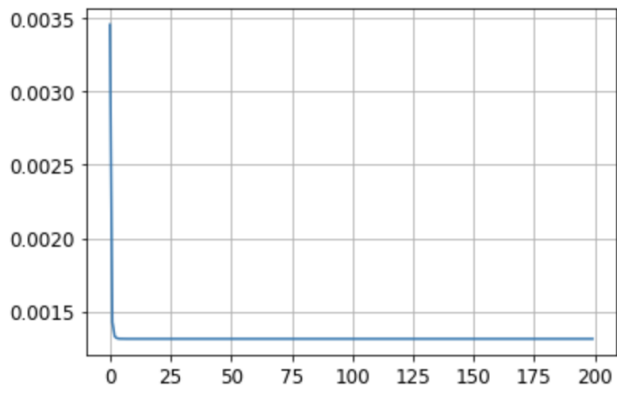
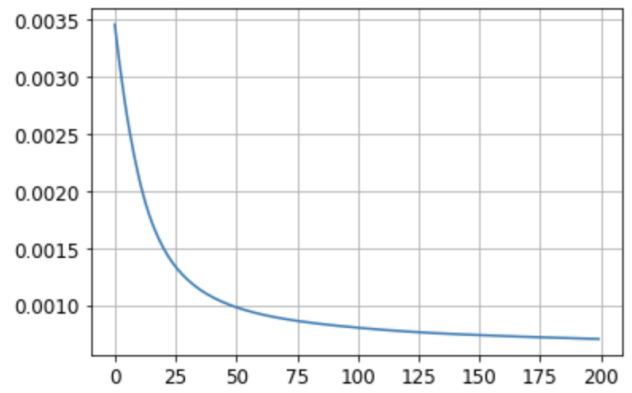


Figure 11: Time series plots of the seven fields at three locations corresponding to nodes 1, 20001 and 40001 (see Figure 6). Orange represents the high-fidelity model results and blue represents the PredAAE prediction with the backpropagation method.



(a) Constrained method



(b) Backpropagation method

Figure 12: Comparison between the two methods of predicting in time with the PredAAE. The mean square error between solutions at successive iterations is shown on the  $y$  axis, the number of iterations on the  $x$  axis.



### 408 3.4. Assimilating data

409 Once trained, the PredAAE reduced-order model is used to assimilate data. We investigate two scenarios  
410 in this paper: (1) assimilating data from the test dataset, where the sensor values are provided by  
411 the high-fidelity model results from 2885 to 3600 seconds, and the POD coefficients are taken from  
412  $t \in \{5, 10, \dots, 40\}$  and marched forward to 720s (sometimes referred to as a dual twin experiment);  
413 (2) assimilating observations collected from a classroom.

#### 414 3.4.1. Scenario 1: dual-twin experiment

415 To test the capability of the data assimilation algorithm, we perform a dual-twin experiment. The AAE  
416 expects an input of POD coefficients and sensor values at one particular time level. For the dual-twin  
417 experiment, the input of the AAE is consists of sensor data from time 2885s to 3600s (taken from the test  
418 dataset) and the POD coefficients taken from 5s to 720s (using “initial conditions” from 5s to 40s and  
419 predicting/assimilating data with the AAE up to 720s). After assimilating the sensor values, the aim is  
420 to see whether the POD coefficients (having been mapped to the solution space) are consistent with the  
421 solutions from 2885s to 3600s. If so, the data assimilation method is judged to be successful. Figure 13  
422 shows the solutions for CO<sub>2</sub>, temperature and humidity after assimilating data from the sensors taken  
423 during the test dataset. This illustrates clearly that before data assimilation, the predictions (from the  
424 POD coefficients, mapped to the solution space) are aligned with the CFD model from the time interval  
425 [5, 720], whereas after assimilation, the predictions are closely aligned with the CFD solution from the  
426 time interval [2885, 3600]. This tells us that the data assimilation algorithm is working well. Results are  
427 also shown at a number of nodes, see Figure 14, from which the same conclusion can be drawn, that before  
428 data assimilation, the predictions are consistent with the initial conditions and the CFD solutions from  
429 the time interval [5, 720], whereas after assimilation, the predictions are closely aligned with the behaviour  
430 associated with the observations, i.e. the CFD solution from the time interval [2885, 3600]. The prediction  
431 for CO<sub>2</sub> (Figure 14a) broadly follows the trend seen in the relative humidity (Figure 14c) and the viral  
432 load (Figure 14f), which has a similar source from the breath to CO<sub>2</sub> but has a half-life of SARS-CoV-2 of  
433 35 minutes. A little more variation can be seen in CO<sub>2</sub> than the viral load or relative humidity though.  
434 This suggests that relative humidity, CO<sub>2</sub> and viral load all share the same source (exhalation) whereas  
435 temperature, which does not follow the same trend, has a different predominant source (heat from the  
436 radiator). The spatial variation for the dual twin simulations can be seen in Figure 15, which shows  
437 an iso-surface of CO<sub>2</sub> for the value 1000 ppm at two times: at 100s (or 2985s) and at 500s (or 3385s).  
438 Three different cases are shown: the predictions of PredAAE (with no data assimilation), predictions of  
439 PredAAE with data assimilation (DA-PredAAE) and the original high-fidelity model. This again shows  
440 that, after data assimilation, the prediction agrees closely with the CFD model taken from the time  
441 interval that is consistent with the observations rather than the initial conditions.

#### 442 3.4.2. Scenario 2: assimilating observations from experiments

443 In this section, we assimilate observations collected with sensors from a classroom in Houndsfield Primary  
444 School into predictions of the surrogate model. Sensor values for CO<sub>2</sub>, humidity and temperature are  
445 available, however, sensor values for the velocity fields were not measured and are taken from the high-  
446 fidelity model results instead. These values are required in the input (and output) of the AAE, but, by



