# Data Assimilation in the Latent Space of a Neural Network

Maddalena Amendola[a], Rossella Arcucci[a,*], Laetitia Mottet[c], César Quilodrán Casas[a], Shiwei Fan[b], Christopher Pain[a,c], Paul Linden[d], Yi-Ke Guo[a]

[a]*Data Science Institute, Department of Computing, Imperial College London, UK*

[b]*Department of Chemistry, University of Cambridge, UK*

[c]*Department of Earth Science & Engineering, Imperial College London, UK*

[d]*Department of Applied Mathematics and Theoretical Physics, University of Cambridge, UK*

## Abstract

There is an urgent need to build models to tackle Indoor Air Quality issue. Since the model should be accurate and fast, Reduced Order Modelling technique is used to reduce the dimensionality of the problem. The accuracy of the model, that represent a dynamic system, is improved integrating real data coming from sensors using Data Assimilation techniques. In this paper, we formulate a new methodology called Latent Assimilation that combines Data Assimilation and Machine Learning. We use a Convolutional neural network to reduce the dimensionality of the problem, a Long-Short-Term-Memory to build a surrogate model of the dynamic system and an Optimal Interpolated Kalman Filter to incorporate real data. Experimental results are provided for $CO_2$ concentration within an indoor space. This methodology can be used for example to predict in real-time the load of virus, such as the SARS-COV-2, in the air by linking it to the concentration of $CO_2$.

*Keywords:* Data Assimilation, Autoencoder, LSTM, Reduced Order Modelling, Indoor Air Pollution, Latent space

---

*Corresponding author

*Email address:* r.arcucci@imperial.ac.uk (Rossella Arcucci)

## 1. Introduction

Urbanisation is the process for which people move from rural zone to urban zone changing their habits. This process grows year by year: about half of the global population already lives in urban areas and by 2050 two-thirds of the world's people are expected to live in urban areas. Urbanisation process has led to an increase in building, human activities and energy consumption causing environmental degradation. High building densities and the low presence of vegetation impair the air quality and circulation. People who live in such areas are hesitant to open the windows of their house thinking that this can led to an increment of pollution in their habitation. A solution from their point of view is to use air-conditioning increasing in this way the energy consumption. This is a vicious cycle of increased urban emissions of heat, pollutants and greenhouse gases and an associated increase in energy demand.

The scope of the MAGIC[1] (Managing Air for Green Inner Cities) project is to study and build systems to assist reduction in energy demand through natural ventilation [1]. To this aim, there is the need to use systems with high accuracy in predicting air flows and air pollution concentration. These systems use the Large Eddy Simulation method within the Computational Fluids Dynamics (CFD) software: Fluidity [2]. Fluidity is an open source, general purpose, multi-phase computational fluid dynamics code capable of numerically solving the Navier-Stokes equations and advection-diffusion equations on arbitrary unstructured finite-element meshes. Fluidity is used in a number of different scientific areas including geophysical fluid dynamics, ocean modelling, mantle convection and air pollution.

Numerical simulation has been widely applied in many fields including environmental sciences, aerospace engineering, bio-medicine and industrial design. It provides powerful technical support for solving industrial problems and making scientific research in these fields. However, high fidelity numerical simulations of complex systems consume vast time and computing resources. When real data collected by instruments (i.e. sensors) are available, it is possible to use them to improve the accuracy of the prediction.

---

[1]http://www.magic-air.uk/home.html

The integration is made up by Data Assimilation techniques.

Data Assimilation (DA) is an approach for fusing data (observations) with prior knowledge (e.g., mathematical representations of physical laws; model output) to obtain an estimate of the distribution of the true state of a process [3]. In order to perform DA, one needs observations (i.e., a data or measurement model), a background (i.e., a priori state or process model) and information about the distribution of the errors on these two. For those applications, where the background is defined in big computational grids which lead to a big data problem sometimes impossible to handle without introducing approximations or space reductions, Reduced Order Modelling (ROM) techniques are used [4, 5].

ROM allows to speed up the dynamic model and the DA process. Popular approaches to reduce the domain are the Principal Component Analysis (PCA) and the Empirical Orthogonal Functions (EOF) technique both based on a Truncated Singular Value Decomposition (TSVD) analysis [6]. The simplicity and the analytic derivation of those approaches are the main reasons behind their popularity in atmospheric and ocean science. However, despite those powerful approaches, the accuracy of the obtained solution exhibits a severe sensibility to the variation of the value of the truncation parameters. This issue introduces a severe drawback to the reliability of these approaches, hence their usability in operative software in different scenarios [7].

An approach to reduce the dimensionality maintaining information of the data is the Neural Network (NN), precisely the AutoEncoders [8, 9]. NNs have the ability to fit functions and they can fit almost any unknown function theoretically. That is the ability which makes it possible for neural networks to face complex problems. AutoEncoders with non-linear encoder functions and non-linear decoder functions can thus learn a more powerful non-linear generalisation of methods based on TSVD. In the latent space, the evolution of the transformed state variables defined in time, can be learned using Recurrent Neural Networks (RNN) [10, 11]. In the present work, we propose a new methodology which we called Latent Assimilation (LA). It consists in reducing the dimensionality with NN and perform both prediction through a surrogate dynamic model and DA directly in the latent space. In the latent space, the surrogate dynamic system is built by a RNN.

## 2. Related Work and contribution of the present work

The future challenges of Numerical Weather Prediction (NWP) is to include more accurate initial conditions that take advantage of the increasing volume of real-time observations, and improve the post-processing of model outputs, amongst others [12]. To answer this need, Neural network (NN) for correction of error in forecasting have been extensively studied [13, 14, 15]. However, the error correction by NN does not have a direct relation with the updated model system at each step and the training is not on the results of the assimilation process.

A framework for integration of NN with physical models by Data Assimilation (DA) algorithms is described in [16]: the NNs are iteratively trained when observed data are updated. Unfortunately, this approach presents a limit due to the time complexity of the numerical models involved, which limits the use of the forecast model for large data problems. An approach for employing artificial neural networks (NNs) to emulate the Local Ensemble Transform Kalman Filter (LETKF) as a method of data assimilation is presented in [17]. Deep learning and Data Assimilation technologies are also combined to predict the production of gas from mature gas wells in [18]. The authors used a modified deep Long Short-Term Memory (LSTM) model as their prediction model in the Ensemble KF framework for parameter estimation. A Neural Network is integrated into a conventional DA in [16]: deep learning shows great advantage in function approximations which have unknown model and strong non-linearity. The authors used NNs to characterise the structural model uncertainty. The NN is implemented in an End-to-End (E2E) approach and its parameters are iteratively updated with coming observations by applying the DA method.

A framework which performs fast data assimilation with sufficient accuracy for open ocean is proposed in [19]. Speed improvement is achieved by performing the data assimilation on a reduced-space rather than on a full-space. A dimension reduction of the full-state is made by an Empirical Orthogonal Function (EOF) analysis while retaining most of the explained variance. Analysis of EOFs can be used to identify structures in geophysical data which hold a large part of the variance. In this framework, the assimilation is performed in the control space. EOFs analysis has become a fundamental tool in atmosphere, ocean, and climate science for data diagnostics and dynamical mode reduction. Each of these applications exploits the fact that EOFs

allow a decomposition of a data function into a set of orthogonal functions, which are designed so that only a few of these functions are needed in lower-dimensional approximations. Furthermore, since EOFs are the eigenvectors of the error co-variance matrix, its condition number is reduced as well. Nevertheless, the accuracy of the solution obtained by truncating EOFs exhibits a severe sensibility to the variation of the value of the truncation parameter, so that a suitably choice of the number of EOFs is strongly recommended. This issue introduces a severe drawback to the reliability of EOFs truncation, hence to the usability of the operative software in different scenarios. A powerful solution to this is to use a Tikhonov regularisation which reveals to be more appropriate than truncation of EOFs [4].

Neural networks have tremendous ability to fit functions and they can fit almost any unknown function theoretically. That is the ability which makes it possible for neural networks to model complex flows. The complex computations involving matrices is reduced by factorising the representation deriving a latent state used from the Kalman Filter in [20]. The authors also used a linear dynamic model to compute, i.e predict, the next timestep. A variational AutoEncoder capable to generate trajectories from a latent space where the dynamics is linear is presented in [21].

In this paper, we propose a new methodology that use the NNs to reduce the space and perform the assimilation of the sensors data in the latent space. Specifically, we use a Convolutional AutoEncoder to reduce the domain and we perform an Optimal Interpolated Kalman Filter in the latent space.

In this paper, we make the following contributions:

- We have designed a novel data assimilation technology, we called Latent Assimilation (LA), mainly composed by an AutoEncoder, a surrogate model and an Optimal Kalman Filter. The Latent Assimilation model performs the prediction of the flows and the assimilation of observed data through a Kalman Filter in the latent space.

- We have developed a Convolutional AutoEncoder to reduce the space where the surrogate model will work and where we perform the assimilation of the observation using the Optimal Interpolated Kalman Filter. We have chosen to use an encoder-decoder model instead of Principal Component Analysis (PCA) since neural networks maintain non-linearities and perform better in modelling flows;

5

- We have built a Recurrent Neural Network (LSTM) to emulate a Computational Fluid Dynamics (CFD) simulation in the latent space of an AutoEncoder: the trained LSTM represents the surrogate model to predict the $CO_2$ concentration in a room;

- We prove that our novel Latent Assimilation model answers the needs of accuracy, stability and efficiency required by real-time applications.

