

## Journal Pre-proof

A comparison of combined data assimilation and machine learning methods for offline and online model error correction

Alban Farchi, Marc Bocquet, Patrick Laloyaux, Massimo Bonavita, Quentin Malartic



PII: S1877-7503(21)00143-5  
DOI: <https://doi.org/10.1016/j.jocs.2021.101468>  
Reference: JOCS 101468

To appear in: *Journal of Computational Science*

Received date: 23 July 2021  
Revised date: 7 September 2021  
Accepted date: 28 September 2021

Please cite this article as: A. Farchi, M. Bocquet, P. Laloyaux et al., A comparison of combined data assimilation and machine learning methods for offline and online model error correction, *Journal of Computational Science* (2021), doi: <https://doi.org/10.1016/j.jocs.2021.101468>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2021 Published by Elsevier B.V.

1 A comparison of combined data assimilation and machine learning methods  
2 for offline and online model error correction

3 Alban Farchi<sup>a,\*</sup>, Marc Bocquet<sup>a</sup>, Patrick Laloyaux<sup>b</sup>, Massimo Bonavita<sup>b</sup>, Quentin Malartic<sup>a,c</sup>

4 <sup>a</sup>*CEREA, École des Ponts and EDF R&D, Île-de-France, France*

5 <sup>b</sup>*ECMWF, Shinfield Park, Reading, United Kingdom*

6 <sup>c</sup>*LMD/IPSL École Normale Supérieure and PSL University, École Polytechnique, Université Paris-Saclay, Sorbonne*  
7 *Université, CNRS, Paris, France*

---

8 **Abstract**

Recent studies have shown that it is possible to combine machine learning methods with data assimilation to reconstruct a dynamical system using only sparse and noisy observations of that system. The same approach can be used to correct the error of a knowledge-based model. The resulting surrogate model is hybrid, with a statistical part supplementing a physical part. In practice, the correction can be added as an integrated term (*i.e.* in the model resolvent) or directly inside the tendencies of the physical model. The resolvent correction is easy to implement. The tendency correction is more technical, in particular it requires the adjoint of the physical model, but also more flexible. We use the two-scale Lorenz model to compare the two methods. The accuracy in long-range forecast experiments is somewhat similar between the surrogate models using the resolvent correction and the tendency correction. By contrast, the surrogate models using the tendency correction significantly outperform the surrogate models using the resolvent correction in data assimilation experiments. Finally, we show that the tendency correction opens the possibility to make online model error correction, *i.e.* improving the model progressively as new observations become available. The resulting algorithm can be seen as a new formulation of weak-constraint 4D-Var. We compare online and offline learning using the same framework with the two-scale Lorenz system, and show that with online learning, it is possible to extract all the information from sparse and noisy observations.

9 *Keywords:* data assimilation, machine learning, model error, surrogate model, neural networks

---

10 **1. Introduction: machine learning for model error correction**

11 Over the past decade, data-driven methods, and in particular machine learning (ML), have shown  
12 remarkable success in reproducing complex spatiotemporal processes, and have therefore been used in an  
13 increasing number of applications (LeCun et al., 2015; Goodfellow et al., 2016; Chollet, 2018). In the

---

\*Corresponding author

Email address: [alban.farchi@enpc.fr](mailto:alban.farchi@enpc.fr) (Alban Farchi)

14 geosciences only, there is a fairly recent wealth of studies dealing with the problem of inferring the dynamics  
15 of a system from observations. Typical examples include the use of analogs, delay coordinates embedding,  
16 random forests, echo state networks and other neural networks such as residual, recurrent, or convolutional  
17 neural networks (Brunton et al., 2016; Hamilton et al., 2016; Lguensat et al., 2017; Pathak et al., 2018;  
18 Dueben and Bauer, 2018; Fablet et al., 2018; Scher and Messori, 2019; Weyn et al., 2019; Arcomano et al.,  
19 2020). Most, if not all, of these examples implement a type of supervised learning where the goal is to  
20 minimise the loss function, a measure of the discrepancy between the statistical model (also called surrogate  
21 model) predictions and the observation dataset. The underlying assumption is that the system is fully  
22 observed without or with very little noise. In order to handle sparse and noisy observations, which is the case  
23 in most realistic systems in the geosciences, more and more studies consider the possibility of hybridising ML  
24 and data assimilation (DA) techniques (Abarbanel et al., 2018; Bocquet et al., 2019; Brajard et al., 2020;  
25 Bocquet et al., 2020a; Arcucci et al., 2021). In practice, DA tools are used, with the surrogate model, to  
26 estimate the state of the system from the observations while ML tools are used to estimate the surrogate  
27 model from the analysis (estimated) state. This method has been reformulated using a unifying Bayesian  
28 formalism by Bocquet et al. (2020a).

29 In the geosciences, even though models are affected by errors (*e.g.*, misrepresented physical phenomena,  
30 unresolved small-scale processes, numerical integration errors, etc), they benefit from a long history of  
31 modelling and therefore they already provide a solid baseline. For this reason, recent studies focus on using  
32 ML techniques for model error correction instead of full model emulation (Rasp et al., 2018; Bolton and  
33 Zanna, 2019; Jia et al., 2019; Watson, 2019; Bonavita and Laloyaux, 2020; Brajard et al., 2021; Gagne  
34 et al., 2020; Wikner et al., 2020; Farchi et al., 2021). The idea is to build a hybrid model with a physical,  
35 knowledge-based part, and a statistical part to supplement it. This means that the statistical model is trained  
36 to learn the error of the physical model. The underlying rationale is that model error correction should be  
37 an easier inference problem than full model emulation (Jia et al., 2019; Watson, 2019; Farchi et al., 2021).

38 From a technical perspective, the geoscientific models are based on a set of physical laws, usually  
39 represented as ordinary or partial differential equations (ODEs or PDEs). These equations define the  
40 *tendencies* of the model. A numerical scheme is used to integrate them for a small time step, and several  
41 integration steps are composed to define the *resolvent* between two forecast times. Following Farchi et al.  
42 (2021), two strategies are possible for a correction term: (i) apply an integrated correction between two  
43 forecast times, *i.e.* in the resolvent, or (ii) apply a correction directly in the tendencies. The first method is  
44 by far the simplest to implement, which is why it is the most widely applied, but it faces some limitations,  
45 in particular when using the hybrid model for DA experiments. The first objective of the present paper is  
46 hence to make an exhaustive comparison of the two methods for both forecast and assimilation experiments  
47 in a simplified modelling framework.

48 Beyond the design of the model error correction — or more generally of any surrogate model — the

49 question of the use of observations arise. In most cases, the statistical model is only trained once the entire  
 50 observation dataset is available: this is called *offline learning*. The other option, *online*, or *sequential learning*,  
 51 *i.e.* improving the surrogate model as new observations become available, is also possible in ML, even if it  
 52 is less common because the methods usually require very large datasets to achieve good performance. In  
 53 a context where information is only available through sparse and noisy observations, this means that we  
 54 have to learn both the state of the system and the surrogate model at the same time. This is the topic of  
 55 several recent studies (Bocquet et al., 2020b; Gottwald and Reich, 2021), which emphasise the connections  
 56 between this problem and classical parameter estimation in DA (Ruiz et al., 2013; Pulido et al., 2018). In the  
 57 geosciences, online learning is more natural because observations are acquired sequentially, and improvements  
 58 can be expected before having a long series of observations since the training begins from the first observation.  
 59 Therefore, the second objective of the present paper is to explore the possibility to use online learning for  
 60 model error correction.

61 The paper is organised as follows. Section 2 introduces the main methodological aspects for offline  
 62 learning. We start with a brief overview of the Bayesian framework for combining DA and ML and how  
 63 it can be used for model error correction. We then discuss the advantages and drawbacks of applying a  
 64 correction term in the resolvent or in the tendencies, with an emphasis on the implications for forecast and  
 65 assimilation applications. The two methods are compared in section 3 using the two-scale Lorenz model  
 66 (L05III, Lorenz, 2005). Section 4 further develops the methodology to enable online learning for model error  
 67 correction. Section 5 illustrates the use of online learning with the same L05III model, and compares it to  
 68 offline learning. Finally, conclusions are drawn in section 6.

## 69 2. Offline learning of model error with resolvent or tendency correction

### 70 2.1. A Bayesian framework for data assimilation and machine learning

71 The starting point of the present work is a series of observations  $\mathbf{y}_k \in \mathbb{R}^{N_y}$  of a system at discrete times  
 72  $t_k$  for  $k \in \mathbb{N}$ . The state of the system is represented by a vector  $\mathbf{x}_k \in \mathbb{R}^{N_x}$ . The observations are related to  
 73 the state through the observation equation

$$74 \quad \mathbf{y}_k = \mathcal{H}_k(\mathbf{x}_k) + \mathbf{v}_k, \quad (1)$$

75 where  $\mathcal{H}_k : \mathbb{R}^{N_x} \rightarrow \mathbb{R}^{N_y}$  is the observation operator and  $\mathbf{v}_k \in \mathbb{R}^{N_y}$  the observation error at time  $t_k$ . We  
 76 assume that the time evolution of the state is governed by the state equation

$$77 \quad \mathbf{x}_{k+1} = \mathcal{M}_k^t(\mathbf{x}_k) + \mathbf{w}_k, \quad (2)$$

78 where  $\mathcal{M}_k^t : \mathbb{R}^{N_x} \rightarrow \mathbb{R}^{N_x}$  is the resolvent of the (unknown) true dynamical model from  $t_k$  to  $t_{k+1}$ , and  
 79  $\mathbf{w}_k \in \mathbb{R}^{N_x}$  is the corresponding model error (*e.g.*, related to sub-scale processes). To simplify the presentation,  
 80 we make the following assumptions:

- 81 • observations are available at regular intervals  $t_k = k\Delta t$ ;
- 82 • the observation operator is constant over time  $\mathcal{H}_k \equiv \mathcal{H}$ ;
- 83 • the observation error is uncorrelated in time and normally distributed  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ , where  $\mathbf{R}$  is the
- 84 observation error covariance matrix;
- 85 • the model error  $\mathbf{w}_k$  is uncorrelated to the observation error  $\mathbf{v}_k$ .