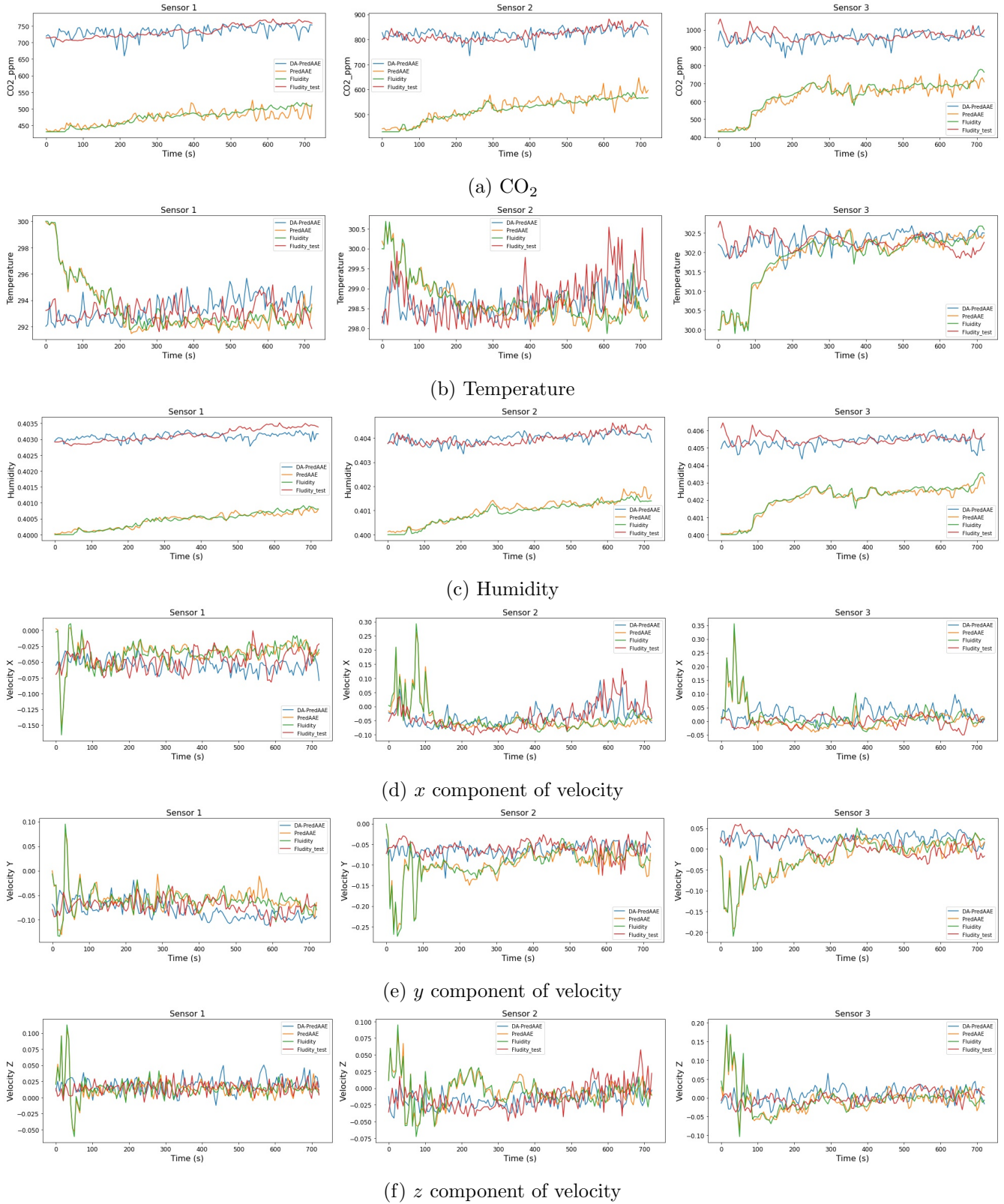


Figure 13: Time histories plotted at sensor locations for the dual twin experiment. The AAE is used to predict for time interval [5,720] seconds but given the sensor data from the time interval [2885,3600]. The above plots show that after assimilating the sensor data the predictions of the AAE match closely those from [2885,3600].

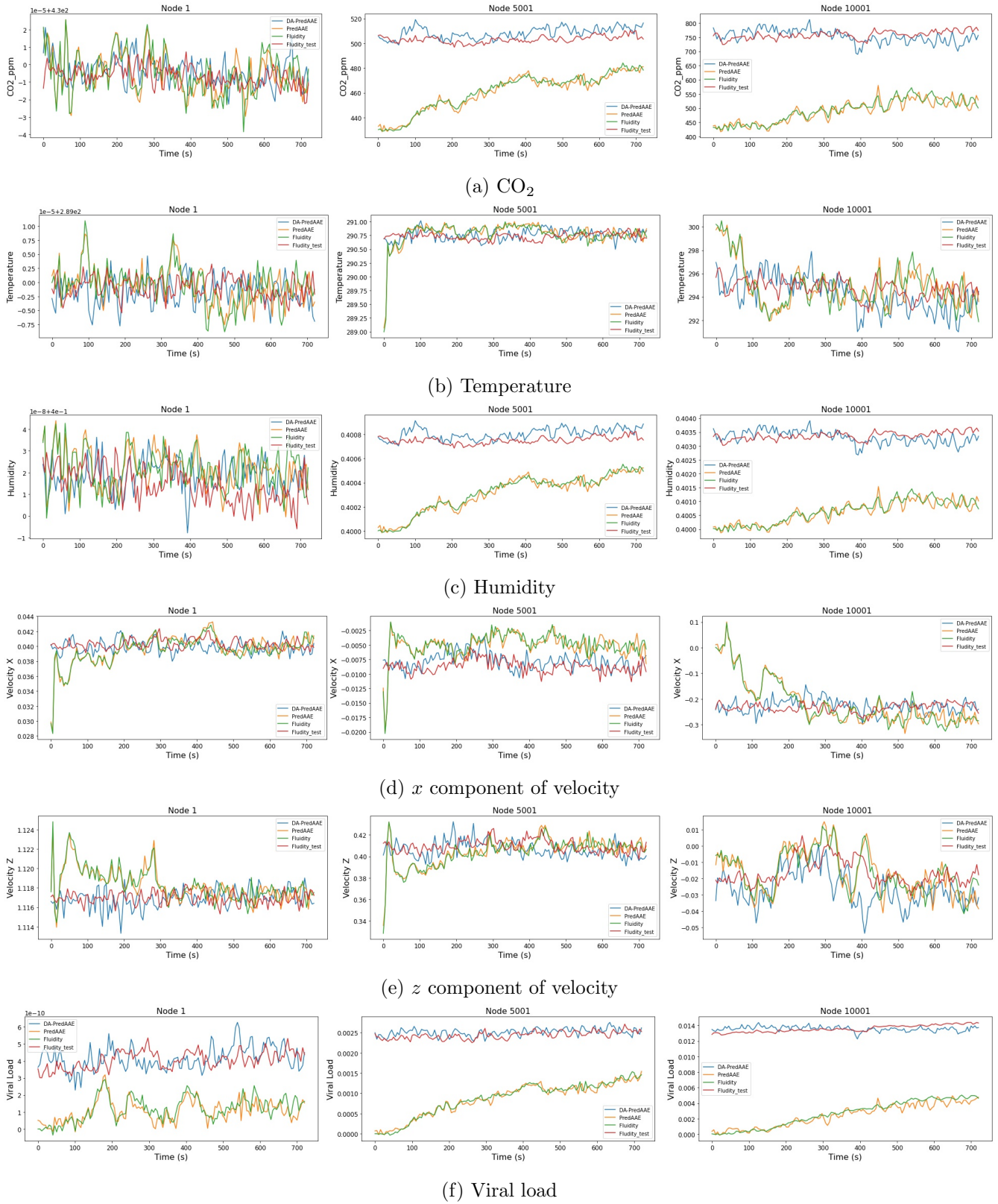


Figure 14: Time histories plotted at a number of nodal locations for the dual twin experiment (see Figure 6c). The AAE is used to predict for time interval [5, 720] seconds but given the sensor data from the time interval [2885, 3600]. The above plots show that after assimilating the sensor data the predictions of the AAE match closely those from [2885, 3600].