- We have developed a software written in python to test the Latent Assimilation model. The LA code and the pre-processed data can be downloaded using the link:

  `https://github.com/DL-WG/LatentAssimilation`.

Experimental results are provided for pollutant dispersion within an indoor space. This methodology can be used for example to predict in real-time the load of virus, such as the SARS-COV-2, in indoor spaces by linking it to the concentration of $CO_2$ [22].

## 3. Data Assimilation

In this section, we introduce the concept of Data Assimilation (DA) and the Kalman Filter (KF) which is one of the most used approach for DA.

DA merges the estimated state $x_t \in \mathbb{R}^n$ of a discrete-time dynamic process at time $t$:

$$x_{t+1} = M_{t+1}x_t + w_t \tag{1}$$

with an observation $y_t \in \mathbb{R}^m$:

$$y_t = H_t x_t + v_t \tag{2}$$

where $M_{t+1}$ is a dynamic linear operator and $H_t$ is the observation operator. The vectors $w_t$ and $v_t$ represent the process and observation errors, respectively. They are usually assumed to be independent, white-noise processes with Gaussian probability distributions:

$$w_t \sim \mathcal{N}(0, Q_t), \quad v_t \sim \mathcal{N}(0, R_t)$$

where $Q_t$ and $R_t$ are called errors covariance matrices of the model and the observations, respectively.

6

DA tries to answer questions such as "what can be said about the value of an unknown variable $x_t$ that represents the evolution of a system, if we have some measured data $y_t$ and a model $M$ of the underlying mechanism that generated the data?". This is the Bayesian context, where we seek a quantification of the uncertainty in our knowledge of the parameters that, according to Bayes' rule takes the form

$$p\left(x_t|y_t\right) = \frac{p\left(y_t|x_t\right)p\left(x_t\right)}{p\left(y_t\right)} \tag{3}$$

Here, the physical model is represented by the conditional probability (also known as the likelihood) $p\left(y_t|x_t\right)$, and the prior knowledge of the system by the term $p\left(x_t\right)$. The denominator is considered as a normalising factor and represents the total probability of $y_t$. DA is a Bayesian inference that combines the state $x_t$ with $y_t$ at each given time. The Bayes theorem conducts to the estimation of $x_t^a$ which maximise a probability density function given the observation $y_t$ and a prior from $x_t$. This approach is implemented in one of the most popular DA methods which is the Kalman Filter (KF) [23] which mainly consists of two steps: a prediction (equation (4)) and a correction (equations (5)-(6)) steps. The goal of the KF is to compute an optimal *a posteriori* estimate, $x_t^a$, which is a linear combination of an *a priori* estimate, $x_t$, and a weighted difference between the actual measurement, $y_t$, and the measurement prediction, $H_t x_t$ as described in equation (6).

1. **Prediction**:
$$x_{t+1} = M_{t+1} x_t^a \tag{4}$$

2. **Correction**:

$$K_{t+1} = Q_{t+1} H_{t+1}^T (H_{t+1} Q_{t+1} H_{t+1}^T + R_{t+1})^{-1} \tag{5}$$

$$x_{t+1}^a = x_{t+1} + K_{t+1}(y_{t+1} - H x_{t+1}) \tag{6}$$

For big data problems, KF is usually implemented in a simplified version as an Optimal Interpolation method [24] for which the covariance matrix $Q_t = Q$ is fixed at each timestep $t$.

The prediction-correction cycle is complex and time-consuming and it mandates the introduction of simplifications, approximations or data reductions techniques. In the next section, we present the Latent Assimilation approach which consists in performing KF in the latent space of an Autoencoder with

7

nonlinear encoder and nonlinear decoder functions. In the latent space, the dynamic system in equation (4) is replaced by a surrogate model built with a RNN.

## 4. Latent Assimilation

Latent Assimilation is a model that implements the idea of assimilating real data in the Latent Space of a Neural Network (NN). Instead of using PCA or others mathematical approaches to reduce the space, we model the reduction with non-linear transformations using Deep NNs. Specifically, we choose to use Convolutional Autoencoder to reduce the space. The model is divided into four main parts:

1. Dimensionality reduction: the physical space is transformed in a latent space of smaller dimension by a Convolutional Autoencoder;
2. Surrogate model: a surrogate of the CFD is built in the latent space by a Recurrent Neural Network;
3. Data Assimilation: observed data are assimilate in the surrogate of the CFD by a Kalman Filter;
4. Physical space: the results of the DA in the latent space are then reported in the physical space through a Decoder.

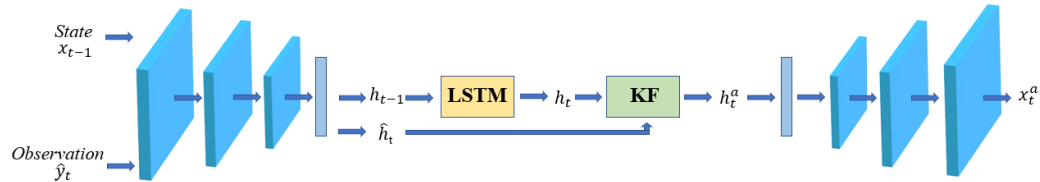Figure 1 shows the work flows of the Latent Assimilation model.



Figure 1: Latent Assimilation model workflow. Let assume that we want to predict the state of the system at time $t$ and we assume that the LSTM needs one observation back to predict the next timestep. The input of the system is the state $x_{t-1}$. We encode $x_{t-1}$ producing its encoded version $h_{t-1}$. From $h_{t-1}$ we compute $h_t$ through LSTM. To perform the Kalman Filter, we need the observation $\hat{y}_t$ at timestep $t$. We encode $y_t$ and we combine the result, $\hat{h}_t$, with the prediction $h_t$ through the KF. The result $h_t^a$ is the updated prediction. We report the updated prediction in its physical space through the Decoder, producing $x_t^a$.

*4.1. Dimensionality reduction*

The dimensionality reduction is implemented by an AutoEncoder (AE). AEs are usually used for dimensionality reduction or feature learning. To use the autoencoder for dimensionality reduction, the encoder function must returns an output with lower dimension with respect to the input. This kind of autoencoder are called *undercomplete*. Learning an undercomplete representation forces the autoencoder to capture the most salient features of the training data. One type of undercomplete autoencoder is the Convolutional autoencoder. As we can deduce from the name, this autoencoder uses the Convolutional operation. Thanks to the convolutional operation, the network takes into account the spatial information: they are specially used with images or grid data. Usually, the Convolutional AutoEncoders are composed by more than one convolutional layers, each followed by pooling layer to reduce the input [25]. Latent Assimilation implements a Convolutional Autoencoder which produces a representation of the state vector $x_t \in \mathbb{R}^n$ in (1) in a "latent" state vector $h_t \in \mathbb{R}^p$ defined in a Latent Space where $p < n$. We denote with $f : \mathbb{R}^n \to \mathbb{R}^p$ the Encoder function

$$h_t = f(x_t) \tag{7}$$

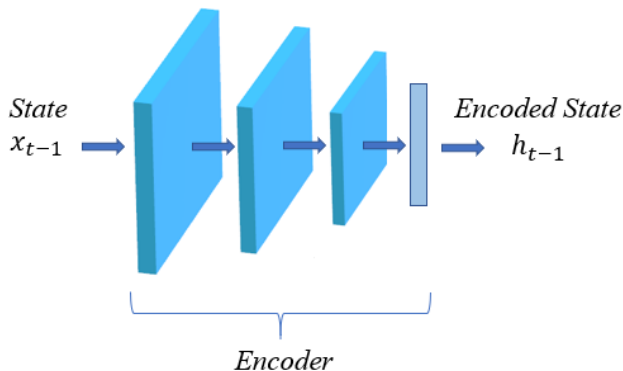which transforms the state $x_t$ in a latent variable $h_t$.



Figure 2: Latent Assimilation - Encoder: $x_{t-1}$ is the input state and $h_{t-1}$ is the corresponding encoded state at timestep $t-1$.

*4.2. Surrogate model*

In the latent space we perform a regression through a Long Short Term Memory (LSTM) function $l : \mathbb{R}^{p \times q} \to \mathbb{R}^p$

$$h_{t+1} = l(\boldsymbol{h_{t,q}}) \tag{8}$$

where $\boldsymbol{h_{t,q}} = \{h_i\}_{i=t,\ldots,t-q}$ is a sequence of $q$ encoded timesteps up to time $t$. The LSTM is a Recurrent Neural Network (RNN) with good performance with time-series data [26]. It is composed by gates and cells as shown if Figure 3. The gates decide which information should pass using a sigmoid function. The LSTM is composed by four elements described below. In all formulas, $b$, $U$ and $W$ denote respectively the biases, input weights and recurrent weights for corresponding gate.

1. **Forget gate**: it decides what information should pass via a logistic sigmoid $\sigma(x) = (1/(1 + e^{-x}))$ unit

$$d_i^{(t)} = \sigma(b_i^d + \sum_j U_{i,j}^d h_j^{(t)} + \sum_j W_{i,j}^d u_j^{(t-1)}) \tag{9}$$

   where $h^{(t)}$ is the input and $u^{(t)}$ is the hidden layer.