86 In particular, the third point implies that the observations are not biased, which helps to attribute correctly  
 87 the model errors. Furthermore, we also make the assumption that the true dynamical model is autonomous,  
 88 in which case  $\mathcal{M}_k^t \equiv \mathcal{M}_{\Delta t}^t$  the resolvent of the surrogate model for a  $\Delta t$  integration. The extension of the  
 89 present work to non-autonomous dynamics is not trivial and briefly discussed in section 5.6.

90 Our goal is to derive a surrogate of the true model, which can be used to predict  $\mathbf{x}_{k+1}$  from  $\mathbf{x}_k$ . Let  $\mathbf{p}$  be  
 91 the set of parameters defining the surrogate model. The discrepancy between the surrogate model predictions  
 92 and the observations is measured with a cost function. A traditional ML approach to this problem is to use  
 93 dense observations (*i.e.*,  $\mathcal{H} = \mathbf{I}$  the identity operator) and to neglect the observation errors (*i.e.*, assuming  
 94 that  $\mathbf{R} = \mathbf{0}$ ), which yields the following cost function:

$$95 \quad \mathcal{J}(\mathbf{p}) \triangleq \mathcal{L}(\mathbf{p}) + \frac{1}{2} \sum_{k=0}^{N_t-1} \|\mathbf{y}_{k+1} - \mathcal{M}_{\Delta t}(\mathbf{p}, \mathbf{y}_k)\|_{\mathbf{Q}_k^{-1}}^2, \quad (3)$$

96 where  $\mathcal{L}$  is a regularisation (prior) term on  $\mathbf{p}$ ,  $N_t$  is the number of observation batches used to define  $\mathcal{J}$ , and  
 97  $\mathbf{x} \mapsto \mathcal{M}_{\Delta t}(\mathbf{p}, \mathbf{x})$  is the resolvent of the surrogate model for a  $\Delta t$  integration. The matrix norm notation  
 98  $\|\mathbf{v}\|_{\mathbf{A}}^2$  stands for  $\mathbf{v}^\top \mathbf{A} \mathbf{v}$ , and  $\mathbf{Q}_k$  is the model error covariance matrix at time  $t_k$ .

99 With sparse observations, the problem is more complex because in order to derive the surrogate model,  
 100 we need to estimate the true state. A rigorous Bayesian approach to this problem consists in extending eq. (3)  
 101 to include the system trajectory  $\mathbf{x}_0, \dots, \mathbf{x}_{N_t}$  in the control variables (Hsieh and Tang, 1998; Abarbanel et al.,  
 102 2018; Bocquet et al., 2019, 2020a). The joint cost function reads

$$103 \quad \mathcal{J}(\mathbf{p}, \mathbf{x}_0, \dots, \mathbf{x}_{N_t}) \triangleq \mathcal{L}(\mathbf{p}, \mathbf{x}_0) + \frac{1}{2} \sum_{k=0}^{N_t-1} \|\mathbf{x}_{k+1} - \mathcal{M}_{\Delta t}(\mathbf{p}, \mathbf{x}_k)\|_{\mathbf{Q}_k^{-1}}^2 + \frac{1}{2} \sum_{k=0}^{N_t} \|\mathbf{y}_k - \mathcal{H}(\mathbf{x}_k)\|_{\mathbf{R}^{-1}}^2, \quad (4)$$

104 where  $\mathcal{L}$  is a regularisation term on both  $\mathbf{p}$  and  $\mathbf{x}_0$ . The second term in eq. (4) corresponds to the second  
 105 term in eq. (3), in which the observations have been replaced with the state, and the third term is the  
 106 observation error term. Equation (4) is overall very similar to a typical weak-constraint (WC) 4D-Var cost  
 107 function (Trémolet, 2006).

108 Because the size of the trajectory control vector  $N_t \times N_x$  is likely to be large, an efficient minimisation  
 109 method relies on a coordinate descent technique, alternating DA steps to estimate the state with ML steps to  
 110 estimate the surrogate model (Brajard et al., 2020; Bocquet et al., 2020a). This combined DA-ML method,

111 illustrated in fig. 1, explicitly exploits the different nature between the arguments of  $\mathcal{J}$  (state of the system  
 112 and surrogate model parameters) and is highly flexible since the DA and ML steps are independent. A  
 113 comprehensive description of the DA-ML method is given by Bocquet et al. (2020a).

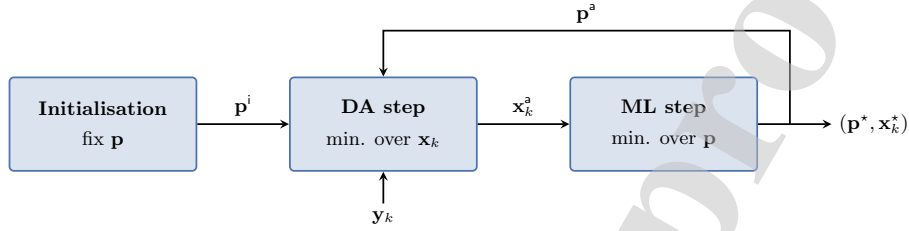


Figure 1: Illustration of the DA-ML method for the minimisation strategy of the cost function, eq. (4): alternate DA steps with ML steps to estimate the model parameters  $\mathbf{p}$  and the state trajectory  $\mathbf{x}_0, \dots, \mathbf{x}_{N_t}$  with an increasing accuracy.

114 The DA-ML method has first been used for full model emulation, *e.g.* by Brajard et al. (2020). In their  
 115 example, the surrogate model is a neural network (NN) which represents the model tendencies. It is combined  
 116 with an integration scheme to define the resolvent between two time steps. In this case,  $\mathbf{p}$  corresponds to  
 117 the set of weights and biases of the NN. The method has then been used to correct an imperfect physical  
 118 model by Brajard et al. (2021); Farchi et al. (2021). For this problem, the formalism is simply obtained by  
 119 replacing the resolvent of the surrogate model  $\mathcal{M}_{\Delta t}$  with the resolvent of the corrected model, in particular  
 120 in eq. (4). In general, model error correction should be an easier inference problem than full model emulation,  
 121 which means that smaller surrogate models (smaller in number of parameters) and less training data are  
 122 necessary. Moreover, using a physical model is likely to be beneficial to the method, in particular during the  
 123 DA steps. It also solves the issue of the initialisation: the first step of the method is to perform DA with  
 124 the (non-corrected) physical model. The advantages of model error correction over full model emulation are  
 125 further investigated in Farchi et al. (2021).

## 126 2.2. A typical geophysical model architecture

127 In the present work, we investigate model error correction with the DA-ML method. Before we introduce  
 128 any model error correction, we need to discuss the characteristics of the model to correct. The geophysical  
 129 models rely on physical laws, which most of the time take the form of ODEs or PDEs.

130 The core of a model is a numerical code computing the model *tendencies*  $\phi$ , which are defined as a  
 131 discretised version of the differential equations in  $\mathbb{R}^{N_x}$ :

$$132 \quad \phi(\mathbf{x}) \triangleq \frac{d\mathbf{x}}{dt}. \quad (5)$$

133 The model tendencies are integrated over time step  $\delta t$  using a dedicated integration scheme, for example the

134 explicit Euler scheme:

$$135 \quad \mathcal{I}(\mathbf{x}) \triangleq \mathbf{x} + \delta t \cdot \phi(\mathbf{x}), \quad (6)$$

136 or more elaborate schemes such as Runge–Kutta methods. Finally, several integration steps are composed to  
137 define the *resolvent*<sup>1</sup> from one time step to the next:

$$138 \quad \mathcal{M}_{\Delta t}(\mathbf{x}) \triangleq \mathcal{I} \circ \dots \circ \mathcal{I}(\mathbf{x}). \quad (7)$$

139 Two strategies can be used to correct such physical model. The first one is to include a correction in the  
140 resolvent, eq. (7). This is called resolvent correction (RC). The other strategy is to include the correction  
141 directly in the differential equations, in other words in the model tendencies, eq. (5) (Bocquet et al., 2019).  
142 This is called tendency correction (TC). According to Farchi et al. (2021), both strategies have advantages  
143 and drawbacks. Let us illustrate the difference using a simple univariate example.

### 144 2.3. Resolvent or tendency correction in a simple univariate example

145 Suppose that we follow the evolution of a process  $x \in \mathbb{R}$  over two  $\delta t$ -integration steps with the explicit  
146 Euler scheme. The true two-step resolvent is given by

$$147 \quad \mathcal{M}_2^t(x) = x + \delta t \cdot f(x) + \delta t \cdot f\{x + \delta t \cdot f(x)\}, \quad (8)$$

148 where  $f$  represents the true model tendencies. Our imperfect physical model has tendencies  $g$  and a two-step  
149 resolvent given by

$$150 \quad \mathcal{M}_2^p(x) = x + \delta t \cdot g(x) + \delta t \cdot g\{x + \delta t \cdot g(x)\}. \quad (9)$$

151 To simplify the expressions, in the following we take  $\delta t = 1$ .

152 When using a TC, we assume that the corrected model has tendencies  $g + \alpha$ , whereas when using RC, we  
153 assume that the two-step resolvent of the corrected model is  $\mathcal{M}_2^p + \beta$ . The optimal  $\alpha$  and  $\beta$  corrections are  
154 given by

$$155 \quad \alpha^*(x) = f(x) - g(x), \quad (10)$$

$$156 \quad \beta^*(x) = \mathcal{M}_2^t(x) - \mathcal{M}_2^p(x) \quad (11)$$

$$157 \quad = f(x) - g(x) + f\{x + f(x)\} - g\{x + g(x)\}, \quad (12)$$

159 where the difference is highlighted in red. Obviously, the optimal  $\beta$  is likely to be more complex than the  
160 optimal  $\alpha$ . The expression suggests that it will also be more nonlinear if  $f$  or  $g$  (or both) are nonlinear.