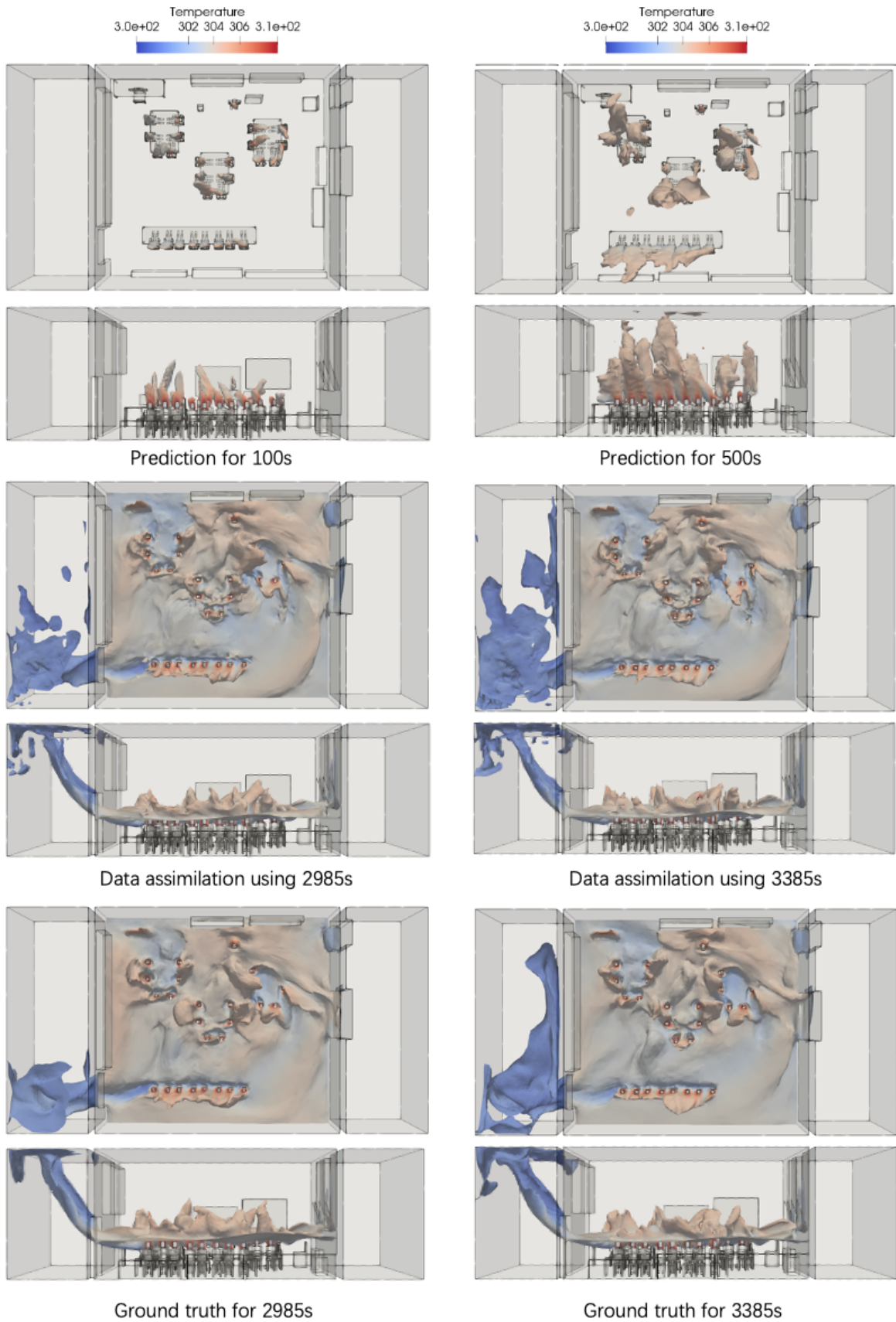


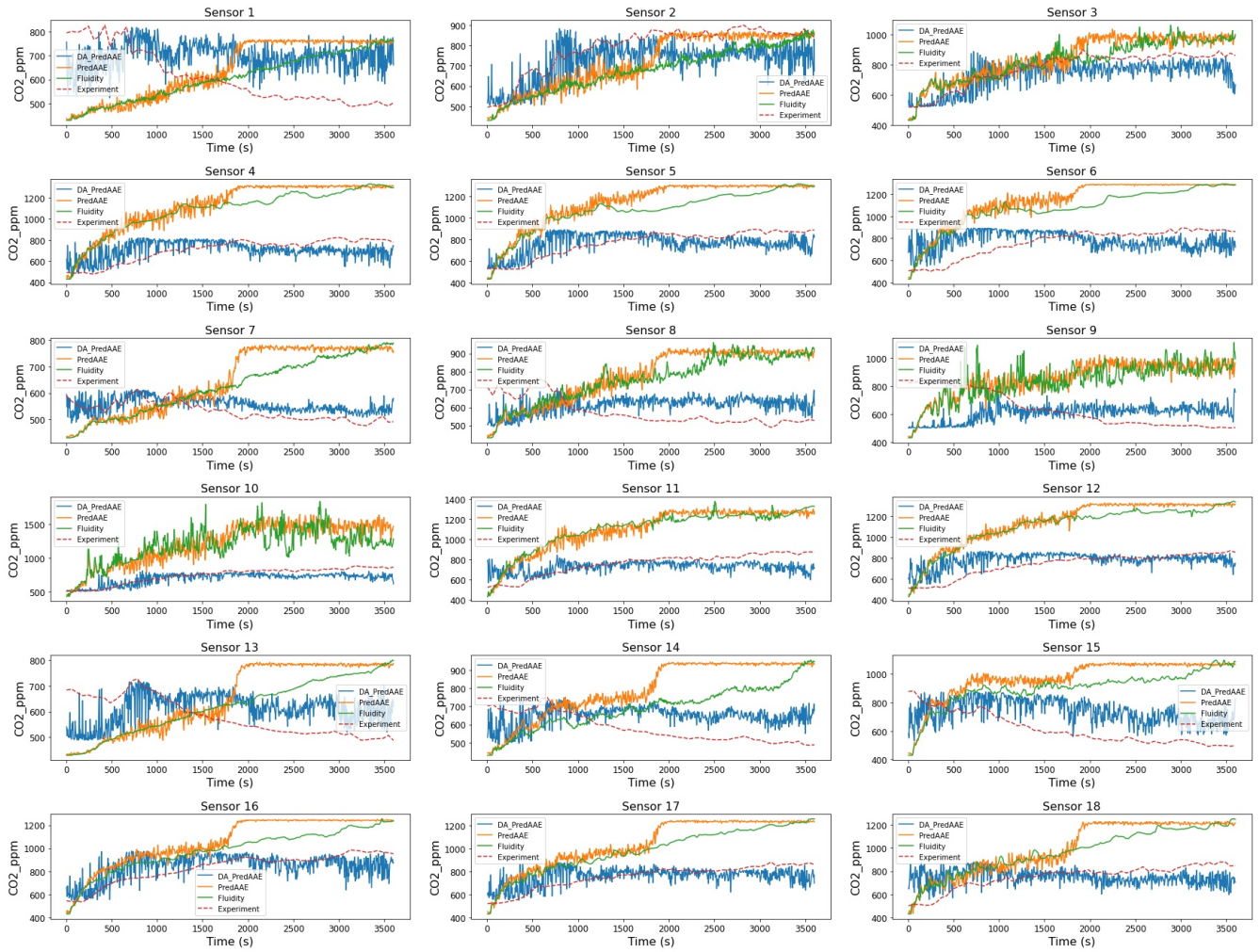
Figure 15: Spatial variation for the dual twin experiment at 100s/2985s (left) and 500s/3385s (right). An iso-surface of  $\text{CO}_2$  at 1000 ppm is coloured by temperature. At two points in time (shown on the left and right), we compare the prediction of PredAAE (top plots) with results from data assimilation DA-PredAAE (centre) with the CFD results (bottom).