2. **Input Gate**: it is similar to the Forget Gate but with its own parameters

$$c_i^{(t)} = \sigma(b_i^c + \sum_j U_{i,j}^c h_j^{(t)} + \sum_j W_{i,j}^c u_j^{(t-1)}) \tag{10}$$

3. **Cell State**: the cell state is then updated using sigmoid and hyperbolic tangent functions

$$s_i^i = d_i^{(t)} s_i^{(t-1)} + c_i^{(t)} \sigma(b_i + \sum_j U_{i,j} h_j^{(t)} + \sum_j W_{i,j} u_j^{(t-1)}) \tag{11}$$

4. **Output Gate**: it decides what the value of the next hidden state using a sigmoid function on the previous hidden state and the current input and an hyperbolic tangent function on the newly modified cell state.

$$q_i^{(t)} = \sigma(b_i^o + \sum_j U_{i,j}^o h_j^{(t)} + \sum_j W_{i,j}^o u_j^{(t-1)}) \tag{12}$$

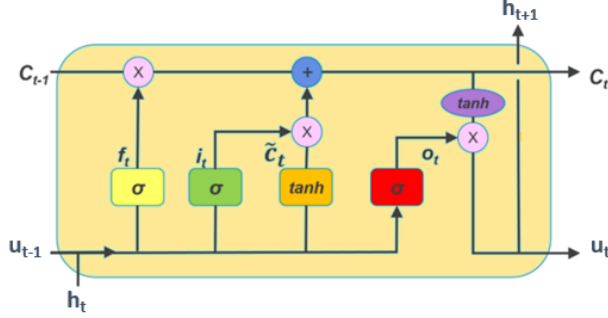$$h_i^{(t+1)} = \tanh(s_i^{(t)}) q_i^{(t)} \tag{13}$$

10

Figure 3: Latent Assimilation - LSTM.

### 4.3. Data Assimilation

The assimilation is performed in the latent space. In order to merge the observations in (2) with the "latent" state vector $h_t$, the observations are processed by the Encoder in the same way as the state vector. As $y_t \in \mathbb{R}^m$ where usually $m \leq n$, i.e. the observations are usually held or measured in just few point in space, the observations vector $y_t$ is interpolated in the state space $\mathbb{R}^n$ obtaining $\hat{y}_t \in \mathbb{R}^n$. The observations $\hat{y}_t$ are then processed in the same way as the state vector trough $f$:

$$\hat{h}_t = f(\hat{y}_t) \tag{14}$$

The "latent" observations $\hat{h}_t$, transformed by the Encoder in the latent space, are then assimilated by the prediction-correction steps as described in equations (15)-(17) ad as shown in Figure 4:

1. **Prediction**:
$$h_{t+1} = l(\boldsymbol{h_{t,q}}) \tag{15}$$

2. **Correction**:
$$\hat{K}_{t+1} = \hat{Q}\hat{H}^T(\hat{H}\hat{Q}\hat{H}^T + \hat{R}_{t+1})^{-1} \tag{16}$$
$$h_{t+1}^a = h_{t+1} + \hat{K}_{t+1}(\hat{h}_{t+1} - \hat{H}h_{t+1}) \tag{17}$$

where $l$ in (15) is the surrogate model defined in (8) computed by the LSTM, $\hat{Q}$ and $\hat{R}$ are the errors covariance matrices of the transformed background $h_t$ and observations $\hat{h}_t$ respectively: they are computed directly in the latent space. The background covariance matrix $\hat{Q}$ is computed with a sample of $s$

model state forecasts $\boldsymbol{h}$ that we set aside as background such that:

$$\boldsymbol{h} = [h_1, ..., h_s] \in \mathbb{R}^{p \times s}, \quad V = (\boldsymbol{h} - \bar{h}) \in \mathbb{R}^{p \times s} \tag{18}$$

where $\bar{h}$ is the mean of the sample of background states, then $\hat{Q} = VV^T$. The observations errors covariance matrix $\hat{R}$ can be computed with the same process than in equation (18) by replacing $h_t$ with $\hat{h}_t$ $\forall t$

$$\hat{\boldsymbol{h}} = [\hat{h}_1, ..., \hat{h}_s] \in \mathbb{R}^{p \times s}, \quad \hat{V} = (\hat{\boldsymbol{h}} - \bar{\hat{h}}) \in \mathbb{R}^{p \times s} \tag{19}$$

where $\bar{\hat{h}}$ is the mean of the sample of observations, then $\hat{R} = \hat{V}\hat{V}^T$. The covariance matrix $\hat{R}$ can be estimated by evaluations of measurements (instruments) errors

$$\hat{R} = \sigma I \in \mathbb{R}^{p \times p} \tag{20}$$

where $0 < \sigma < 1$ and $I \in \mathbb{R}^{p \times p}$ denotes the identity matrix [24]. $\hat{K}$ is the Kalman Gain matrix defined in the latent space and $\hat{H}$ is the observation operator.
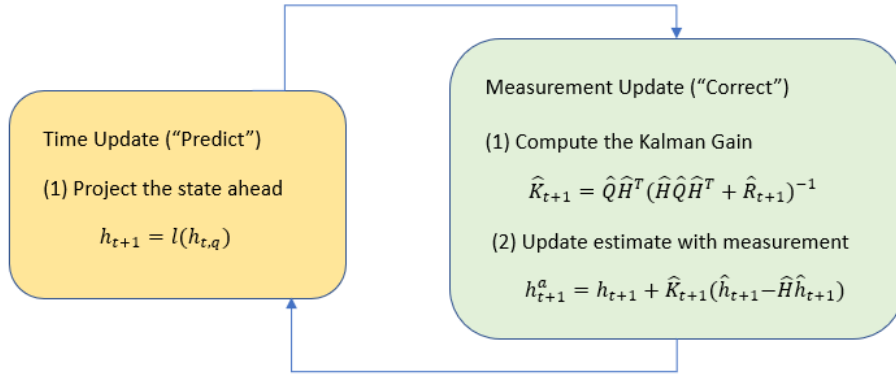


Figure 4: Latent Assimilation - Kalman Filter loop: the yellow block is the prediction phase (LSTM) and the green block is the correction phase (KF).

*4.4. Physical space*

The results of the DA in the latent space are then reported in the physical space through the Decoder, applying the function $g : \mathbb{R}^p \to \mathbb{R}^n$ to compute

$$x^a_{t+1} = g(h^a_{t+1}). \tag{21}$$

The Decoder is almost a mirror of the Encoder: it is composed of a Fully Connected Layer followed by some Convolutional Layers.
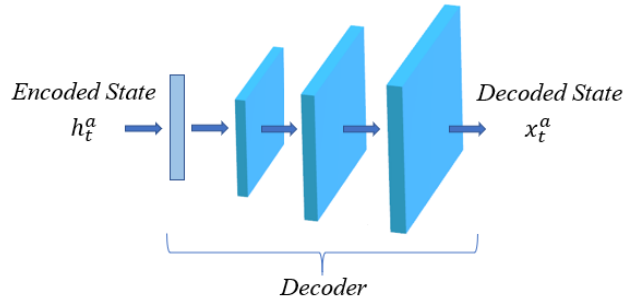


Figure 5: Latent Assimilation - Decoder: $h_t^a$ is the encoded updated state and $x_t^a$ is the decoded updated state.

## 4.5. Latent Assimilation code

The code is written in Python and is available at the following link: `https://github.com/DL-WG/LatentAssimilation`. The LatentAssimilation folder is composed by different subfolders:

- DataSet: it contains the Structured dataset divided in train and test;

- PreProcess: it contains all the code written to extract the Structured dataset starting from the unstructured meshes. We used the python libraries math, numpy, vtktools and pyvista.

- AutoEncoder and LSTM: both folders contain the code used to find the structure of the model and the hyper-parameters for the Structured dataset. All results are stored and also visualized in AnalysisLS7 jupyter notebook. We used python libraries such as numpy, sklearn, pandas and tensorflow.

- Data Assimilation: here there is all the observation data preprocessing, the Kalman Filter and the LatentAssimilation module which performs the assimilation in the Latent Space and it prints the table of the results.

In the next Section we apply Latent Assimilation to the problem of assimilating data to improve the prediction of air flows and indoor pollution transport in a real scenario [1]. We show the performance of the model step by step and we compare results with a standard DA performed in the physical space.
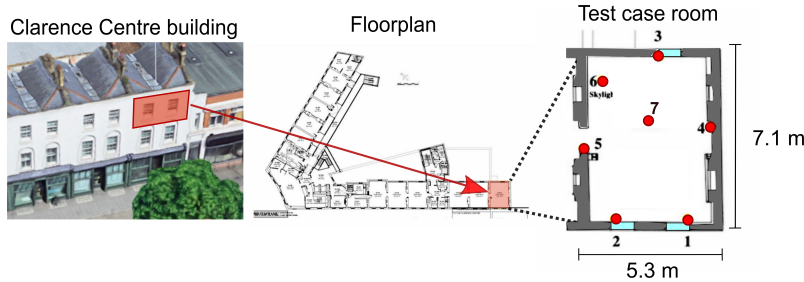
13

Figure 6: Test case room located in Clarence Centre building, London, UK. Red dots denote the location of the 7 sensors used during the field experiment. Blue rectangles show the location of the three windows [1].

## 5. Set up of a real test case

### 5.1. CFD simulation

The LA model presented in Section 4 is applied to real data collected in the context of the MAGIC project [1]: external and internal air quality measurements were performed within a naturally ventilated office room located at the top floor of the three-storey Clarence Centre building, Borough of Southwark, London, UK (Figure 6). The room has two windows facing a busy road (London Road), one window facing a traffic-free courtyard and a skylight in the ceiling. Seven sensors located in different positions were used to record, amongst others, the indoor temperature and $CO_2$ concentration, with a sampling rate of 1 minute. The three windows were opened during 25 minutes to look at cross ventilation effect on the decay of temperature and $CO_2$ concentration. During the whole period of the experiment, the predominant wind was a south-westerly wind.