---

<sup>1</sup>The term *resolvent* is usual in the context of integral or differential equations. The same operator is often called *flow*, or *flow map* in dynamical systems and *propagator* in theoretical physics.

161 To further understand the difference, we derive the two-step resolvent with RC and TC, respectively  
 162 written  $\mathcal{M}_2^\beta$  and  $\mathcal{M}_2^\alpha$ :

$$163 \quad \mathcal{M}_2^\beta(x) = x + g(x) + g\{x + g(x)\} + \beta(x) \quad (13)$$

$$164 \quad \mathcal{M}_2^\alpha(x) = x + g(x) + g\{x + g(x) + \alpha(x)\} + \alpha(x) + \alpha\{x + g(x) + \alpha(x)\}, \quad (14)$$

166 where the difference is highlighted in red. From this perspective, it is clear that with TC,  $\mathcal{M}_2^\alpha(x)$  is marked  
 167 by the interaction between the physical model and the correction term  $\alpha$ . While this interaction is beneficial  
 168 because it enhances  $\mathcal{M}_2^\alpha(x)$ , the downside is that inferring  $\alpha$  from data is technically more difficult than  
 169 inferring  $\beta$ . Let us see why.

170 Suppose that both  $\alpha$  and  $\beta$  depend on a coefficient  $p \in \mathbb{R}$ . Observation data usually come in the form of  
 171 pairs  $(x_0, x_2)$  with  $x_2 = \mathcal{M}_2^\dagger(x_0)$ , possibly with some observation noise. Therefore, a learning step based on  
 172 some kind of gradient descent would require the gradient of the *corrected* two-step resolvent with respect to  
 173  $p$ , which is given by

$$174 \quad \frac{\partial \mathcal{M}_2^\beta}{\partial p}(x) = \frac{\partial \beta}{\partial p}(x), \quad (15)$$

$$175 \quad \frac{\partial \mathcal{M}_2^\alpha}{\partial p}(x) = \frac{\partial \alpha}{\partial p}(x) \cdot \left[ 1 + g'\{x + g(x) + \alpha(x)\} + \frac{\partial \alpha}{\partial p}\{x + g(x) + \alpha(x)\} \right], \quad (16)$$

177 where the difference is once again highlighted in red. In particular, it depends on  $g'$ , the derivative of  $g$ .  
 178 The equivalent for a geophysical numerical model would be the tangent linear (TL) operator, which may be  
 179 difficult to compute.

180 To summarise, compared to RC, TC is more difficult to program because the correction term ( $\alpha$  in the  
 181 present example) is intrusive, meaning that it requires to modify deeply the code of the physical model. It is  
 182 also more difficult to train, as illustrated by the difference between eq. (15) and eq. (16). On the other hand,  
 183 once it is implemented, the TC has the potential to yield richer dynamics through the interaction with the  
 184 physical model. Furthermore, by construction the RC can only correct the two-step resolvent, while the TC  
 185 can also correct the one-step resolvent. This would make a difference when using the corrected model in a  
 186 DA experiment with observations at every step, because then the one-step resolvent is explicitly needed. The  
 187 simplest workaround is to assume a linear growth of errors in time (Brajard et al., 2021; Farchi et al., 2021).  
 188 In this case, the one-step resolvent with RC would be given by

$$189 \quad \mathcal{M}_1^\beta(x) = x + g(x) + \frac{1}{2}\beta(x), \quad (17)$$

190 since  $\beta$  is the correction term for the two-step resolvent  $\mathcal{M}_2^\beta(x)$ . However, even with the optimal  $\beta$  correction  
 191 from eq. (12), this one-step resolvent would still differ from the true one-step resolvent, given by

$$192 \quad \mathcal{M}_1^\dagger(x) = x + f(x). \quad (18)$$



193 In their experiments, Farchi et al. (2021) found that this hypothesis was the main limitation for improving  
194 the accuracy of DA experiments with the corrected model. They concluded that the best strategy to correct  
195 a model to be used in DA experiments would probably be the TC.

196 Finally, let us mention that the model error correction considered in this section is autonomous and  
197 additive. The autonomous hypothesis can be relaxed, for example by including time in the set of predictors.  
198 However, one must keep in mind that in this case, the training dataset should capture the time evolution of  
199 the model error. The additive hypothesis can also be relaxed. Without prior knowledge on the model error  
200 form, using an additive correction is the simpler option but other choices are possible, *e.g.* a multiplicative  
201 correction. Also note that if the physical model explicitly depends a set of parameters, the same framework  
202 can be used to calibrate these parameters.

#### 203 2.4. Comparing resolvent and tendency correction

204 Farchi et al. (2021) chose to focus on the RC because it is easier to implement. In the following section,  
205 we illustrate the difference between RC and TC. First hints in favour of the TC approach were gained from  
206 the comparison of the results of Bocquet et al. (2019) and of Brajard et al. (2020) on the Lorenz 40-variable  
207 model. To address the inference problem, we use the combined DA-ML method described in section 2.1. We  
208 start by a DA step with an imperfect physical model to assimilate the observations. We then use a ML step  
209 to train a model error correction from the analysis of the DA step. Pushing the DA-ML method further,  
210 we could iterate in place: use the corrected model to get a more accurate analysis in further DA steps and  
211 learn from this more accurate analysis to get an improved model error correction in further ML steps, as  
212 illustrated by fig. 1.

213 However, we choose to stop after the first DA-ML cycle for two reasons. First, DA experiments with  
214 realistic models are numerically expensive, and it may not be realistic to perform more than one DA step  
215 if the size of the trajectory  $N_t$  is large. It is also worth noting that operational centres usually compute  
216 *reanalyses*, which means that the first DA step is a product which is likely to be already available (Hersbach  
217 et al., 2020). Second, if the physical model (without correction) is reasonably accurate, the analysis of the  
218 first DA step should be reasonably accurate and hence the improvement of the first ML step should be much  
219 larger than the improvement of further ML steps. However, even though we stop after the first DA-ML cycle,  
220 we perform a second DA step, but only for evaluation purposes.

221 Finally, we emphasise again the offline nature of the DA-ML method previously discussed. As is, the  
222 ML step starts only when the DA analysis is available, *i.e.* once the entire observation dataset has been  
223 assimilated in the DA step. An alternative, online approach is proposed in section 4.

### 224 3. Numerical illustration with the two-scale Lorenz model (part I)

#### 225 3.1. Models description

226 In our experiments, the true model is the L05III model, which describes the evolution of two sets of  
 227 variables: the slow variables  $x_n$  for  $n \in \{1, \dots, N_x\}$  and the fast variables  $u_m$  for  $m \in \{1, \dots, N_x \times N_u\}$ .  
 228 These two-scale dynamics are given by

$$229 \quad \frac{dx_n}{dt} = x_{n-1}(x_{n+1} - x_{n-2}) - x_n + F - \frac{hc}{b} \sum_{m=1}^{N_u} u_{m+(n-1)N_u}, \quad (19a)$$

$$230 \quad \frac{du_m}{dt} = \frac{c}{b} \{b^2 u_{m+1}(u_{m-1} - u_{m+2}) - bu_m\} + \frac{hc}{b} x_{1+(m-1)N_u}, \quad (19b)$$

232 where  $//$  is the integer division and where the indices are applied periodically:  $x_{N_x+n} = x_n$  and  
 233  $u_{N_x \times N_u + m} = u_m$ . The idea is that each slow variable  $x_n$  is coupled to the  $N_u$  fast variables  $u_m$  for  
 234  $m \in \{1 + (n-1)N_u, \dots, nN_u\}$ .

235 A first order approximation of the L05III model is the one-scale Lorenz model (L96, Lorenz and Emanuel,  
 236 1998), which only describes the evolution of the slow variables  $x_n$ . The model is defined by

$$237 \quad \frac{dx_n}{dt} = x_{n-1}(x_{n+1} - x_{n-2}) - x_n + F, \quad (20)$$

238 where the indices once again apply periodically:  $x_{N_x+n} = x_n$ . This model is used in our experiments as the  
 239 (imperfect) physical model to correct.

240 Both L05III and L96 models are integrated using a fourth-order Runge–Kutta scheme, and the parameter  
 241 values are reported in table 1. With this setup, the true model dynamics is chaotic, with a leading Lyapunov  
 242 exponent of 1.3775 (Mitchell and Carrassi, 2015) and the model variability, defined as the standard deviation  
 243 of the climatological distribution of the state, averaged over the slow variables, is 3.5372. When using the  
 244 L96 model in place of the L05III model, two sources of model error are introduced:

- 245 1. the fast variables  $u_m$  generate unresolved processes;
- 246 2. the integration time step  $\delta t$  is 0.05 instead of 0.005.

247 Moreover, even though the forcing coefficient  $F$  differs in both models, this cannot strictly be considered as a  
 248 third source of model error as  $F = 10$  is chosen for the L05III model to better match the dynamics of the  
 249 L96 model with  $F = 8$ .