447 setting the prioritisation parameter to zero for these fields, are not assimilated. The algorithm we propose  
448 has the flexibility to switch on and off which variables are assimilated. The solutions resulting from the DA  
449 process might be assumed to be a good estimate of the fields within the classroom and an interpolation  
450 of the sensor data in space and time. Figure 16 shows the results from the high-fidelity model, the  
451 reduced-order model without DA (PredAAE), the reduced-order model with DA (DA-PredAAE) and  
452 the sensor values in time series plots, at the sensor locations. In all cases, the predictions of PredAAE  
453 are close to the results from the high-fidelity model. For CO<sub>2</sub>, temperature and humidity, assimilating  
454 data (DA-PredAAE) successfully brings the model predictions closer to the sensor values, especially for  
455 CO<sub>2</sub> and temperature. Figure 17 plots time series predictions from the high-fidelity model, predictions  
456 from the reduced-order model without data assimilation (PredAAE) and predictions from the reduced-  
457 order model with data assimilation (DA-PredAAE), at a number of nodal locations. Here one can see  
458 that the predictions of PredAAE are close to the high-fidelity model, and that, after assimilating the  
459 data, the predictions deviate from the high-fidelity model, showing that the data assimilation as had an  
460 effect. Figure 18 shows the spatial variation at time 1000s for the high-fidelity model, PredAAE and  
461 DA-PredAAE. Although we do not know the conditions associated with the sensor values from the school  
462 (e.g. wind direction), we can see from Figure 16 that after assimilation, the predictions are closer to  
463 the sensor values, and we expect that the predictions shown in Figure 18 are more consistent with the  
464 conditions in the classroom, given the previous results which demonstrate the ability of the proposed  
465 method to assimilate data (Section 3.4.1).

#### 466 4. Conclusions

467 We have proposed a new method which assimilates data whilst marching forwards and backwards in time,  
468 as is done in 4D-Variational approaches. To make the method computationally tractable a surrogate  
469 model is used, which combines proper orthogonal decomposition, (POD) to obtain a low-dimensional  
470 space, and an adversarial autoencoder (AAE), to predict how the POD coefficients evolve in time. The  
471 input (and output) of the AAE consists of POD coefficients and sensor values, both over a sequence  
472 of  $m + 1$  time levels. Two methods of predicting with the AAE were presented: the “constrained”  
473 method involves repeatedly passing solutions through the AAE until convergence is reached; the second  
474 “backpropagation” method uses a backpropagation or adjoint algorithm to find which latent variables  
475 produce the first  $m$  values of the output. Once convergence is reached, the remaining values are taken  
476 to be the POD coefficients at the future time level. A school classroom is used as a test case here.  
477 A Computational Fluid Dynamics Large Eddy Simulation (CFD-LES) code is used to generate high-  
478 resolution solutions over the period of an hour, on which the surrogate model is based. The results were  
479 good for both the constrained method and the backpropagation method, although some deviation in the  
480 solution is seen towards the end of the one hour period for both surrogate models when compared with  
481 the CFD model. For the data assimilation, a type of dual-twin experiment was performed, in which the  
482 AAE was given POD coefficients from  $m$  time levels starting from time level 1 and marched forward in  
483 time for 120 time levels whilst assimilating observations from the test dataset at the end of the hour time  
484 period (time levels 577 to 720). The algorithm is able to assimilate the data and predict outputs that  
485 are very close to the CFD results in the test dataset (despite starting with POD coefficients from time  
486 level 1). After this successful demonstration of the data assimilation algorithm, observations that had