To replicate the field study experiment, a numerical simulation has been performed using the Computational Fluid Dynamics (CFD) software Fluidity (`http://fluidityproject.github.io/`). The same CFD simulation has been used in a previous paper [27] and only the main details of the CFD setup are re-called here. The computational domain includes the Clarence Centre building as well as the immediate upwind building, and the test room office as shown in Figure 7 in order to replicate the cross ventilation scenario done during the field study. The mesh generated is an unstructured tetrahedral mesh composed by 400,000 nodes (Figure 7). The initial and boundary conditions are set to replicate the experimental conditions and are derived
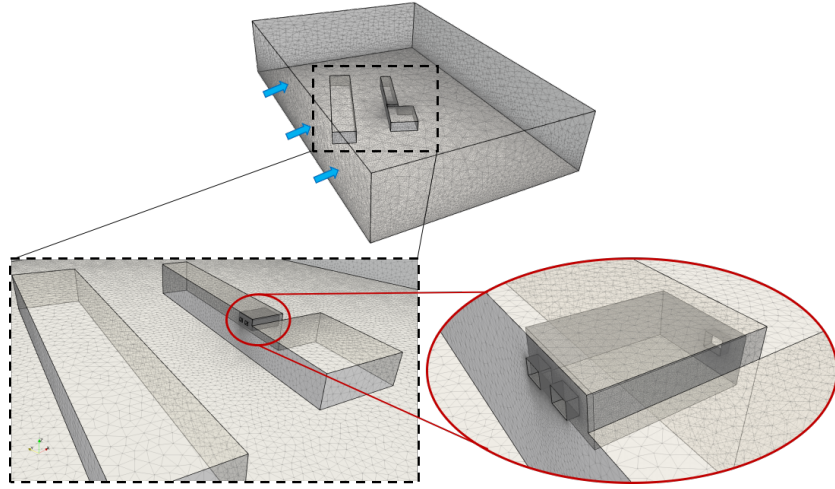
14

Figure 7: Computational domain and surface mesh of the area of interest showing the Clarence Centre and the upwind building as well as the test case room. The blue arrows denote the wind direction.

from the indoor sensors and the weather station used during the field study. The initial indoor $CO_2$ concentration is set equal to 1420 ppm, while the outdoor background $CO_2$ level is set equal to 400 ppm. The initial indoor and outdoor temperatures are equal to 19.5 $^{o}C$ and 9.1 $^{o}C$, respectively. The inlet velocity is following a log-law profile reaching 2.58 m/s at 28.5 m height. The simulation was run in parallel on 20 CPU and 15 minutes were simulated rendering approximately 3,500 timesteps. In this paper, the working variable of interest is the $CO_2$ concentration. It is worth noting that after the timestep 2,500 the concentration of $CO_2$ is low everywhere in the room since the room is completely ventilated.

### 5.2. CFD data pre-processing

The data generated by the CFD simulation are stored on an unstructured mesh and need to be converted into structured data in order to apply the LA model presented in Section 4. Indeed, convolutional kernels work on the assumption that adjacent states are equally spaced. As a first step, only the nodes located within the test room were selected to work with, thus excluding the rest of the domain to the working dataset as shown in Figure 8. As a second step, we choose to extract and work with the data from a 2D slice located at half height of the room: this location being a good compromise
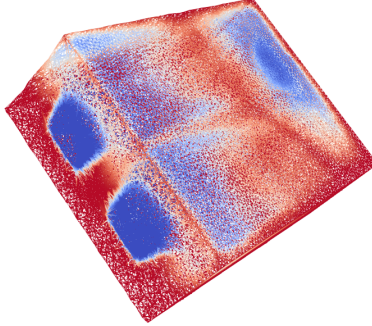
15

Figure 8: Unstructured mesh of the room coloured by $CO_2$ concentration. Blue colour means low $CO_2$ concentration, i.e. 400 ppm, and red colour indicates high value of concentration, i.e. 1420 ppm.

between the different heights of sensors used during the field experiment. Finally, two different pre-processing approaches were adopted:

- **"Structured dataset"** Data from the unstructured 2D slice are projected on a structured grid. The "Structured dataset" is generated by interpolating the $CO_2$ concentration values of the unstructured grid on a structured grid. The values stored in the final matrix corresponds to actual values of $CO_2$ concentration as shown in Figure 9.

- **"RGB dataset"** Data from the unstructured 2D slice are directly converted into a RGB image: a screenshot of the 2D slice coloured based on the $CO_2$ concentration values is created. The scalar bar of the RGB images is set based on the minimum and maximum $CO_2$ concentration, i.e. 400 ppm and 1420 ppm, respectively. This transformation allows to move from unstructured mesh to structured data since the RGB image is a 3D structured matrix of pixel values. The values stored in the final matrix corresponds to RGB values being between 0 and 255 as shown in Figure 10.

Both pre-processing approaches are performed for each timestep and the final size of the working matrix is a 180×250 regular grid.

As a final step, all the data are normalised between 0 and 1. The "RGB dataset" is divided by 255, while the "Structured dataset" is normalised based on the minimum and maximum values of $CO_2$ concentration.
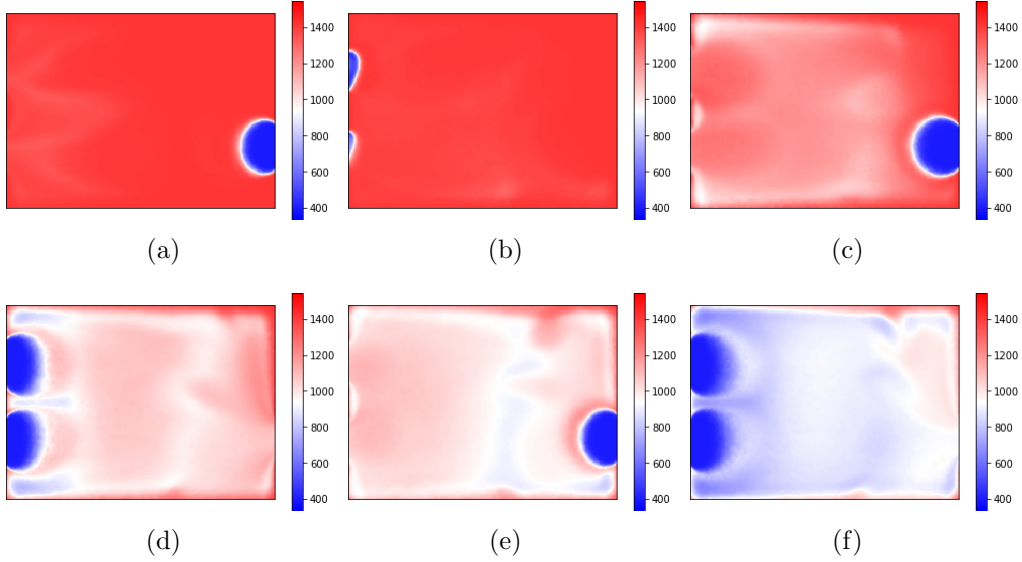
Figure 9: "Structured dataset" coloured by $CO_2$ concentration at timesteps (a) 509, (b) 643 (c) 1062, (d) 1359, (e) 1485 and (f) 1741. Blue colour means low $CO_2$ concentration, i.e. 400 ppm, and red colour indicates high value of concentration, i.e. 1420 ppm.

## 5.3. Sensors data pre-processing

The location of the sensors during the field experiment are shown in Figure 6 and are reported in our 180×250 matrix at the corresponding timestep. Based on the sampling rate of the sensors, 10 CFD output were selected corresponding to time levels for which we have sensors data. As a first pre-processing step, considering that the area of influence of one sensor has a radius of about 15 cm, a zone of 10 pixels × 10 pixels centred on the sensor location is defined and the value given by the sensor at that location is assigned to this whole area. The rendering of this process is shown in Figure 11. The second step consists in interpolating linearly the values of the sensors to the entire 2D structured grid as shown in Figure 12. This pre-processing is performed to be consistent with both the "Structured dataset" and the "RGB dataset", i.e. is done in terms of $CO_2$ concentration and RGB values scaled using the $CO_2$ concentration, respectively. As a final step, all the data are normalised between 0 and 1 as for the CFD pre-processing.

Figure 10: "RGB dataset" colored by RGB values scaled using the $CO_2$ concentration at timesteps (a) 509, (b) 643 (c) 1062, (d) 1359, (e) 1485 and (f) 1741. Blue colour means low $CO_2$ concentration, i.e. 400 ppm, and corresponds to RGB values equal to 0. Red colour indicates high value of concentration, i.e. 1420 ppm, and corresponds to RGB values equal to 255.

Figure 11: Rendering of the sensor values expanding to a radius of 15 cm coloured by $CO_2$ concentration (in ppm) at timesteps (a) 509, (b) 643, (c) 1062, (d) 1359, (e) 1485 and (f) 1741.