250 The accuracy of the physical (L96) model in reproducing the dynamics of the true (L05III) model is  
 251 measured using the forecast skill (FS) defined as the average root-mean-squared error (RMSE) of the forecast  
 252 after a given lead time:

$$253 \quad \text{FS}(k\delta t) \triangleq \frac{1}{N_e} \sum_{i=1}^{N_e} \text{RMSE} [\mathbf{\Pi} \circ \mathcal{M}_{k\delta t}^t(\mathbf{x}_i, \mathbf{u}_i), \mathcal{M}_{k\delta t}^p(\mathbf{x}_i)]. \quad (21)$$

Table 1: Parametrisation for the true (L05III) and physical (L96) models.

Parameter	Symbol	L05III	L96
Number of slow variables	$N_x$	36	36
Number of fast variables per slow variable	$N_u$	10	
Forcing	$F$	10	8
Coupling	$h$	1	
Time-scale ratio	$c$	10	
Space-scale ratio	$b$	10	
Integration time step	$\delta t$	0.005	0.05

254 In this equation,  $\mathcal{M}_{k\delta t}^t$  and  $\mathcal{M}_{k\delta t}^p$  are the resolvents of the true and physical models for a  $k\delta t$  integration,  
 255 respectively,  $\Pi$  is the projection operator onto the set of slow variables  $\Pi(\mathbf{x}, \mathbf{u}) = \mathbf{x}$ , and  $(\mathbf{x}_i, \mathbf{u}_i)$  for  
 256  $i \in \{1, \dots, N_e\}$  is a set of  $N_e$  initial conditions representative of the true model climatology. The FS,  
 257 normalised by the model variability, is shown in fig. 2a and illustrates the poor accuracy of the physical  
 258 model. In the following sections, we will see how model error corrections can be used to improve the FS, but  
 259 one must keep in mind that there is an intrinsic limit to potential improvements, because it is presumably  
 260 impossible to exactly reproduce the dynamics of the true model with only  $N_x = 36$  variables.

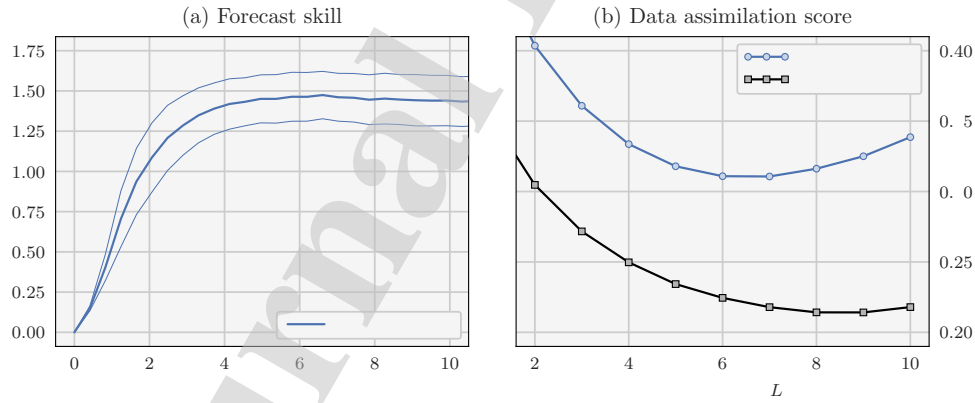


Figure 2: Left panel (a): forecast skill of the physical model (in units of the model variability) as a function of the lead time (in units of the Lyapunov time). The thick line shows the average over the  $N_e = 1024$  initial conditions and the thin lines indicate plus or minus one standard deviation. Right panel (b): accuracy of the DA step as a function of the length of the DAW  $L$  with the physical model (in blue) and with the true model (in black). The sRMSE is averaged over at least  $8192//L$  cycles after a spin-up period of at least  $1024//L$  cycles, and over 16 repetitions of each experiment. For each value of  $L$ ,  $b$  is optimally tuned to yield the lowest sRMSE.

261 3.2. Data assimilation with the physical model

262 The first step of the DA-ML method is to perform DA with the physical model. The truth  $(\mathbf{x}_k^t, \mathbf{u}_k^t)$  is  
 263 generated using the true model. Observations are taken every  $\Delta t = 0.05$  from the slow variables only, using

$$264 \quad \mathbf{y}_k = \mathbf{x}_k^t + \mathbf{v}_k, \quad \mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (22)$$

265 In other words, the observation operator is  $\mathcal{H} = \mathbf{I}$ , the observations are not biased, and the observation error  
 266 covariance matrix is  $\mathbf{R} = \mathbf{I}$ . Numerical illustrations with sparse observation operators are provided, *e.g.*, by  
 267 Brajard et al. (2020); Bocquet et al. (2020a); Brajard et al. (2021); Farchi et al. (2021). Except in the case  
 268 of very sparse observations, the results are qualitatively similar over a wide range of observation density,  
 269 which is why, for the present study, we have chosen to use dense observations for simplicity.

270 As explained in section 2.1, the goal of the DA step is to minimise eq. (4) with respect to the state  
 271 trajectory  $\mathbf{x}_0, \dots, \mathbf{x}_{N_t}$ , which is a WC 4D-Var problem. However, solving a WC minimisation is probably  
 272 unaffordable for large trajectories, as discussed by Bocquet et al. (2020a). To overcome this issue, we choose  
 273 to assimilate the observations using the cycled strong-constraint (SC) 4D-Var algorithm, with consecutive  
 274 DA windows (DAWs) of  $L$  batches of observations. This will provide an approximate solution to the WC  
 275 4D-Var problem. More specifically, each 4D-Var problem consists in minimising the cost function

$$276 \quad \mathcal{J}(\mathbf{x}_k) = \frac{1}{2} \|\mathbf{x}_k - \mathbf{x}_k^b\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2} \sum_{l=0}^{L-1} \|\mathbf{y}_{k+l} - \mathcal{M}_{l\Delta t}^p(\mathbf{x}_k)\|_{\mathbf{R}^{-1}}^2, \quad (23)$$

277 where  $\mathbf{x}_k^b$  is the background,  $\mathbf{B}$  is the background error covariance matrix,  $\{\mathbf{y}_k, \dots, \mathbf{y}_{k+L-1}\}$  is the set of  
 278 assimilated observations, and  $\mathcal{M}_{l\Delta t}^p$  is the resolvent of the physical model for an integration of  $l\Delta t$ . The  
 279 analysis is performed at time  $t_k$  (the time of the first batch of assimilated observations) and, in this cycled  
 280 context, it is used to obtain the background for the next analysis which is performed at time  $t_{k+L}$ , using

$$281 \quad \mathbf{x}_{k+L}^b = \mathcal{M}_{L\Delta t}^p(\mathbf{x}_k^a). \quad (24)$$

282 For the first cycle, the background state is obtained by perturbing the truth:

$$283 \quad \mathbf{x}_0^b = \mathbf{x}_0^t + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (25)$$

284 Finally, the background error covariance matrix  $\mathbf{B}$  is set to  $b^2\mathbf{I}$ , where  $b$  is an algorithmic parameter to  
 285 specify.

286 At each cycle, the cost function  $\mathcal{J}$ , eq. (23), is minimised using the L-BFGS algorithm (Byrd et al.,  
 287 1995), a quasi-Newton minimisation algorithm. The gradient of  $\mathcal{J}$  is computed exactly using automatic  
 288 differentiation, and the starting point of the minimisation is  $\mathbf{x}_k^b$ . The accuracy of the DA step is measured  
 289 using the RMSE of the analysis (analysis minus truth) at the start of the DAW, hereafter called the *smoothing*  
 290 RMSE (sRMSE), averaged over a sufficiently large number of cycles to ensure the convergence of the statistical  
 291 indicators.

292 In order to choose an appropriate value for the length of the DAW  $L$ , we first study the evolution of the  
 293 sRMSE as a function of  $L$ . The results are shown in fig. 2b. As expected, the sRMSE starts by decreasing  
 294 with  $L$ . It reaches an optimum for  $L = 6$ , and then increases with  $L$  as the impact of model error grows. For  
 295 comparison, fig. 2 also shows the results when using the true model in place of the physical model. Note that  
 296 in this case the 4D-Var cost function  $\mathcal{J}$ , eq. (23), depends on both the slow and the fast variables  $\mathbf{x}_k$  and  
 297  $\mathbf{u}_k$ . The evolution of the sRMSE as a function of  $L$  is very similar, with the exception that the scores are  
 298 overall much lower, and that the sRMSE increase for large values of  $L$  does not come from model error but  
 299 from optimisation issues. Indeed, for long DAWs, the cost function  $\mathcal{J}$  is likely to have several local minima,  
 300 which would make the L-BFGS algorithm not suited for the minimisation. Using a quasi-static formulation  
 301 of 4D-Var could mitigate this issue (Pires et al., 1996; Fillion et al., 2018).

### 302 3.3. Model error correction with a univariate polynomial regression

303 The present model error setup, as described in section 3.1, has already been addressed outside the scope  
 304 of ML, for example by Wilks (2005). The idea is to replace the physical model tendencies, eq. (20), by

$$305 \frac{dx_n}{dt} = x_{n-1}(x_{n+1} - x_{n-2}) - x_n + F + g(x_n), \quad (26)$$

306 where  $g$  is a univariate fourth-order polynomial correction, shared between all  $N_x = 36$  slow variables. The  
 307 five coefficients of  $g$  are computed using a least-square regression of the difference between eq. (20) and the  
 308 empirical tendencies

$$309 \frac{x_n^t(t + \delta t) - x_n^t(t)}{\delta t} \quad (27)$$

310 computed from a trajectory  $\mathbf{x}^t(t)$  of the true model.