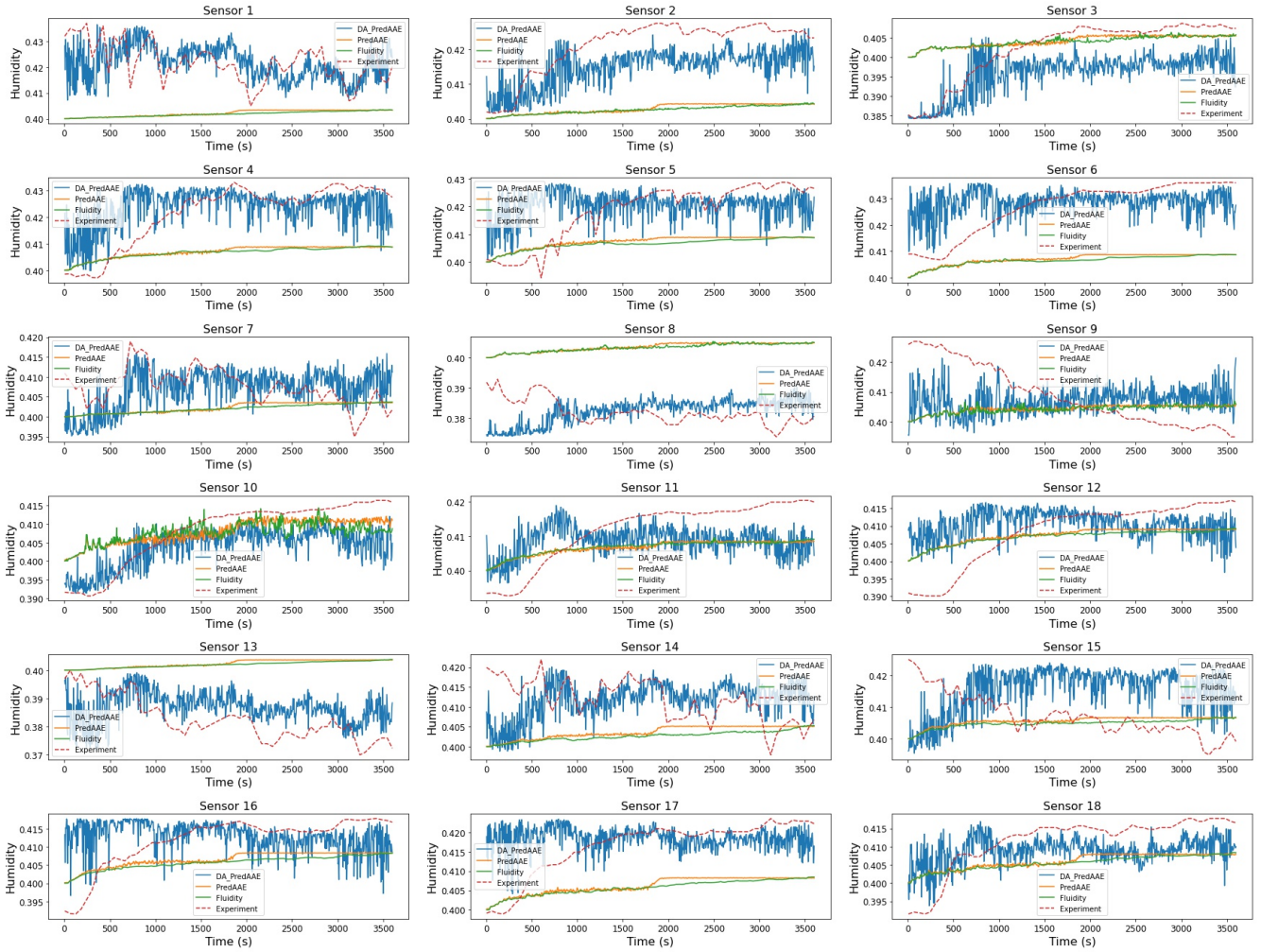




(a) CO<sub>2</sub> time series at the 18 sensors



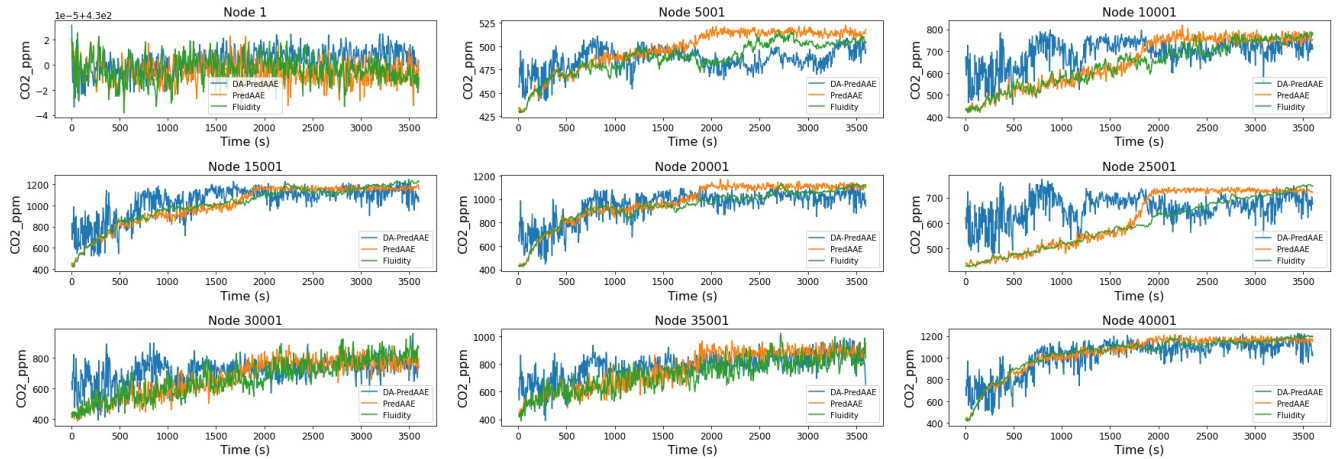
(b) Temperature time series at the 18 sensors



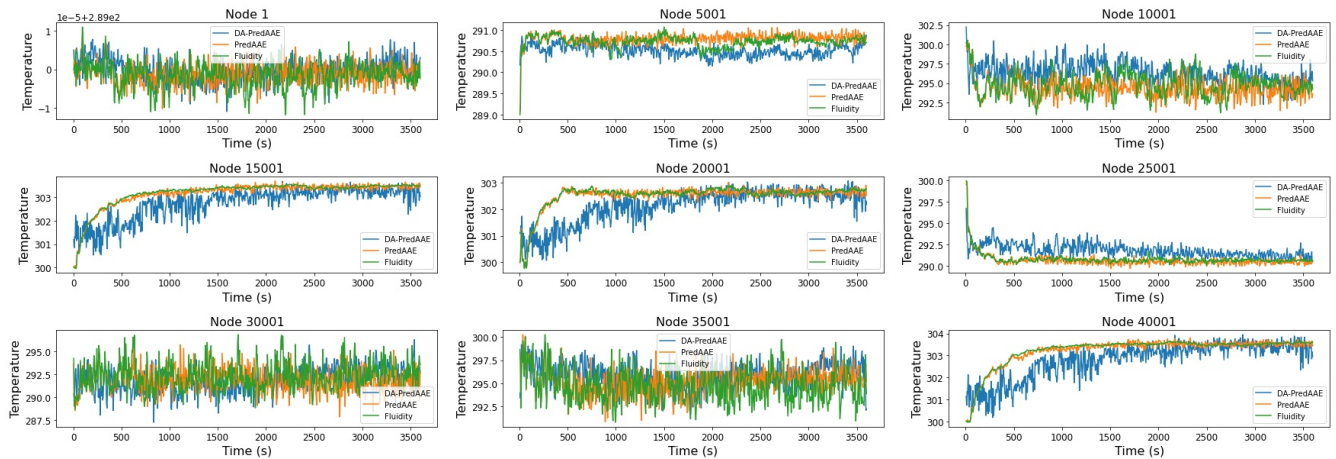
(c) Humidity time series at the 18 sensors

Figure 16: Assimilating data from the classroom in Houndsfield primary school. The above plots compare results from the high-fidelity model (green) with PredAAE predictions (orange), PredAAE prediction with data assimilation (blue) and the sensor observations (red). See Figure 6 for the location of the sensors.

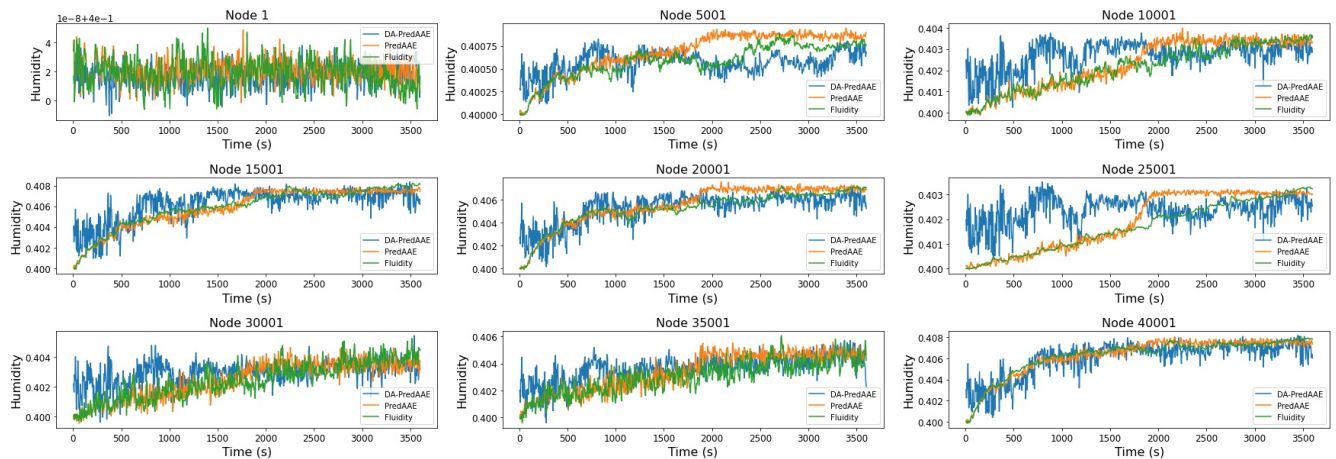




(a) CO<sub>2</sub> time series



(b) Temperature time series



(c) Humidity time series

Figure 17: Assimilating data from the classroom in Houndsfield primary school. The above plots compare results from the high-fidelity model (green) with PredAAE predictions (orange) and PredAAE prediction with data assimilation (blue). The time series are taken at locations of various nodes — nodes 1, 5001, 10001, 15001, 20001, 25001, 30001, 35001 and 40001 (see Figure 6).