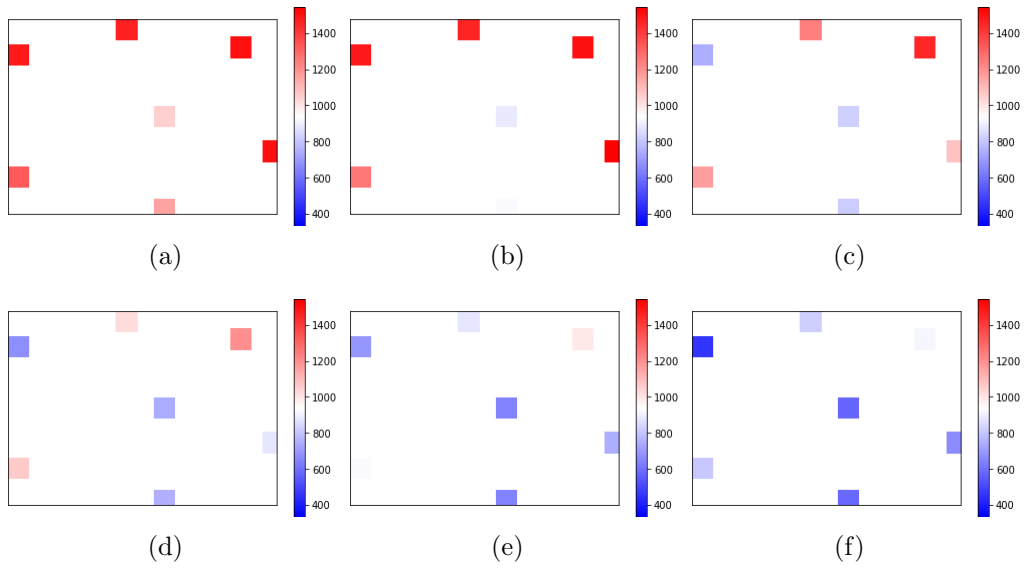Figure 12: Rendering of the linear interpolation of the sensor values coloured by $CO_2$ concentration (in ppm) at timesteps (a) 509, (b) 643, (c) 1062, (d) 1359, (e) 1485 and (f) 1741.

## 6. Results and Discussions

In this section, the procedure to determine the optimal network architectures of both the AutoEncoder and the LSTM is first presented. Then the results of the novel Latent Assimilation model developed in this paper are discussed based on the assimilation of the sensors data in both latent and physical space.

The data set is decomposed into training, validation and testing sets. In the CFD simulation, the flow field and the associated $CO_2$ concentration does not change much between consecutive steps. For this reason, we decide to divide the data in training, validation and testing sets making jumps. First, the CFD output at timesteps corresponding to sensors data are excluded and are assigned to the testing set. For the remaining data, two consecutive timesteps are considered for the training, then a jump is performed. The jumped data, i.e the ones not considered yet, are assigned to the validation and testing sets alternately. Considering a jump equal to 1, this process is summarised in Figure 13.

20

| Train | Train | Val | Train | Train | Test | • • • | Train | Train | Val | Train | Train | Test |
|-------|-------|-----|-------|-------|------|-------|-------|-------|-----|-------|-------|------|

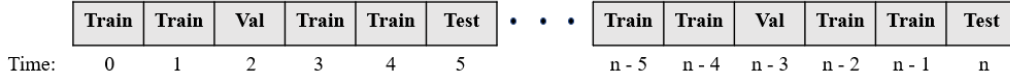Time:  0    1    2    3    4    5        n - 5   n - 4   n - 3   n - 2   n - 1   n

Figure 13: Training, Validation and Testing sets. Two consecutive timesteps are considered for the training. The jumped data are assigned to the validation and testing sets alternately.

## 6.1. AutoEncoder network architecture

The AutoEncoder implemented in this paper is a Convolutional AutoEncoder (CAE). Specifically, the encoder is composed by several convolutional layers followed by a flattened layer then a regular densely-connected layer which determine the shape of the latent space. In our CAE, the decoder architecture has almost the same structure than the encoder one: indeed, an additional convolutional layer is used in the decoder. Finding the optimal construction of the CAE architecture is divided in two steps:

1. Finding the optimal numbers of layers
2. Grid search: Finding the optimal hyperparameters

For each CAE network architecture tested, 5-fold cross-validation are performed for which the data are shuffled to make the neural network independent from the order of the data. Both the training and validation sets are used for the cross-validation. The evaluation of the CAE network architecture is estimated based on the mean and the standard deviation of the Mean Squared Error between the CAE prediction and the CFD output, i.e Mean-MSE and Std-MSE, respectively; the mean and the standard deviation of the Mean Absolute Error between the CAE prediction and the CFD output, i.e Mean-MAE and Std-MAE; and the mean and the standard deviation of the CAE execution time, i.e. Mean-Time and Std-Time. A low MSE/MAE standard deviation reflects that the model is stable and does not depend on the data used to train and validate it, while a low MSE/MAE means that the prediction is close to the real input, i.e has a good accuracy.

### 6.1.1. Optimal structure and number of layers

The *baseline* CAE network architecture is using the following fixed parameters:

- **Convolutional layers parameters**

21

– Number of filters: 32

– Activation function: Last decoder layer: sigmoid function to restrict the output in a range [0, 1] as the input; Rectified Linear Unit (ReLU) otherwise.

- **Regular densely-connected layer** Latent Space: 7

- **Training configuration**

  – Losses and metrics optimiser: Adam, learning rate of $1.10^{-3}$. Adam optimisation is a stochastic gradient descent method based on adaptive estimation of first-order and second-order moments well suited for problems with large data.

  – Number of epochs: 300

  – Batch size: 32

The following different CAE network architectures are tested using the "Structured dataset" in order to find the optimal numbers and structure of layers:

1. Encoder: 3 convolutional layers. Decoder: 4 convolutional layers. 3×3 kernel size.
2. Encoder: 4 convolutional layers. Decoder: 5 convolutional layers. 3×3 kernel size.
3. Encoder: 5 convolutional layers. Decoder: 6 convolutional layers. 3×3 kernel size.
4. Encoder: 4 convolutional layers. Decoder: 5 transpose convolutional layers. 3×3 kernel size.
5. Encoder: 4 convolutional layers. Decoder: 5 convolutional layers. 5×5 kernel size.
6. Encoder: 4 convolutional layers; 5×5, 5×5, 3×3 and 3×3 kernel sizes. Decoder: 5 convolutional layers; 3×3, 3×3, 5×5, 5×5 and 5×5 kernel sizes.

Results of the evaluation of the six CAE network architectures are reported in Table 1 where the configuration number N are given by the list above.

| N | Mean-MSE | Std-MSE | Mean-MAE | Std-MAE | Mean-Time | Std-Time |
|---|----------|---------|----------|---------|-----------|----------|
| 1 | 1.324e-02 | 2.072e-02 | 5.573e-02 | 7.575e-02 | 775.665 | 5.606e+00 |
| **2** | **2.435e-04** | **3.851e-05** | **1.034e-02** | **8.776e-04** | **812.293** | **1.436e+00** |
| 3 | 2.293e-02 | 2.763e-02 | 8.887e-02 | 9.818e-02 | 828.328 | 1.574e+00 |
| 4 | 2.587e-04 | 4.114e-05 | 1.041e-02 | 9.384e-04 | 746.804 | 3.089e+00 |
| 5 | 1.055e-02 | 2.099e-02 | 4.460e-02 | 7.922e-02 | 1222.251 | 6.627e+00 |
| 6 | 1.058e-02 | 2.094e-02 | 4.650e-02 | 7.826e-02 | 1164.464 | 4.605e+00 |

Table 1: Convolutional AutoEncoder performance evaluation for 6 network architectures. N denotes the configuration number as listed in the main text. Time is given in seconds. The bold row, i.e. configuration 2, is the configuration highlighting the best overall performance.

**Number of layers** Comparing configurations 1, 2 and 3 for which only the number of convolutional layers is changing, configuration 2, i.e. 4 convolutional layers for the encoder and 5 convolutional layers for the decoder, is the one highlighting the best performance in terms of both Mean-MSE and Mean-MAE with MSE two order of magnitude lower than configurations 1 and 3. Moreover, configuration 2 is the most stable regarding the standard deviations, reflecting well that this CAE network architecture does not depend on the data used to train and validate it. In addition, the execution time of configuration 2 is relatively acceptable to answer real-time problems. Hence, in the following, the number of layers is taken as the same than configuration 2: 4 for the encoder and 5 for the decoder.

**Convolutional vs transpose convolutional layers in the decoder** The accuracy (Mean-MSE/Mean-MAE) and the stability (Std-MSE/Std-MAE) are slightly better, while the execution time is slightly longer, when using convolutional layers (config. 2) rather than transpose convolutional layers (config. 4) in the decoder. As no major improvements in terms of MSE/MAE is observed when switching from convolutional (config. 2) to transpose convolutional layers in the decoder (config. 4), convolutional layers are used for the decoder.

**Size of the kernel** Configurations 2, 5 and 6 have the same layer number and the size of the kernel is changed. Using a 5×5 kernel size for all the layers (config. 5) or using a mix of 3×3 and 5×5 kernel sizes (config. 6) both increase the MSE/MAE by two order of magnitude compared to using

a 3×3 kernel size for all the layers (config. 2). In addition, the execution cost is considerably increased, about 50 % less efficient, when the kernel size is larger as the complexity scales with $k^3$ where $k$ is the kernel size. Overall, 3×3 is then used as the optimal kernel size in the following.

*6.1.2. Grid search for the optimal hyperparameters*
The grid search is now performed in order to find the optimal hyperparameters for both the "Structured dataset" and the "RGB dataset". The CAE network architecture is using the following fixed parameters:

- **Convolutional layers parameters**
    - Number of encoder/decoder convolutional layers: 4 and 5
    - Kernel size: 3x3
- **Regular densely-connected layer** Latent Space: 7
- **Training configuration** Adam optimiser: learning rate of $1.10^{-3}$.