311 Offering a baseline score for later comparison, fig. 3 shows the FS and the DA score for the model with  
 312 the polynomial regression  $g$ . For this illustration, following the approach of Wilks (2005), the coefficients of  
 313  $g$  are computed using 2000 pairs of snapshots ( $\mathbf{x}^t(t)$ ,  $\mathbf{x}^t(t + \delta t)$ ) with  $\delta t = 0.005$ , the integration time step  
 314 of the true model. The time interval between two consecutive pairs of snapshots is set to 1000 integration  
 315 steps. The results show that this simple correction is effective, both in forecast and DA experiments. In  
 316 particular, the DA score is very close to the one obtained with the true model. It is even better for  $L \geq 10$ .  
 317 This probably comes from the fact that a small amount of model error regularises the cost function eq. (23)  
 318 and mitigates the numerical issues discussed at the end of section 3.2.

319 Of course, this method cannot be applied to realistic models because the regression requires the truth, at  
 320 a very high frequency ( $\delta t = 0.005$ ). We have checked that using  $\delta t = 0.05$  (the integration time step of the  
 321 physical model) yields an ineffective correction. Nevertheless, it shows that the model error structure in this  
 322 setup can be effectively represented with a small number of parameters. The TC framework presented in  
 323 section 2 can be seen as a generalisation of the method of Wilks (2005) to (i) any kind of correction  $g$ , in

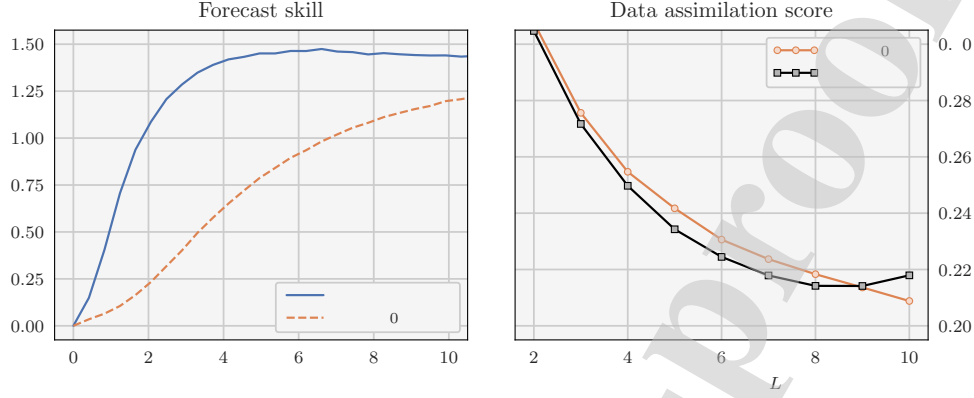


Figure 3: Same as fig. 2 for the physical model (in blue), the model with polynomial regression (in orange) and the true model (in black).

particular multivariate ones, and (ii) sparse and noisy observations for the training. A more complex but less scalable model error correction scheme has also been proposed for the same model by Pulido et al. (2018).

### 3.4. Resolvent and tendency correction with the DA-ML method

#### 3.4.1. The data assimilation step

Given the results of section 3.2, we start the DA-ML method by a long DA experiment with the physical model and with  $L = 6$ . At each cycle, we only keep the analysis at the start of the DAW. The result is a time series of analysis snapshots  $\mathbf{x}_{kL}^a$ , where the time interval between two snapshots is  $L\Delta t = 0.3$ . This trajectory is used to build the training dataset for the ML step. In other words, the surrogate models are trained to reproduce the map

$$\mathbf{x}_{kL}^a \mapsto \mathbf{x}_{(k+1)L}^a. \quad (28)$$

Another trajectory, resulting from a distinct long DA experiment, is used to build the validation dataset. Finally, since the ultimate goal is to predict the true dynamics and not the dynamics of the analysis snapshots, we compute an additional trajectory, this time with the true model. This third trajectory  $\mathbf{x}_{kL}^t$  is used to build the test dataset, and hence to evaluate the ability of the surrogate models to reproduce the map

$$\mathbf{x}_{kL}^t \mapsto \mathbf{x}_{(k+1)L}^t. \quad (29)$$

#### 3.4.2. Designing the surrogate models

The second step of the DA-ML consists in defining and training a surrogate model with the analysis of the first DA step. In this section, three different surrogate models are tested to correct the physical model. All three of them are autonomous and use NNs. For the first surrogate, the correction is computed using a

343 NN called CNN-a and then *added to the resolvent* of the physical model, following the RC approach, which  
 344 yields

$$345 \quad \mathcal{M}_{L\Delta t}^a(\mathbf{p}, \mathbf{x}) \triangleq \mathcal{M}_{L\Delta t}^p(\mathbf{x}) + \mathcal{F}^a(\mathbf{p}, \mathbf{x}). \quad (30)$$

346 In this equation,  $\mathcal{M}_{L\Delta t}^p$  is the resolvent of the physical model for an integration of  $L\Delta t$  (one DAW),  $\mathcal{F}^a$  is  
 347 the map encoding CNN-a,  $\mathbf{p}$  is the set of parameters of CNN-a (the weights and biases of the NN), and  
 348  $\mathcal{M}_{L\Delta t}^a$  is the resolvent of the resulting surrogate model, called RC-CNN-a. For the second surrogate, the  
 349 correction is computed using a NN called CNN-b and then *added to the tendencies* of the physical model,  
 350 following the TC approach:

$$351 \quad \phi^b(\mathbf{p}, \mathbf{x}) \triangleq \phi^p(\mathbf{x}) + \mathcal{F}^b(\mathbf{p}, \mathbf{x}). \quad (31)$$

352 In this equation,  $\phi^p$  represents the physical model tendencies, given by eq. (20),  $\mathcal{F}^b$  is the function encoding  
 353 CNN-b,  $\mathbf{p}$  is the set of parameters of CNN-b, and  $\phi^b$  represents the tendencies of the resulting surrogate  
 354 model, called TC-CNN-b. To compute the resolvent of this model,  $\mathcal{M}_{L\Delta t}^b$ , we keep the integration scheme  
 355 and time step of the physical model. Finally, the third surrogate model, called TC-CNN-c, is similar to  
 356 TC-CNN-b with CNN-b replaced with another NN called CNN-c, which uses a different activation function.

357 As explained in section 3.3, the model error structure is not overly complex. For this reason, we want  
 358 to keep the NNs as simple as possible. We have experimented with several NNs configurations and have  
 359 selected the following sequential (or feed-forward) architecture with:

- 360 1. the input layer;
- 361 2. a sequence of convolutional layers;
- 362 3. a final convolutional layer as output layer (without activation).

363 All intermediate convolutional layers share the same number of filters, the same convolutional window, and  
 364 the same activation function. They also use periodic padding to preserve the input and output shape of the  
 365 layers. The last convolutional layer uses only one filter, a convolution window of only one variable, and no  
 366 activation function. The purpose of this layer is not to actually perform a convolution, but to project the  
 367 output of the previous layer to the output variables. The settings of the intermediate convolutional layers  
 368 are reported in table 2 for CNN-a, CNN-b, and CNN-c, alongside the total number of parameters.

### 369 3.4.3. Neural networks initialisation

370 When working with NNs, the parameter initialisation step is important. A common method is to use  
 371 random values for the initial weights and to set the initial biases to zero. The underlying idea is that there is  
 372 no reason for the optimal weights to display any specific symmetry. Because such symmetries are preserved  
 373 during the training, even with stochastic gradient descent, they need to be broken during the initialisation,  
 374 hence the use of random initial values (Goodfellow et al., 2016).

Table 2: Settings of the convolutional layers of the NNs used in the DA-ML method. The absence of an activation function is tantamount to a linear activation function.

Setting	CNN-a	CNN-b	CNN-c
Number of layers	4	1	1
Number of filters per layer	16	16	16
Size of the convolutional window	5	5	5
Activation function	tanh		tanh
Total number of parameters	4001	113	113

375 In our case however, the situation is different because the surrogate models are hybrid. Since the  
376 corrections are additive, all three surrogate models are equivalent to the physical model when  $\mathbf{p} = \mathbf{0}$ , and it  
377 is highly probable that a random  $\mathbf{p}$  would make the model predictions worse. Initialising the NNs with  $\mathbf{p} = \mathbf{0}$   
378 would hence make sense, but for the reasons aforementioned, this could yield suboptimal surrogate models  
379 after the training. Therefore, we initialise the NNs using the following approach, which we found to be a  
380 good compromise. The intermediate convolutional layers are initialised using the classical method in ML  
381 (random weights and zero biases) and the last convolutional layer is initialised to zero (both zero weights and  
382 zero biases). This approach is very similar to the ReZero method developed by Bachlechner et al. (2020).

#### 383 3.4.4. Training the surrogate models

384 We now start the ML step. The surrogate model parameters  $\mathbf{p}$  are optimised using the Adam algorithm,  
385 a variant of the stochastic gradient descent (Kingma and Ba, 2015). The loss function is the mean-squared  
386 error (MSE) over the training dataset, made of analysis snapshots. The training consists of 1024 epochs with  
387 a learning rate of  $1 \times 10^{-3}$  and a batch size of 32. After the entire training step, we keep the model which  
388 yields the lowest MSE over the validation dataset, also made of analysis snapshots. This is necessary since the  
389 cost functions of the NNs do not include any internal mechanism to mitigate overfitting (*e.g.* regularisation).  
390 Finally, we evaluate the trained model by computing the MSE over the test dataset, made of truth snapshots,  
391 hereafter called test MSE (tMSE). For comparison, we also train and evaluate the surrogate models using the  
392 exact same method but with snapshots from the truth (instead of the analysis) in the training and validation  
393 datasets. This is equivalent to using dense and noiseless observations.