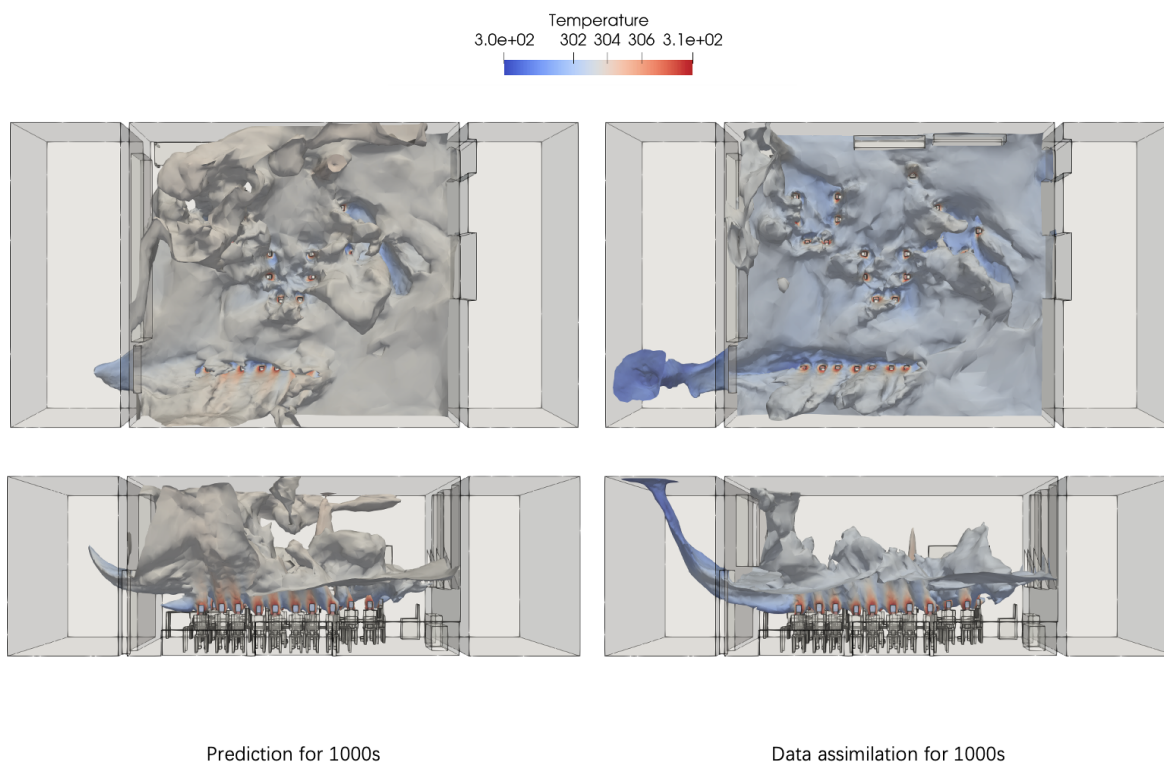


Figure 18: Assimilating data from the classroom in Houndsfield Primary School. Spatial variation of the iso-surface of CO<sub>2</sub> at 1000 ppm coloured by the temperature. The prediction without data assimilation is shown on the left at 1000s and the prediction with data assimilation is shown on the right also at 1000 s.

487 been collected from a classroom in Houndsfield Primary School were assimilated and results presented.

## 488 Acknowledgements

489 We would like to acknowledge the Applied Computational Science & Engineering MSc course and the  
 490 Environmental Data Science and Machine Learning MSc course at Imperial College London. Seven of  
 491 the co-authors were students on these courses and worked on projects related to the methods presented  
 492 in this article. We also thank the University of Surrey’s Global Centre for Clean Air Research (GCARE)  
 493 team members Dr Arvind Tiwari, Nidhi Rawat, Dr Sarkawt Hama, Dr Gopinath Kalaiarasan and Dr  
 494 Ana Paula Mendes Emygdio for collecting the classroom data [38]. We are grateful to Imperial College  
 495 Research Computing Service (<http://doi.org/10.14469/hpc/2232>) for providing computational resources.

496 The authors would like to acknowledge the following EPSRC grants: RELIANT, Risk EvaLuatIon fAst  
 497 iNtelligent Tool for COVID19 (EP/V036777/1); CO-TRACE, COvid-19 Transmission Risk Assessment  
 498 Case Studies — education Establishments (EP/W001411/1); INHALE, Health assessment across biolog-  
 499 ical length scales (EP/T003189/1); MAGIC, Managing Air for Green Inner Cities (EP/N010221/1); the  
 500 PREMIERE programme grant (EP/T000414/1); and MUFFINS (EP/P033180/1).

501 **CRedit authorship contribution statement**

502 **CEH:** methodology, software, writing (original draft, review and editing), supervision. **JT, JY:** software,  
 503 writing (review and editing), analysis, visualisation. **DG, JI, SK, YL, DS:** software, writing (review and  
 504 editing), analysis. **BC:** software, writing (original draft, review and editing), supervision. **LM:** writing  
 505 (review and editing), supervision, software. **PK:** data acquisition, writing (review and editing), funding  
 506 acquisition. **CCP:** conceptualisation, methodology, software, writing (review and editing), supervision,  
 507 funding acquisition.

508 **Appendix A. Hyperparameters of the Adversarial Autoencoder**

509 The hyperparameter settings used in the adversarial autoencoder are shown in Tables A.4 and A.5.

	optimal values	grid search
latent space	512	128, 256, 512, 1024
encoder layers	3	2, 3, 4, 5
decoder layers	2	2, 3, 4, 5
kernel size	(3, 3)	(2,2), (3,3), (5,5)
dropout rate	0.2	0, 0.2, 0.3
optimizer	Adam	Adam, Adamax, RMSProp, SGD
activation functions:		LeakyReLU, ReLU, sigmoid
hidden layers in encoder and decoder	LeakyReLU	
hidden layers in discriminator	ReLU	
output layers	sigmoid	
epochs	10,000	5,000, 10,000, 20,000
learning rate:		$10^{-3}$ , $10^{-4}$ , $4 \times 10^{-5}$ , $10^{-5}$ , $5 \times 10^{-5}$
autoencoder	$10^{-4}$	
encoder-discriminator	$5 \times 10^{-5}$	

Table A.4: Hyperparameter settings for the predictive AAE. The values over which a grid search was performed are on the right, the optimal values are on the left.

510 **References**

511 [1] P. Bauer, A. Thorpe, G. Brunet, The quiet revolution of numerical weather prediction, Nature 525  
 512 (2015) 47–55. doi:10.1038/nature14956.

513 [2] E. Kalnay, Atmospheric Modeling, Data Assimilation and Predictability, Cambridge University Press,  
 514 2003. doi:10.1017/CB09780511802270.

515 [3] D. G. Cacuci, I. M. Navon, M. Ionescu-Bujor, Computational methods for data evaluation and  
 516 assimilation, Chapman and Hall (CRC), 2016.

autoencoder		discriminator	
input	$9 \times 258$		
conv2D	$9 \times 258 \times 64$		
conv2D	$5 \times 129 \times 32$		
conv2D	$5 \times 65 \times 16$		
reshape	5200		
dense	512	→	input 512
dense	5200		dense 256
reshape	$5 \times 65 \times 16$		dense 128
conv2D	$5 \times 65 \times 32$		dense 64
conv2D	$5 \times 129 \times 64$		output 1
output	$9 \times 258$		

Table A.5: Architecture of the predictive AAE. The size of the 2D convolutional layers ('conv2D') written in the table gives the two dimensions of the feature maps, and the third number denotes the number of channels or filters used.

- 517 [4] M. Asch, M. Bocquet, M. Nodet, Data Assimilation: Methods, Algorithms and Applications, SIAM,  
518 2016.
- 519 [5] R. N. Bannister, A review of operational methods of variational and ensemble-variational data  
520 assimilation, Quarterly Journal of the Royal Meteorological Society 143 (2017) 607–633. doi:10.  
521 1002/qj.2982.
- 522 [6] R. Maulik, V. Rao, J. Wang, G. Mengaldo, E. Constantinescu, B. Lusch, P. Balaprakash,  
523 I. Foster, R. Kotamarthi, Efficient high-dimensional variational data assimilation with machine-  
524 learned reduced-order models, Geoscientific Model Development 15 (2022) 3433–3445. doi:10.5194/  
525 gmd-15-3433-2022.
- 526 [7] S. Hatfield, M. Chantry, P. Dueben, P. Lopez, A. Geer, T. Palmer, Building Tangent-Linear and  
527 Adjoint Models for Data Assimilation With Neural Networks, Journal of Advances in Modeling  
528 Earth Systems 13 (2021) e2021MS002521. doi:10.1029/2021MS002521.
- 529 [8] A. J. Geer, Learning earth system models from observations: machine learning or data assimilation,  
530 Philosophical Transactions of the Royal Society A 379 (2021) 20200089. doi:10.1098/rsta.2020.  
531 0089.
- 532 [9] M. Bocquet, J. Brajard, A. Carrassi, L. Bertino, Bayesian inference of chaotic dynamics by merging  
533 data assimilation, machine learning and expectation-maximization, Foundations of Data Science 2  
534 (2020) 55–80. doi:10.3934/fods.2020004.
- 535 [10] M. Sonnewald, R. Lguensat, D. C. Jones, P. D. Dueben, J. Brajard, V. Balaji, Bridging observations,  
536 theory and numerical simulation of the ocean using machine learning, Environmental Research Letters  
537 17 (2021). doi:10.1088/1748-9326/ac0eb0.