The hyperparameters tested for the grid search are as follows:

- **Convolutional layers parameters**
    - Number of filters: 16, 32, 64
    - Activation function: Rectified Linear Unit (ReLU), Exponential Linear Unit (ELU).
- **Training configuration**
    - Number of epochs: 250, 300, 400
    - Batch size: 16, 32, 64

Table 2 shows the optimal hyperparameters found for each input dataset, while the evaluation performance are reported in Table 3. The optimal hyperparameters are the same for both dataset, i.e. 64 number of filters, a ReLU activation function and 400 epochs. Only the batch size differs: 32 for the "Structured dataset" and 16 for the "RGB dataset". The fact that "RGB dataset" needs less batch size than the "Structured dataset" can be potentially attributed to the fact that the former has 3 channels (R, G and B colours), so requiring less batch size. The results show that both datasets

24

have very similar accuracy and stability: low MSE and low standard deviation, of the order of $10^{-5}$, meaning that the CAE is not dependent on the set of input chosen to train it. Using "RGB dataset" highlights better performance, with Mean-MSE 57 % lower than when using "Structured dataset". However, using the "RGB dataset", more time is needed to train the CAE because an element of this dataset is composed by three channels, i.e. the R, G and B colour values.

| Dataset | Filters | Activation | Epochs | Batch size |
|---|---|---|---|---|
| **Structured** | 64 | ReLU | 400 | 32 |
| **RGB** | 64 | ReLU | 400 | 16 |

Table 2: Optimal Convolutional AutoEncoder hyperparameters determined by a grid search when using "Structured dataset" or "RGB dataset" as input.

| Dataset | Mean-MSE | Std-MSE | Mean-MAE | Std-MAE | Mean-Time | Std-Time |
|---|---|---|---|---|---|---|
| **Structured** | 8.509e-05 | 1.577e-05 | 6.182e-03 | 8.714e-04 | 1887.612 | 6.845e+00 |
| **RGB** | 3.670e-05 | 9.261e-06 | 3.601e-03 | 4.333e-04 | 2490.920 | 1.825e+01 |

Table 3: Convolutional AutoEncoder performance using "Structured dataset" or "RGB dataset" as input and the optimal hyperparameters found with the grid search (Table 2. Time is given in seconds.

*6.2. LSTM network architecture*

In this model, the training set is used for the fitting step and the validation set for the validation step. An extra splitting of the training set is performed: the data are split in small sequences such that one timestep is predicted and 3 timesteps are used as "look back" values as shown in Figure 14. All data are encoded with the AutoEncoder: each input sample of the LSTM is then a vector of 7 scalars. Finding the optimal construction of the LSTM architecture is divided in two steps:

1. Finding the optimal numbers of layers
2. Grid search: Finding the optimal hyperparameters

For each LSTM network architecture tested, the fitting and the evaluation of the model is repeating 5 times. As for the CAE network architecture, the

LSTM network architecture performance is evaluated based on the Mean-MSE, Std-MSE, Mean-MAE, Std-MAE, Mean-Time and Std-Time.
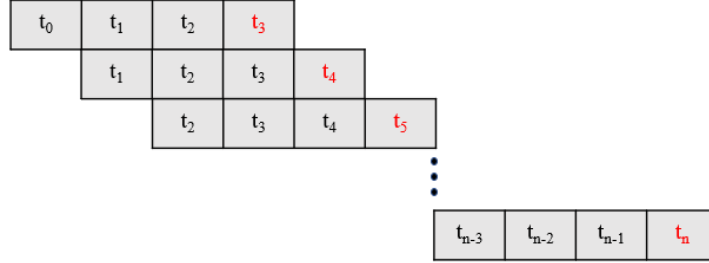


Figure 14: Data splitting for LSTM model. Timesteps in black are the ones used to make the prediction and the red ones are the timesteps predicted.

*6.2.1. Optimal number of layers*

The *baseline* LSTM network architecture, tested with the "Structured dataset", is using the following fixed parameters:

- **LSTM layers parameters**
    - Number of neurons: 30, i.e the dimensionality of the output space
    - Activation function: ReLU
    - Number of steps: 3
- **Regular densely-connected layer** Latent Space: 7
- **Training configuration**
    - Optimiser: Adam with a learning rate of $1.10^{-3}$
    - Number of epochs: 300
    - Batch size: 32

LSTMs are stacked from 1 to 5 times in order to see if the model gains in accuracy, stability and efficiency: the results are shown in Table 4. The single layer LSTM is the one highlighting the best accuracy with the lowest Mean-MSE and Mean-MAE values. Indeed, the input of the LSTM consists of a 7×1 vector and adding more LSTM layer introduces overfitting bias. In addition, the standard deviation, reflecting the stability, of the single layer

26

LSTM are about one order of magnitude lower than the other tested LSTM. Finally, as expected, the single layer LSTM is also the most efficient in term of computation cost.

| N | Mean-MSE | Std-MSE | Mean-MAE | Std-MAE | Mean-Time | Std-Time |
|---|----------|---------|----------|---------|-----------|----------|
| **1** | **1.634e-02** | **2.510e-03** | **7.822e-02** | **8.553e-03** | **230.355** | **1.050e+00** |
| 2 | 2.822e-02 | 7.244e-03 | 1.060e-01 | 1.679e-02 | 360.877 | 6.618e-01 |
| 3 | 4.619e-02 | 1.942e-02 | 1.381e-01 | 3.034e-02 | 494.254 | 2.258e+00 |
| 4 | 5.020e-02 | 1.675e-02 | 1.414e-01 | 2.778e-02 | 658.039 | 2.632e+00 |
| 5 | 4.742e-02 | 1.183e-02 | 1.459e-01 | 2.075e-02 | 806.001 | 5.921e+00 |

Table 4: LSTM performance evaluation for 5 network architectures. N denotes the number of stacked LSTMs. Time is given in seconds. The bold row, i.e. 1 LSTM, is the configuration highlighting the best overall performance.

*6.2.2. Grid Search*

The grid search is now performed in order to find the optimal hyperparameters for both the "Structured dataset" and the "RGB dataset". The LSTM network architecture is using the following fixed parameters:

- **LSTM layers parameters** Number of layers: 1

- **Regular densely-connected layer** Latent Space: 7

- **Training configuration** Optimiser: Adam, learning rate of $1.10^{-3}$

The hyperparameters tested for the grid search are as follows:

- **LSTM layers parameters**

    - Number of neurons: 30, 50, 70

    - Activation function: ReLU, ELU

    - Number of steps: 3, 5, 7

- **Training configuration**

    - Number of epochs: 200, 300, 400

    - Batch size: 16, 32, 64

Table 5 shows the optimal hyperparameters found for each input dataset, while the evaluation performance are reported in Table 6. Exponential Linear Unit (ELU) appears to be the optimal activation function, with an epochs of 400 and a batch size of 16 for both the input dataset. The results show that the "RGB dataset" needs more neurons and more back observations than the "Structured dataset". From Table 6, it can be seen that the LSTM with "RGB dataset" as input has better accuracy and takes also less time than when using "Structured dataset" input. Indeed, the accuracy is about 45 % higher when using input "RGB dataset", while the execution time is reduced by approximately 27 %.

| Dataset | Neurons | Activation | Steps | Epochs | Batch size |
|---|---|---|---|---|---|
| **Structured** | 30 | Elu | 3 | 400 | 16 |
| **RGB** | 50 | Elu | 7 | 400 | 16 |

Table 5: Optimal LSTM hyperparameters determined by a grid search when using "Structured dataset" or "RGB dataset" as input.

| DataSet | Mean-MSE | Std-MSE | Mean-MAE | Std-MAE | Mean-Time | Std-Time |
|---|---|---|---|---|---|---|
| **Structured** | 1.233e-02 | 1.398e-03 | 6.743e-02 | 3.428e-03 | 949.328 | 7.508e+00 |
| **RGB** | 6.712e-03 | 6.704e-04 | 4.661e-02 | 3.004e-03 | 690.594 | 2.743e+00 |

Table 6: LSTM performance using "Structured dataset" or "RGB dataset" as input and the optimal hyperparameters found with the grid search (Table 5). Time given in seconds.

### 6.3. Latent Assimilation model

In this section, results of our novel Latent Assimilation (LA) model are presented: the assimilation takes place in the latent space. The Testing set is considered and both dataset, i.e. "Structured dataset" and "RGB dataset", are encoded using the AutoEncoders with optimal network architecture as presented in Section 6.1.2. The predictions are performed through the LSTM and are updated using the corresponding observations through Optimal Interpolated Kalman Filter (KF).

In the KF, the error covariance matrix $\hat{Q}$ is computed as $\hat{Q} = VV^T$, where $V$ is as defined in equation (18). Since both predictions of the model and observations are values of $CO_2$ or pixels, i.e. the observations do not have

to be transformed, the operator $\hat{H}$ is an identity matrix. We studied how KF improves the accuracy of the prediction by testing different forms of the observation error covariance matrix $\hat{R}$: computed using equation (18) or, fixed as $\hat{R} = 0.01I$, $0.001I$, $0.0001I$ where $I \in \mathbb{R}^{p \times p}$ denotes the identity matrix. This last assumption is usually made to give higher fidelity and trust to the observations [24].