394 Figure 4 shows the training results for datasets of increasing size. Note that, in order to build a dataset  
395 of size  $N_t$ , the total number of cycles (or DAWs) required in the preliminary DA step is  $2 \times (N_t + 1)$ :

- 396 •  $N_t + 1$  cycles to produce the  $N_t$  pairs of analysis snapshots  $(\mathbf{x}_{kL}^a, \mathbf{x}_{(k+1)L}^a)$  for the training dataset;
- 397 • and the same for the validation dataset.



398 Let us first discuss the training with the truth, because it shows the full potential of each model. First,  
 399 all three models do improve over the physical model and yield very low tMSEs, but the scores with TC  
 400 are significantly better than with RC. As explained in section 2.3, the models with TC benefit from the  
 401 interaction between the physical model and the correction term (the NN). This is even more important here  
 402 than in the univariate example of section 2.3, because here the number of interactions for a prediction is 24:  
 403 4 interactions per integration step (through the fourth-order Runge–Kutta scheme) times 6 integration steps  
 404 between two DAWs. Furthermore, the training dataset needs to be much larger to get an accurate model with  
 405 RC than with TC. This is consistent with the total number of trainable parameters for each model, reported  
 406 in table 1: 4001 for RC-CNN-a and only 113 for TC-CNN-b and TC-CNN-c. It is remarkable that with  
 407 a training dataset of only one pair of truth snapshots (the smallest possible training dataset), the models  
 408 with TC are already very accurate, more than RC-CNN-a trained with the largest dataset considered in the  
 409 present study! Obviously, this is only possible because the correction is autonomous. The overall accuracy of  
 410 TC-CNN-b is also remarkable, given the fact that the correction provided by CNN-b is linear. The only  
 411 difference between CNN-b and CNN-c is their activation function for the intermediate layers. This means  
 412 that the difference between TC-CNN-b and TC-CNN-c illustrates the nonlinearity of the error in the model  
 413 tendencies: weak but non-negligible. For comparison, we have checked that with RC, the accuracy of the  
 414 surrogate models built using linear NNs is not satisfactory. This shows that the nonlinearity of the dynamics  
 415 over one DAW, eq. (29), is significant and that estimating and correcting model error as a tendency forcing  
 416 is a good first-order strategy to mitigate the effects of nonlinearities. For completeness, we mention that  
 417 better scores could have been obtained with RC, for example with larger or deeper NNs. However, this would  
 418 increase the number of trainable parameter, which means that the training dataset should be even larger.

419 When training with the analysis, the accuracy of the surrogate models is much lower, which was expected,  
 420 but all three models are still able to improve over the physical model. Most of the conclusions hold: the  
 421 scores with TC are better than with RC and require much smaller datasets. However this time, the linear  
 422 TC-CNN-b outperforms the nonlinear TC-CNN-c. This is probably a sign that nonlinear NNs are harder to  
 423 train.

#### 424 3.4.5. Forecast skill of the surrogate models

425 The tMSE is a measure of the accuracy of a model for an integration of one DAW of  $L\Delta t = 0.3$ . In  
 426 this section, we measure the accuracy of the surrogate models for longer forecast, by computing the FS as  
 427 defined in section 3.1. In order not to penalise RC-CNN-a over TC-CNN-b and TC-CNN-c, we evaluate the  
 428 surrogate models which have been trained with the largest dataset.

429 The results are shown in fig. 5 and demonstrate that the models, trained for an integration of one  
 430 DAW, remain effective for much longer integrations. When training with the truth, the ranking of the three  
 431 surrogate models is clear and consistent with the tMSE results: TC-CNN-c is the most accurate, followed

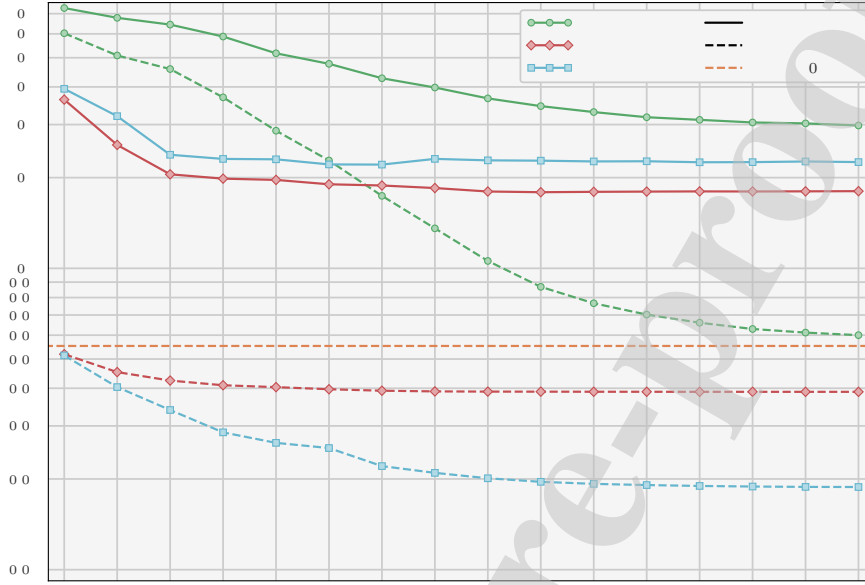


Figure 4: Evolution of the tMSE of the trained surrogate model as a function of the size of the training and validation datasets for RC-CNN-a (in green), TC-CNN-b (in red), and TC-CNN-c (in cyan). The tMSE is computed using a test dataset of size  $N_t = 8192$  and then normalised by the tMSE of the physical model (0.27785). Each experiment (initialisation, training and evaluation) is repeated 16 times with different training, validation, and test datasets and the final score is averaged over the 16 repetitions. The surrogate models are trained either with the analysis (continuous lines) or with the truth (dashed lines). For comparison, the horizontal dashed orange line indicates the score for the model with the polynomial regression of Wilks (2005).

432 by TC-CNN-b, and the least accurate is RC-CNN-a. When training with the analysis, the ranking of the  
 433 models for long integrations (*e.g.* longer than 6 DAWs) is the opposite of the ranking for the tMSE results:  
 434 RC-CNN-a is the most accurate, very closely followed by TC-CNN-c, and the least accurate is TC-CNN-b,  
 435 the only model built using a linear NN. There is a crossover in the forecast error curves after 2 or 3 DAWs,  
 436 with the errors of the linear NN (TC-CNN-b) becoming larger than those of the nonlinear NNs (RC-CNN-a  
 437 and TC-CNN-c). This issue is not specific to the L05III model nor to the use of nonlinear NNs since Farchi  
 438 et al. (2021) also faced the same kind of issues with a quasi-geostrophic model and with linear NNs. We do  
 439 not have yet a convincing explanation for this behaviour, which is specific to the use of the analysis for the  
 440 training, and understanding it better will require further work. In any case, we conclude that, when the  
 441 surrogate models are trained with the analysis, the accuracy for long forecasts is somewhat similar with RC  
 442 and with TC and requires nonlinear corrections.

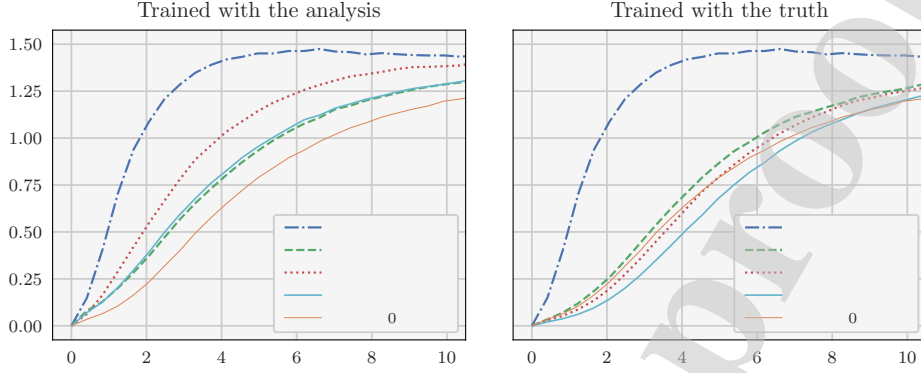


Figure 5: Same as fig. 2, left panel, for the physical model (in blue) and the trained surrogate models: RC-CNN-a (in green), TC-CNN-b (in red), and TC-CNN-c (in cyan). The surrogate models are trained either with the analysis (left panel) or with the truth (right panel). For comparison, the thin orange line indicates the scores for the model with the polynomial regression of Wilks (2005).

#### 3.4.6. Data assimilation experiments with the surrogate models

In this section, the goal is to measure the accuracy of the surrogate models in DA experiments. To do that, we reimplement the 4D-Var setup described in section 3.2 replacing the physical model by one of the surrogate models. In particular, in the cost function eq. (23),  $\mathcal{M}_{l\Delta t}^p$ , the resolvent of the physical model for an integration of  $l\Delta t$ , is replaced with  $\mathcal{M}_{l\Delta t}^a$ ,  $\mathcal{M}_{l\Delta t}^b$ , or  $\mathcal{M}_{l\Delta t}^c$ , the resolvents of RC-CNN-a, TC-CNN-b, or TC-CNN-c for an integration of  $l\Delta t$ . While the construction of  $\mathcal{M}_{l\Delta t}^b$  and  $\mathcal{M}_{l\Delta t}^c$  is similar to that of  $\mathcal{M}_{l\Delta t}^p$  because TC-CNN-b and TC-CNN-c use the TC method, the construction of  $\mathcal{M}_{l\Delta t}^a$  requires an additional assumption because RC-CNN-a uses the RC method. As discussed in section 2.3, the simplest assumption is the linear growth of errors in time, for which

$$\mathcal{M}_{\Delta t}^a \triangleq \mathcal{M}_{\Delta t}^p + \frac{1}{L} \mathcal{F}^a, \quad (32)$$

where  $\mathcal{F}^a$  is defined in section 3.4.2 as the function encoding CNN-a. This assumption is the standard assumption in the current implementation of WC 4D-Var (Laloyaux et al., 2020a,b) and it is the assumption we make for our DA experiments. Furthermore, for the same reason as above, we evaluate the models which have been trained with the largest dataset.