- 538 [11] C. Buizza, C. Quilodrán Casas, P. Nadler, J. Mack, S. Marrone, Z. Titus, C. Le Cornec, E. Heylen,  
539 T. Dur, L. Baca Ruiz, C. Heaney, J. A. Daz Lopez, K. S. Kumar, R. Arcucci, Data Learning:  
540 Integrating Data Assimilation and Machine Learning, *Journal of Computational Science* 58 (2022)  
541 101525. doi:10.1016/j.jocs.2021.101525.
- 542 [12] J. Brajard, A. Carrassi, M. Bocquet, L. Bertino, Combining data assimilation and machine learning  
543 to emulate a dynamical model from sparse and noisy observations: a case study with the Lorenz 96  
544 model, *Journal of Computational Science* 44 (202) 101171.
- 545 [13] A. Chennault, A. A. Popov, A. N. Subrahmanya, R. Cooper, A. Karpatne, A. Sandu, Adjoint-  
546 Matching Neural Network Surrogates for Fast 4D-Var Data Assimilation, arXiv preprint,  
547 arxiv:2111.08626 (2021). doi:10.48550/arxiv.2111.08626.
- 548 [14] C. Q. Casas, R. Arcucci, P. Wu, C. Pain, Y.-K. Guo, A Reduced Order Deep Data Assimilation  
549 model, *Physica D: Nonlinear Phenomena* 412 (2020) 132615. doi:10.1016/j.physd.2020.132615.
- 550 [15] M. E. Levine, A. M. Stuart, A framework for machine learning of model error in dynamical systems,  
551 arXiv preprint, arxiv:2107.06658 (2021). doi:10.48550/arXiv.2107.06658.
- 552 [16] J. Brajard, A. Carrassi, M. Bocquet, L. Bertino, Combining data assimilation and machine learning  
553 to infer unresolved scale parametrization, *Philosophical Transactions of the Royal Society A: Math-*  
554 *ematical, Physical and Engineering Sciences* 379 (2021) 20200086. doi:10.1098/rsta.2020.0086.
- 555 [17] A. Farchi, P. Laloyaux, M. Bonavita, M. Bocquet, Using machine learning to correct model error in  
556 data assimilation and forecast applications, *Quarterly Journal of the Royal Meteorological Society*  
557 147 (2021) 3067–3084. doi:10.1002/qj.4116.
- 558 [18] M. R. Ebers, K. M. Steele, J. N. Kutz, Discrepancy Modeling Framework: Learning missing physics,  
559 modeling systematic residuals, and disambiguating between deterministic and random effects, arXiv  
560 preprint, arXiv:2203.05164 (2022). doi:10.48550/arxiv.2203.05164.
- 561 [19] S. Pawar, O. San, Equation-free surrogate modeling of geophysical flows at the intersection of machine  
562 learning and data assimilation, arXiv preprint, arxiv:2205.13410 (2022). doi:10.48550/arxiv.2205.  
563 13410.
- 564 [20] M. Amendola, R. Arcucci, L. Mottet, C. Q. Casas, S. Fan, C. Pain, P. Linden, Y.-K. Guo, Data  
565 Assimilation in the Latent Space of a Convolutional Autoencoder, in: M. Paszynski, D. Kranzlmüller,  
566 V. V. Krzhizhanovskaya, J. J. Dongarra, P. M. Sliot (Eds.), *Computational Science – ICCS 2021,*  
567 *2021*, pp. 373–386.
- 568 [21] M. Peyron, A. Fillion, S. Gürol, V. Marchais, S. Gratton, P. Boudier, G. Goret, Latent space data  
569 assimilation by using deep learning, *Quarterly Journal of the Royal Meteorological Society* 147 (2021)  
570 3759–3777. doi:10.1002/qj.4153.
- 571 [22] R. S. Cintra, H. F. de Campos Velho, Data Assimilation by Artificial Neural Networks for an  
572 Atmospheric General Circulation Model, in: A. El-Shahat (Ed.), *Advanced Applications for Artificial*  
573 *Neural Networks*, IntechOpen, Rijeka, 2018. doi:10.5772/intechopen.70791.

- 574 [23] F. P. Härter, H. F. d. C. Velho, Multilayer Perceptron Neural Network in a Data Assimilation  
575 Scenario, *Engineering Applications of Computational Fluid Mechanics* 4 (2010) 237–245. doi:10.  
576 1080/19942060.2010.11015313.
- 577 [24] S. Pawar, O. San, A. Rasheed, I. M. Navon, A nonintrusive hybrid neural-physics modeling of  
578 incomplete dynamical systems: Lorenz equations, *International Journal on Geomathematics* 12  
579 (2021) 17. doi:10.1007/s13137-021-00185-z.
- 580 [25] M. H. Rammay, S. Alyaev, A. H. Elsheikh, Probabilistic model-error assessment of deep learning  
581 proxies: an application to real-time inversion of borehole electromagnetic measurements, *Geophysical*  
582 *Journal International* 230 (2022) 1800–1817. doi:10.1093/gji/ggac147.
- 583 [26] C. López, E. Hernández-García, Low-dimensional dynamical system model for observed coherent  
584 structures in ocean satellite data, *Physica A: Statistical Mechanics and its Applications* 328 (2003)  
585 233–250. doi:10.1016/S0378-4371(03)00505-3.
- 586 [27] Y. Cao, J. Zhu, I. M. Navon, Z. Luo, A reduced-order approach to four-dimensional variational data  
587 assimilation using proper orthogonal decomposition, *International Journal for Numerical Methods*  
588 *in Fluids* 53 (2007) 1571–1583. doi:10.1002/flid.1365.
- 589 [28] D. N. Daescu, I. M. Navon, Efficiency of a POD-based reduced second-order adjoint model in 4D-  
590 Var data assimilation, *International Journal for Numerical Methods in Fluids* 53 (2007) 985–1004.  
591 doi:10.1002/flid.1316.
- 592 [29] H. Gong, S. Cheng, Z. Chen, Q. Li, C. Quilodrn-Casas, D. Xiao, R. Arcucci, An efficient digital twin  
593 based on machine learning SVD autoencoder and generalised latent assimilation for nuclear reactor  
594 physics, *Annals of Nuclear Energy* 179 (2022) 109431. doi:10.1016/j.anucene.2022.109431.
- 595 [30] V. L. Silva, C. E. Heaney, Y. Li, C. C. Pain, Data Assimilation Predictive GAN (DA-PredGAN):  
596 applied to a spatio-temporal compartmental model in epidemiology, *Journal of Scientific Computing*  
597 94 (2021) 1–31. doi:10.1007/s10915-022-02078-1.
- 598 [31] F. Regazzoni, D. Chapelle, P. Moireau, Combining data assimilation and machine learning to  
599 build data-driven models for unknown long time dynamics—Applications in cardiovascular mod-  
600 eling, *International Journal for Numerical Methods in Biomedical Engineering* 37 (2021) e3471.  
601 doi:10.1002/cnm.3471.
- 602 [32] C. Quilodrán-Casas, R. Arcucci, L. Mottet, Y.-K. Guo, C. C. Pain, Adversarial autoencoders and ad-  
603 versarial LSTM for improved forecasts of urban air pollution simulations (2021). *arXiv:2104.06297*.
- 604 [33] R. Maulik, B. Lusch, P. Balaprakash, Non-autoregressive time-series methods for stable parametric  
605 reduced-order models, *Physics of Fluids* 32 (2020) 087115. doi:10.1063/5.0019884.
- 606 [34] S. Fresca, A. Manzoni, POD-DL-ROM: Enhancing deep learning-based reduced order models for  
607 nonlinear parametrized PDEs by proper orthogonal decomposition, *Computer Methods in Applied*  
608 *Mechanics and Engineering* 388 (2022) 114181. doi:10.1016/j.cma.2021.114181.