*6.3.1. Latent space*

The MSE between the background data and the observed data in the latent space for the "Structured dataset" and the "RGB dataset", without performing data assimilation, are $7.220 \times 10^{-1}$ and $5.447 \times 10^{-1}$, respectively. Table 7 shows values of MSE in the latent space between the assimilated data $h_t^a$ and the observed data as well as the execution time of the assimilation for both input dataset. As expected, we can observe an improvement in the execution time of the assimilation in assuming $\hat{R}$ as a diagonal matrix instead of a full matrix. In addition, the assimilation increases the accuracy of the model, whatever the input dataset used, with MSE values about 2.2 times lower compared to without assimilation, highlighting that our novel Latent Assimilation model is behaving as expected. Using $\hat{R}$ as an identity matrix of the form $0.0001I$ allows to improve the accuracy by up to 4 order of magnitude.

| | | $\hat{R}$ | | | |
| --- | --- | --- | --- | --- | --- |
| | | cov-matrix eq. (18) | 0.01I | 0.001I | 0.0001I |
| MSE | Structured | 3.215e-01 | 1.250e-02 | 1.787e-03 | 3.722e-05 |
| | RGB | 2.444e-01 | 1.002e-02 | 4.409e-03 | 1.640e-03 |
| Time | Structured | 1.541e-03 | 5.410e-04 | 4.854e-04 | 4.823e-04 |
| | RGB | 2.145e-03 | 5.388e-04 | 4.847e-04 | 4.807e-04 |

Table 7: MSE values in the **latent space** and execution time of the assimilation in seconds of the Latent Assimilation model for different form of the observations error covariance matrix $\hat{R}$ in the latent space when using the "Structured dataset" or the "RGB dataset" as input.

*6.3.2. Physical space*

After having performed the DA in the latent space, the results $h_t^a$ are reported in the physical space through the decoder which gives $x_t^a$. Figure 15 shows in the physical space the results of the assimilation for the timesteps 509, 1062 and 1485 using our novel LA model.
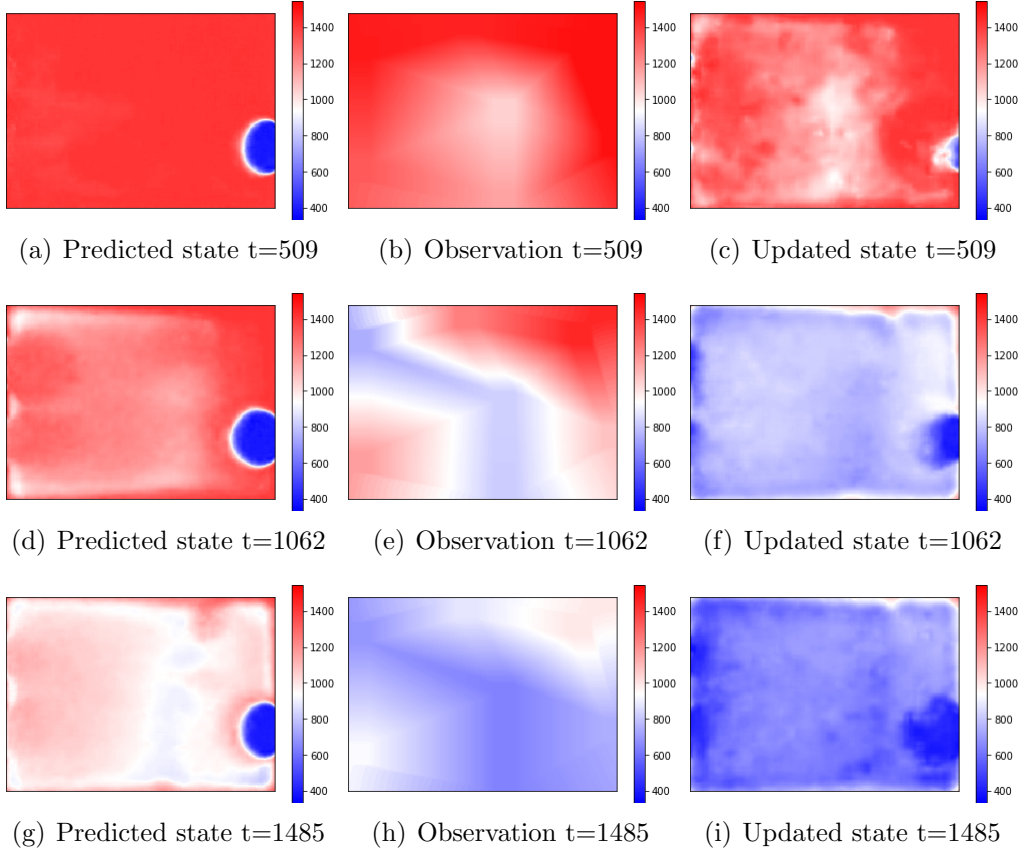
Figure 15: Assimilation results colored by $CO_2$ concentration (in ppm) at timesteps 509, 1062 and 1485. (a) LSTM predicted state, (b) interpolated observation and (c) updated state obtained using the latent assimilation. Blue colour means low $CO_2$ concentration, i.e. 400 ppm, and red colour indicates high value of concentration, i.e. 1420 ppm.

The MSE in the physical space using our LA model is then compared with the one using a standard Data Assimilation (sDA) procedure. sDA is performed in the physical space using a Kalman Filter approach (equations (4)- (6)), where $R \in \mathbb{R}^{n \times n}$ is defined in the physical space. Table 8 shows values of MSE in the physical space between the assimilated data $x_t^a$ and the observed data as well as the execution time of the assimilation for our LA model and the standard methodology (sDA) when using the "Structured dataset" as input. The MSE between the background data and the observed data in the physical space, without performing data assimilation, is $6.491 \times 10^{-2}$. Both

LA and sDA improve the accuracy of the forecasting as shown in Table 8: however it can be observed that the LA model gives 35 % more accuracy than a sDA model. In addition, LA performs better in terms of execution time with respect to a sDA: indeed sDA works directly with big matrices making it slower by six order of magnitude.

| | | $\hat{\mathbf{R}}$ | | | |
|---|---|---|---|---|---|
| | | cov-matrix eq. (18) | 0.01I | 0.001I | 0.0001I |
| MSE | LA | 3.356e-02 | 6.933e-04 | 1.211e-04 | 2.691e-06 |
| | sDA | 5.179e-02 | 6.928e-03 | 6.928e-03 | 6.997e-03 |
| Time | LA | 1.541e-03 | 5.410e-04 | 4.854e-04 | 4.823e-04 |
| | sDA | 2.231e+03 | 2.148e+03 | 2.186e+03 | 2.159e+03 |

Table 8: MSE values of $x_t^a$ using our novel Latent Assimilation (LA) model or using a standard Data Assimilation (sDA) procedure for different form of the observations error covariance matrix $\hat{R}$ when using the "Structured dataset" as input. MSE are computed in the **physical space**. Execution time of the assimilation in seconds.

*6.3.3. Size of the latent space*

In this section, the impact of increasing the size of the latent space is discussed. Results are presented for the "Structured dataset" only. Table 9 and Table 10 give the MSE values of the Latent Assimilation model with different latent space sizes, from 1000 to 20000, computed in the latent and physical space, respectively. The column "No DA" reports the MSE values without the assimilation. Table 11 reports the execution time of the assimilation.

Defining $\hat{R}$ as an identity matrix always highlights better accuracy whatever the latent space size. Increasing the latent space size tends to decrease the MSE, i.e. gain in accuracy, in both the latent and the physical space. Overall, a latent space size equal to 18000 seems optimal for this problem, whatever the form of the observations error covariance matrix $\hat{R}$, which represents about 40 % of the original data. However, the execution time of the assimilation can be up to 5 order of magnitude higher when using the optimal latent space size compared to a lower size. Finding the optimal parameters of our LA depends the expectancy of the user as a balance needs to be taken between accuracy and efficiency. Regarding the small accuracy gain while increasing the latent space size, it is recommended to work with the smallest latent space size as possible in order to benefit of the best efficiency while still keeping a high accuracy.

31

| Latent Space | No DA | 0.01 I | 0.001 I | 0.0001 I |
|---|---|---|---|---|
| 1000 | 4.363e-03 | 5.141e-04 | 4.980e-04 | 4.961e-04 |
| 3000 | 6.472e-04 | 9.933e-05 | 9.566e-05 | 9.501e-05 |
| 5000 | 6.012e-04 | 1.009e-04 | 9.184e-05 | 9.076e-05 |
| 7000 | 7.295e-04 | 1.179e-04 | 1.146e-04 | 1.142e-04 |
| 12000 | 1.727e-04 | 3.651e-05 | 3.540e-05 | 3.531e-05 |
| 15000 | **1.660e-04** | 3.529e-05 | 3.463e-05 | 3.457e-05 |
| 18000 | 2.941e-04 | **3.141e-05** | **3.086e-05** | **3.085e-05** |
| 20000 | 3.465e-04 | 6.506e-05 | 6.177e-05 | 6.073e-05 |

Table 9: MSE values of LA model in the **latent space** with different latent space sizes. The column Latent Space indicates the size of the latent space, the column "No DA" indicates the MSE of the model without the assimilation and the other columns indicate the MSE of the LA model with different observation error covariance matrix $\hat{R}$.

| Latent Space | No DA | 0.01 I | 0.001 I | 0.0001 I |
|---|---|---|---|---|
| 1000 | 3.532e-02 | 1.347e-03 | 1.278e-03 | 1.278e-03 |
| 3000 | 3.261e-02 | 1.828e-03 | 1.694e-03 | 1.653e-03 |
| 5000 | 2.822e-02 | 1.975e-03 | 1.689e-03 | 1.667e-03 |
| 7000 | 3.352e-02 | 1.248e-03 | 1.171e-03 | 1.155e-03 |
| 12000 | 2.479e-02 | 1.248e-03 | 1.171e-03 | 1.155e-03 |
| 15000 | **1.734e-02** | 1.325e-03 | 1.262e-03 | 1.253e-03 |
| 18000 | 3.703e-02 | **1.080e-03** | **9.848e-04** | **9.743e-04** |
| 20000 | 2.514e-02 | 1.621e-03 | 1.424e-03 | 1.365e-03 |