Following the approach of section 3.2, we study the evolution of the sRMSE as a function of the length of the DAW  $L$ . The results are shown in fig. 6 and demonstrate that the improved accuracy of the models also yields more accurate analyses. From these results, two elements could be highlighted. First, the sRMSE is much lower with TC-CNN-b and TC-CNN-c, which both use the TC method, than with RC-CNN-a, which

461 uses the RC method. Additionally, the sRMSE is lower when RC-CNN-a is trained with the analysis than  
 462 with the truth, even though the model error predictions are more accurate, as indicated by the lower tMSE.  
 463 These results confirm the hypothesis of Farchi et al. (2021) that the main obstacle to more accurate analyses  
 464 with the RC method is the assumption of linear growth of errors in time and that the TC method is more  
 465 appropriate for DA experiments where the impact of nonlinearities is significant. Second, when trained with  
 466 the truth, the sRMSE obtained with the TC method is of the same order as that obtained with the true  
 467 model, it is even better for TC-CNN-c! For large windows, typically  $L \geq 8$ , a fraction of the improvement  
 468 comes from the numerical issues with the true model discussed at the end of section 3.2. For smaller windows  
 469 however, we believe that the remaining improvement shows that TC-CNN-c may capture other deficiencies  
 470 than model error.

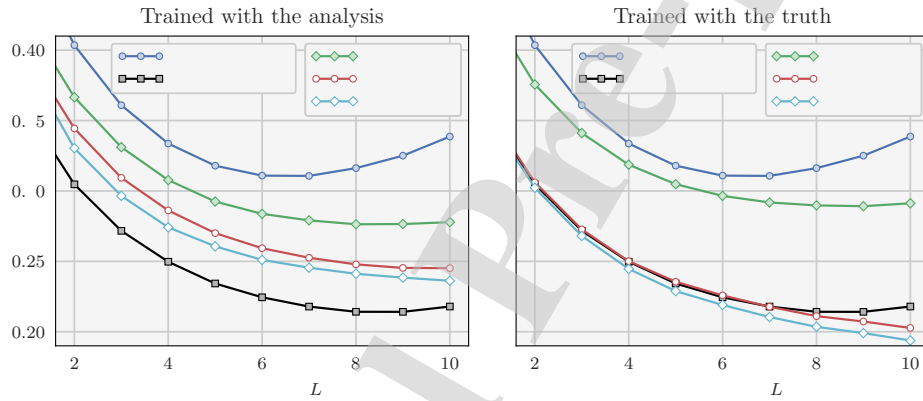


Figure 6: Same as fig. 2, right panel, for the physical model (in blue), the true model (in black), and the trained surrogate models: RC-CNN-a (in green), TC-CNN-b (in red), and TC-CNN-c (in cyan). The surrogate models are trained either with the analysis (left panel) or with the truth (right panel).

#### 471 3.4.7. Learning from less frequent snapshots

472 Because the true model is chaotic, forecasts are highly sensitive to the initial condition. In addition, the  
 473 accuracy of forecasts at longer lead times is generally more sensitive to incorrect model parameters, while  
 474 the accuracy of the training dataset is unchanged. This is beneficial when training the surrogate model with  
 475 the analysis since it reduces the impact of the analysis error in the training dataset. On the other hand,  
 476 longer forecast are inherently harder to predict. Therefore, we expect to see a trade-off between both effects  
 477 when increasing the length of the forecasts.

478 By construction, the surrogate models are trained to emulate the dynamics over one DAW, eq. (28).  
 479 Changing the DAW length  $L$  is not optimal, because  $L = 6$  minimises the sRMSE and hence it ensures the

480 most accurate analysis on average. Another possibility is to keep  $L = 6$  and train the surrogate models to  
 481 emulate the dynamics over multiple DAWs, *i.e.* to reproduce the map

$$482 \quad \mathbf{x}_{kL}^a \mapsto \mathbf{x}_{(k+N)L}^a, \quad (33)$$

483 where  $N$  is the number of DAWs over which the models should emulate the dynamics. This technique has  
 484 been used by Farchi et al. (2021) in the context of the RC method, with only partial success. Indeed, using  
 485  $N > 1$  systematically worsen the DA scores, which is not a surprise because the assumption of linear growth  
 486 of errors in time is less and less valid for longer forecasts. In this section, we evaluate this technique with the  
 487 TC method, which does not suffer from the same limitations as the RC method (in particular it does not  
 488 require additional assumption for DA experiments).

489 Figure 7 shows the forecast skill and the DA score for TC-CNN-c trained with the analysis using  $N = 2$ .  
 490 For comparison, the scores for TC-CNN-c trained with the analysis and with the truth, both using  $N = 1$ ,  
 491 are taken from figs. 5 and 6 and reproduced here. These results confirm that increasing the length of the  
 492 forecasts to predict yields more accurate surrogate models, both for forecast and DA experiments. Moreover,  
 493 we have checked that the scores get worse when further increasing the length of the forecasts to  $N = 4$  (not  
 494 shown here). This indicates that the trade-off aforementioned reaches an optimum for  $N = 2$  or  $N = 3$ .

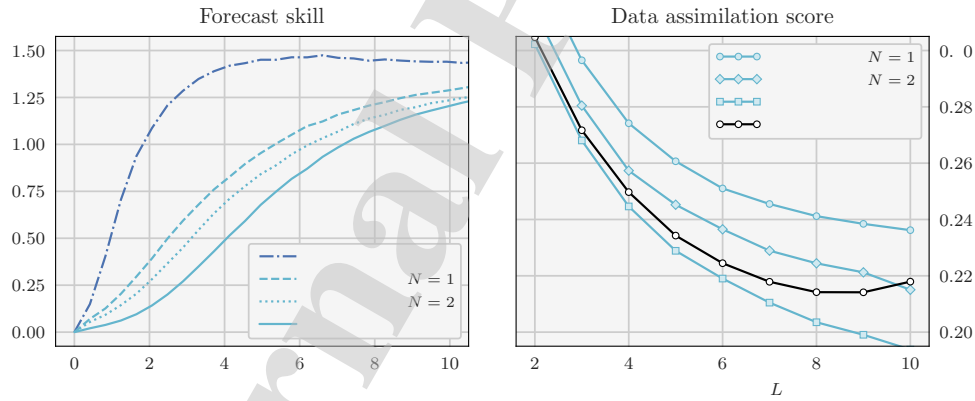


Figure 7: Same as fig. 2 for the physical model (in blue), the true model (in black), and the trained surrogate model TC-CNN-c (in cyan). The surrogate model is trained either with the analysis over one DAW ( $N = 1$ ), over two DAWs ( $N = 2$ ), or with the truth over one DAW ( $N = 1$ ).

#### 495 3.4.8. Iterating the DA-ML method

496 Throughout the previous experiments, we have seen different performances depending on whether the  
 497 surrogate models are trained with the analysis or with the truth, the latter case being equivalent to using

498 dense and noiseless observations. Iterating the DA-ML method, as originally proposed by Brajard et al.  
 499 (2020) and formalised by Bocquet et al. (2020a) as a coordinate descent optimisation, is a way to bring the  
 500 performances of the surrogate models trained with the analysis closer to that of the surrogate models trained  
 501 with the truth.

502 Given the results of the DA experiment with the RC method in section 3.4.6, it is possible that an  
 503 additional iteration of the DA-ML method will not improve the DA score by much. This is not the case with  
 504 the TC method because in this case, the analysis with the surrogate model is much more accurate than that  
 505 of the physical (non-corrected) model. However, for the TC method we would like to show an alternative  
 506 approach in the next section.

#### 507 4. Online learning of model error with tendency correction

##### 508 4.1. From offline to online learning

509 As already mentioned, the DA-ML method presented in section 2 and illustrated in section 3 is an offline  
 510 learning method, meaning that the training starts only once all observations have been assimilated. This  
 511 makes the method simple and flexible because the ML and DA steps are independent from each other. As a  
 512 consequence, it is probably easier to extract all the information from the analysis during the ML step.

513 On the other hand, using an offline approach also has drawbacks. As suggested in section 3.4.8, the  
 514 DA-ML method needs to be iterated to give the best performance. At each iteration, the entire set of  
 515 observations has to be re-assimilated, which can be problematic when the DA step is numerically expensive  
 516 (for example with a realistic model). This is particularly concerning in the case of an operational model for  
 517 which a new correction must be trained each time the model gets updated. Such issues does not affect online  
 518 approaches, since the training could start as soon as the first observation arrives.

519 In our context, online learning means that, each time a new batch of observation becomes available, we  
 520 have to estimate both the state of the system and the surrogate model at the same time. To address this  
 521 problem, the most natural approach is to use the formalism of DA with an augmented state containing the  
 522 current state of the system  $\mathbf{x}$  and the parameters of the surrogate model  $\mathbf{p}$ , following the principle formulated  
 523 by Jazwinski (1970). The resulting inference problem shares many aspects with classical parameter estimation  
 524 in DA (Ruiz et al., 2013).