- 609 [35] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, B. Frey, Adversarial autoencoders, arXiv preprint,  
610 arxiv:1511.05644 (2015). doi:10.48550/arXiv.1511.05644.
- 611 [36] C. E. Heaney, Z. Wolffs, J. A. Tómasson, L. Kahouadji, P. Salinas, A. Nicolle, I. M. Navon, O. K.  
612 Matar, N. Srinil, C. C. Pain, An AI-based non-intrusive reduced-order model for extended domains  
613 applied to multiphase flow in pipes, *Physics of Fluids* 34 (2022) 055111. doi:10.1063/5.0088070.
- 614 [37] C. E. Heaney, X. Liu, H. Go, Z. Wolffs, P. Salinas, I. M. Navon, C. C. Pain, Extending the Capabilities  
615 of Data-Driven Reduced-Order Models to Make Predictions for Unseen Scenarios: Applied to Flow  
616 Around Buildings, *Frontiers in Physics* 10 (2022) 910381. doi:10.3389/fphy.2022.910381.
- 617 [38] P. Kumar, N. Rawat, A. Tiwari, Micro-characteristics of a naturally ventilated classroom air qual-  
618 ity under varying air purifier placements, *Environmental Research* (2022) 114849. doi:10.1016/j.  
619 envres.2022.114849.
- 620 [39] R. K. Bhagat, M. S. Davies Wykes, S. B. Dalziel, P. F. Linden, Effects of ventilation on the indoor  
621 spread of COVID-19, *Journal of Fluid Mechanics* 903 (2020) F1. doi:10.1017/jfm.2020.720.
- 622 [40] P. Dabisch, M. Schuit, A. Herzog, K. Beck, S. Wood, M. Krause, D. Miller, W. Weaver, D. Free-  
623 burger, I. Hooper, B. Green, G. Williams, B. Holland, J. Bohannon, V. Wahl, J. Yolitz, M. Hevey,  
624 S. Ratnesar-Shumate, The influence of temperature, humidity, and simulated sunlight on the  
625 infectivity of SARS-CoV-2 in aerosols, *Aerosol Science and Technology* 55 (2021) 142–153.  
626 doi:10.1080/02786826.2020.1829536.
- 627 [41] P. Holmes, J. Lumley, G. Berkooz, C. Rowley, *Turbulence, Coherent Structures, Dynamical Systems*  
628 *and Symmetry*, Cambridge University Press, 2012.
- 629 [42] Applied Modelling and Computation Group, Imperial College London, Fluidity repository, <https://fluidityproject.github.io/>, 2016. Accessed 2 October 2022.
- 630 [43] Fluidity manual, 2016. <https://fluidityproject.github.io/>, accessed 2 October 2022.
- 631 [44] C. Pain, A. Umpleby, C. de Oliveira, A. Goddard, Tetrahedral mesh optimisation and adaptivity for  
632 steady-state and transient finite element calculations, *Computer Methods in Applied Mechanics and*  
633 *Engineering* 190 (2001) 3771–3796. doi:10.1016/S0045-7825(00)00294-2.
- 634 [45] D. Pavlidis, G. J. Gorman, J. L. M. A. Gomes, C. C. Pain, H. ApSimon, Synthetic-Eddy Method  
635 for Urban Atmospheric Flow Modelling, *Boundary-Layer Meteorology* 136 (2010) 285–299. doi:10.  
636 1007/s10546-010-9508-x.
- 637 [46] H. Woodward, M. Stettler, D. Pavlidis, E. Aristodemou, H. ApSimon, C. Pain, A large eddy sim-  
638 ulation of the dispersion of traffic emissions by moving vehicles at an intersection, *Atmospheric*  
639 *Environment* 215 (2019) 116891. doi:10.1016/j.atmosenv.2019.116891.
- 640 [47] L. Mottet, J. Song, A. C. Short, S. Chen, J. Wu, W. Yu, J. Xiong, Q. Zhang, J. Ge, M. Liu, R. Yao,  
641 B. Li, The hot summer-cold winter region in China: Challenges in the low carbon adaptation of  
642 residential slab buildings to enhance comfort, *Energy and Buildings* 223 (2020) 110181. doi:10.1016/  
643 j.enbuild.2020.110181.
- 644

- 645 [48] L. Mottet, Indoor Geometry Generator Manual, [https://www.researchgate.net/publication/](https://www.researchgate.net/publication/349350000_Indoor_Geometry_Generator_IGG_Manual)  
646 [349350000\\_Indoor\\_Geometry\\_Generator\\_IGG\\_Manual](https://www.researchgate.net/publication/349350000_Indoor_Geometry_Generator_IGG_Manual), 2021. Accessed 2 October 2022.
- 647 [49] C. E. Heaney, Y. Li, O. K. Matar, C. C. Pain, Applying Convolutional Neural Networks to Data  
648 on Unstructured Meshes with Space-Filling Curves, arXiv preprint, arxiv:2011.14820 (2020). doi:10.  
649 48550/arxiv.2011.14820.
- 650 [50] Y. Zhou, C. Wu, Z. Li, C. Cao, Y. Ye, J. Saragih, H. Li, Y. Sheikh, Fully convolutional  
651 mesh autoencoder using efficient spatially varying kernels, in: H. Larochelle, M. Ranzato,  
652 R. Hadsell, M. F. Balcan, H. Lin (Eds.), Advances in Neural Information Processing Systems,  
653 volume 33, 2020, pp. 9251–9262. URL: [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/68dd09b9ff11f0df5624a690fe0f6729-Paper.pdf)  
654 [68dd09b9ff11f0df5624a690fe0f6729-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/68dd09b9ff11f0df5624a690fe0f6729-Paper.pdf).
- 655 [51] J. Tencer, K. Potter, A Tailored Convolutional Neural Network for Nonlinear Manifold Learning of  
656 Computational Physics Data Using Unstructured Spatial Discretizations, SIAM Journal on Scientific  
657 Computing 43 (2021) A2581–A2613. doi:10.1137/20M1344263.
- 658 [52] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean,  
659 M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser,  
660 M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens,  
661 B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals,  
662 P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale machine learning on  
663 heterogeneous systems, 2015. doi:10.5281/zenodo.4724125, software available from tensorflow.org.