Table 10: MSE values of LA model in the **physical space** with different latent space sizes. The column Latent Space indicates the size of the latent space, the column "No DA" indicates the MSE of the model without the assimilation and the other columns indicate the MSE of the LA model with different observation error covariance matrix $\hat{R}$.

| Latent Space | 0.01 I | 0.001 I | 0.0001 I |
| :---: | :---: | :---: | :---: |
| 1000 | 5.637e-01 | 5.525e-01 | 5.637e-01 |
| 3000 | 2.272e+00 | 2.260e+00 | 2.147e+00 |
| 5000 | 5.197e+00 | 5.301e+00 | 5.358e+00 |
| 7000 | 1.104e+01 | 1.125e+01 | 1.136e+01 |
| 12000 | 4.289e+01 | 4.353e+01 | 4.375e+01 |
| 15000 | 7.960e+01 | 8.072e+01 | 8.180e+01 |
| 18000 | 1.432e+02 | 1.441e+02 | 1.443e+02 |
| 20000 | 2.096e+02 | 2.124e+02 | 1.997e+02 |

Table 11: Execution time in seconds of the assimilation in our LA model with different latent space sizes. The column Latent Space indicates the size of the latent space and the other columns indicate the execution time of the assimilation with different observation error covariance matrix $\hat{R}$.

## 7. Conclusion and Future Work

In this paper, we proposed a new methodology called Latent Assimilation (LA) to efficiently and accurately perform Data Assimilation (DA). LA consists in performing the Optimal Kalman Filter in the latent space obtained by a Convolutional AutoEncoder with non-linear encoder functions and non-linear decoder functions. In the latent space, the dynamic system is represented by a surrogate model built by an LSTM network to train a function that emulates the dynamic system in the latent space. The data from the dynamic model and the real data coming from the sensors are both processed through the AutoEncoder.

We applied the methodology to a real test case and we have shown that the LA performs better than a standard DA in terms of both accuracy and efficiency. The data of the real test case was time-series data representing the airflow within a naturally ventilated office room. The data was provided by CFD on an unstructured mesh and we pre-processed these data to extract two different structured datasets: one composed of 2D matrices of $CO_2$ concentration ("Structured dataset") and the other one composed of RGB images colored by $CO_2$ concentration ("RGB dataset"). We pre-processed also the data coming from sensors in the same manner.

We tried different AutoEncoder configurations and we performed a grid search for both input datasets in order to determine the optimal configurations. The same was done for the LSTM: it is the surrogate model. We performed the assimilation in the latent space using the Latent Assimilation model for both datasets as input. We tested also the standard data assimilation in the physical space and we have shown that LA performs better in terms of both efficiency and accuracy.

In conclusion, we have successfully proposed and developed a novel model able to assimilate data in the latent space, thus answering the needs of accuracy, stability and efficiency required by real-time systems. This methodology can be used for example to predict in real-time the load of virus, such as the SARS-COV-2, in indoor spaces by linking it to the concentration of $CO_2$ [22].

There are different improvements that could be applied to the model to be used with more challenging applications:

- Develop an implementation of LA to emulate a variational DA [24] which is often applied to big data problems. In particular, we will focus on a 4D Variational (4DVar) method. 4DVar is a computational expensive method as it is developed to assimilate several observations (distributed in time) for each timestep of the forecasting model. We will develop an extended version of LA able to assimilate set of distributed observations for each timestep and, then, able to perform a 4DVar;

- Add a third dimension, i.e. test the methodology on a 3D space using a 3D Convolutional Autoencoder. Instead of cutting a slice, the 3D Convolutional Autoencoder will work on the complete room space without losing information;

- Recent research studies has started in the direction of working directly with unstructured meshes. It will be challenging developing Latent Assimilation with an Encoder-Decoder which works directly on a 3D adaptive and unstructured mesh;

- Instead of using only indoor data, the methodology could be applied considering the exchange with the outdoor environment or tested in different applications, i.e ocean.

## Acknowledgments

## Bibliography

[1] J. Song, S. Fan, W. Lin, L. Mottet, H. Woodward, M. Davies Wykes, R. Arcucci, D. Xiao, J.-E. Debay, H. ApSimon, et al., Natural ventilation in cities: the implications of fluid mechanics, Building Research & Information 46 (2018) 809–828.

[2] I. C. L. AMCG, Fluidity manual v4.1.12, Imperial College London (2015). URL: https://figshare.com/articles/Fluidity_Manual/1387713.

[3] C. K. Wikle, L. M. Berliner, A bayesian tutorial for data assimilation, Physica D: Nonlinear Phenomena 230 (2007) 1–16.

[4] R. Arcucci, L. Mottet, C. Pain, Y.-K. Guo, Optimal reduced space for variational data assimilation, Journal of Computational Physics 379 (2019) 51–69.

[5] R. Arcucci, C. Q. Casas, D. Xiao, L. Mottet, F. Fang, P. Wu, C. C. Pain, Y.-K. Guo, A domain decomposition reduced order model with data assimilation (dd-roda)., in: PARCO, 2019, pp. 189–198.

[6] P. C. Hansen, J. G. Nagy, D. P. O'leary, Deblurring images: matrices, spectra, and filtering, SIAM, 2006.

[7] A. Hannachi, A primer for eof analysis of climate data, Department of Meteorology, University of Reading (2004) 1–33.

[8] J. Mack, R. Arcucci, M. Molina-Solana, Y.-K. Guo, Attention-based convolutional autoencoders for 3d-variational data assimilation, Computer Methods in Applied Mechanics and Engineering 372 (2020) 113291.

[9] P. Wu, J. Sun, X. Chang, W. Zhang, R. Arcucci, Y. Guo, C. C. Pain, Data-driven reduced order model with temporal convolutional neural network, Computer Methods in Applied Mechanics and Engineering 360 (2020) 112766.

[10] C. Q. Casas, R. Arcucci, P. Wu, C. Pain, Y.-K. Guo, A reduced order deep data assimilation model, Physica D: Nonlinear Phenomena 412 (2020) 132615.

[11] R. Arcucci, L. Moutiq, Y.-K. Guo, Neural assimilation, in: International Conference on Computational Science, Springer, 2020, pp. 155–168.

[12] S.-A. Boukabara, V. Krasnopolsky, J. Q. Stewart, E. S. Maddy, N. Shahroudi, R. N. Hoffman, Leveraging modern artificial intelligence for remote sensing and nwp: Benefits and challenges, Bulletin of the American Meteorological Society 100 (2019) ES473–ES491.

[13] V. Babovic, M. Keijzer, M. Bundzel, From global to local modelling: a case study in error correction of deterministic models, in: Proceedings of fourth international conference on hydroinformatics, 2000.

[14] V. Babovic, R. Caňizares, H. R. Jensen, A. Klinting, Neural networks as routine for error updating of numerical models, Journal of Hydraulic Engineering 127 (2001) 181–193.

[15] V. Babovic, D. R. Fuhrman, Data assimilation of local model error forecasts in a deterministic model, International journal for numerical methods in fluids 39 (2002) 887–918.

[16] J. Zhu, S. Hu, R. Arcucci, C. Xu, J. Zhu, Y.-k. Guo, Model error correction in data assimilation by integrating neural networks, Big Data Mining and Analytics 2 (2019) 83–91.

[17] R. S. Cintra, H. Campos Velho, Data assimilation by artificial neural networks for an atmospheric general circulation model, Advanced Applications for Artificial Neural Networks (2018) 265.

[18] K. Loh, P. S. Omrani, R. van der Linden, Deep learning and data assimilation for real-time production prediction in natural gas wells, arXiv preprint arXiv:1802.05141 (2018).

[19] C. A. Quilodrán Casas, Fast ocean data assimilation and forecasting using a neural-network reduced-space regional ocean model of the north Brazil current, Ph.D. thesis, Imperial College London, 2018.

[20] P. Becker, H. Pandya, G. Gebhardt, C. Zhao, J. Taylor, G. Neumann, Recurrent kalman networks: Factorized inference in high-dimensional deep feature spaces, arXiv preprint arXiv:1905.07357 (2019).

[21] M. Watter, J. Springenberg, J. Boedecker, M. Riedmiller, Embed to control: A locally linear latent dynamics model for control from raw images, in: Advances in neural information processing systems, 2015, pp. 2746–2754.

[22] Z. Peng, J. L. Jimenez, Exhaled co2 as covid-19 infection risk proxy for different indoor environments and activities, medRxiv (2020).

[23] R. E. Kalman, A new approach to linear filtering and prediction problems, Journal of Basic Engineering (1960).

[24] M. Asch, M. Bocquet, M. Nodet, Data assimilation: methods, algorithms, and applications, SIAM, 2016.

[25] I. Goodfellow, Y. Bengio, A. Courville, Y. Bengio, Deep learning, volume 1, MIT press Cambridge, 2016.

[26] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural computation 9 (1997) 1735–1780.

[27] G. Tajnafoi, R. Arcucci, L. Mottet, C. Vouriot, M. M. Solana, C. Pain, Y.-K. Guo, Variational Gaussian Process for Optimal Sensor Placement., Applied Mathematics 66 (2021) In press.