##### 525 4.2. A new formulation of weak-constraint 4D-Var

526 Following the approach described in section 3.2, the observations are assimilated using the 4D-Var  
 527 algorithm, with consecutive windows of  $L$  batches of observations. Since the algorithm now has to estimate  
 528 the model parameters in addition to the model state, each 4D-Var problem consists in minimising the cost

529 function

$$530 \quad \mathcal{J}(\mathbf{p}_k, \mathbf{x}_k) = \frac{1}{2} \|\mathbf{p}_k - \mathbf{p}_k^b\|_{\mathbf{B}_p}^2 + \frac{1}{2} \|\mathbf{x}_k - \mathbf{x}_k^b\|_{\mathbf{B}_x}^2 + \frac{1}{2} \sum_{l=0}^{L-1} \|\mathbf{y}_{k+l} - \mathcal{H} \circ \mathcal{M}_{l\Delta t}(\mathbf{p}_k, \mathbf{x}_k)\|_{\mathbf{R}}^2, \quad (34)$$

531 where  $\mathbf{p}_k^b$  and  $\mathbf{B}_p$  are the background and background error covariance matrix for model parameters,  
 532 respectively. Comparing eq. (34) to the 4D-Var cost function without model parameters, eq. (23), two  
 533 differences can be highlighted. First, the resolvent of the physical model  $\mathcal{M}_{l\Delta t}^p$  is replaced with the resolvent  
 534 of the surrogate model  $\mathcal{M}_{l\Delta t}$ , which now depends on the model parameters. Second, a background (or  
 535 prior) term on model parameters has been added. This background term is very important because, in this  
 536 cycled context, it carries out the information on model parameters from one DAW to the next. Indeed, the  
 537 background on model parameters for the next DAW is equal to the analysis on model parameters of the  
 538 current DAW:

$$539 \quad \mathbf{p}_{k+L}^b = \mathbf{p}_k^a. \quad (35)$$

540 In other words, the forecast model for model parameters is the persistence, which is consistent with the fact  
 541 that the model is autonomous.

542 Like WC 4D-Var, eq. (34) can be inferred from Bayes' rule, but eq. (34) additionally neglects the  
 543 cross-covariances between the background errors for model state and for model parameters. Including  
 544 cross-covariances between model state and model parameters is possible but requires either a prior knowledge  
 545 or the use of an ensemble to estimate them, which is beyond the scope of the present work.

546 Our assimilation method is summarised in algorithm 1 in a cycled context. Rigorously speaking, this  
 547 could be seen as a SC method because it includes only one state vector in the control variables. However  
 548 with our method, by contrast with SC 4D-Var, the model can be updated during the analysis (when the  
 549 model parameters change). Therefore, considering a broader definition of WC methods, algorithm 1 can be  
 550 seen as a specific formulation of WC 4D-Var.

551 Optionally, the model parameters can be pre-trained. For completeness, we would like to mention the  
 552 similarities with the algorithms developed by Bocquet et al. (2020b); Malartic et al. (2021) which solve the  
 553 same kind of problems but in a filtering context (*i.e.* with  $L = 1$ ) with several variants of the ensemble  
 554 Kalman filter. Finally note that, just as the DA-ML method of section 2 can be used for model error  
 555 correction instead of full model emulation, the present formulation of WC 4D-Var can be used for model  
 556 error correction as well.

## 557 5. Numerical illustration with the two-scale Lorenz model (part II)

558 In this section, we illustrate the online learning method developed in section 4 using the same model error  
 559 setup as in section 3. The true model is the L05III model and the physical model is the L96 model. However,  
 560 because the online learning approach is based on the sole formalism of DA, we only use the surrogate models

---

**Algorithm 1** NN formulation of WC 4D-Var in a cycled context

---

**Parameters:** Observation operator  $\mathcal{H}$ , surrogate model  $\mathcal{M}$  with NN, background error covariance matrices for model state and model parameters  $\mathbf{B}_x$  and  $\mathbf{B}_p$ , observation error covariance matrix  $\mathbf{R}$ , DAW length  $L$

**Input:** Initial background for model state and model parameters  $\mathbf{x}_0^b$  and  $\mathbf{p}_0^b$ , observations  $\{\mathbf{y}_0, \dots, \mathbf{y}_{(N_t+1)L}\}$

- 1: **for**  $k = 0$  **to**  $N_t$  **do**
  - 2:   Compute the analysis  $\mathbf{p}_{kL}^a$  and  $\mathbf{x}_{kL}^a$  by minimising  $\mathcal{J}(\mathbf{p}_{kL}, \mathbf{x}_{kL})$ , eq. (34)   ▷ *analysis at time  $t_{kL}$*
  - 3:   Forecast the model parameters  $\mathbf{p}_{(k+1)L}^b = \mathbf{p}_{kL}^a$    ▷ *using persistence*
  - 4:   Forecast the model state  $\mathbf{x}_{(k+1)L}^b = \mathcal{M}_{L\Delta t}(\mathbf{p}_{kL}^a, \mathbf{x}_{kL}^a)$    ▷ *using the surrogate model*
  - 5: **end for**
  - 6: **return**  $(\mathbf{x}_{kL}^a, \mathbf{p}_{kL}^a)$  for  $k \in \{0, \dots, N_t\}$    ▷ *analysis trajectory*
- 

561 TC-CNN-b and TC-CNN-c, built with the TC method, since we have shown in section 3.4.6 that it is the  
562 most consistent and the most efficient for DA experiments.

### 563 5.1. Data assimilation setup

564 For this illustration, we take the same DA problem as in section 3.2. The truth is generated using the  
565 true model, and observations are taken every  $\Delta t = 0.05$  from the slow variables only, using eq. (22).

566 We start the experiments by assimilating the observations using the SC 4D-Var algorithm with the exact  
567 same setup as in section 3.2. In particular, we use the physical model, and we choose  $L = 6$  as it yields the  
568 best results with this model. After a total of 1024 DA cycles, which is enough to ensure the convergence of  
569 the analysis error statistics, we shift to the NN formulation of WC 4D-Var (algorithm 1). In other words, we  
570 switch on model error correction. For the following cycles, we keep  $L = 6$  and  $\mathbf{B}_x = b_x^2 \mathbf{I}$ , although  $b_x$  can be  
571 different than in the 1024 preliminary cycles. In addition, we set  $\mathbf{B}_p = b_p^2 \mathbf{I}$ , where  $b_p$  is another algorithmic  
572 parameter to specify. Finally for consistency, the initial background for model parameters  $\mathbf{p}_0^b$  is constructed  
573 using the NN initialisation method described in section 3.4.3.

574 With WC 4D-Var, at each cycle the cost function  $\mathcal{J}$ , eq. (34), is minimised using the same L-BFGS  
575 algorithm as for the SC 4D-Var. The gradient of  $\mathcal{J}$ , with respect to both model state and model parameters,  
576 is computed exactly using automatic differentiation, and the starting point of the minimisation is  $(\mathbf{p}_k^b, \mathbf{x}_k^b)$ .  
577 The accuracy of the analysis is measured using the RMSE *on model state* (analysis minus truth) at the start  
578 of the DAW, which corresponds to the sRMSE defined in section 3.2. In this cycled context, the sRMSE only  
579 improves if the model is getting more accurate. Nevertheless, we also measure the accuracy of the model by  
580 computing the tMSE defined in section 3.4.4.

581 Finally note that the 1024 preliminary cycles with the physical model are not mandatory. We have  
582 checked that the results are qualitatively similar without them. In fact, adding these preliminary cycle is a



583 way to get in real condition, where we need to correct a physical model which is already running since a  
 584 while. In the following section, we do not discuss the preliminary cycles.

### 585 5.2. Results with TC-CNN-b as surrogate model

586 Let us start by applying the method to TC-CNN-b. In other words, in this case  $\mathcal{M}_{l\Delta t}$  in eq. (34) is the  
 587 resolvent of TC-CNN-b for an integration of  $l\Delta t$ . For this experiment, we use the following algorithmic  
 588 parameters, which we have empirically found to yield good performances:

$$589 \quad b_x = 0.28 + 0.15 \times \exp(-t/256), \quad (36a)$$

$$590 \quad \widehat{b}_p = 0.001 + 0.1 \times \exp(-t/1024), \quad (36b)$$

$$591 \quad b_p = \min \left[ 0.05, \widehat{b}_p \right], \quad (36c)$$

592  
 593 where  $t$  is the time measured in number of DAWs after the 1024 preliminary cycles. The rationale behind  
 594 this choice is that the optimal values of  $b_x$  and  $b_p$  should be larger at the start of the experiment than at  
 595 the end (when the model is more accurate). To come up with eq. (36), we first chose the shape of  $b_x$  and  $b_p$   
 596 (exponential decay) and we then tuned the coefficients until we got satisfying results. Also note that we have  
 597 checked that the results are qualitatively similar when replacing the exponential decay of  $b_x$  and  $b_p$  with a  
 598 linear decay.

599 Figure 8 shows the time series of sRMSE and tMSE throughout the experiment, after the 1024 preliminary  
 600 cycles. For comparison, the scores obtained when TC-CNN-b is trained using the offline DA-ML approach,  
 601 both with the analysis and with the truth, are taken from figs. 4 and 6 and reproduced here. First of all, the  
 602 experiment is a success: the algorithm is working as expected and steadily improves the model, which can be  
 603 seen both in the sRMSE (DA score) and in the tMSE (forecast score). After 128 to 256 cycles the model  
 604 becomes more accurate than if trained offline with the analysis. Finally, the accuracy converges after 2048  
 605 to 4096 cycles. In the end, the model is almost as accurate as if trained offline with the truth! This is a  
 606 strong result because, as mentioned in section 3.4.4, training with the truth illustrates the full potential of a  
 607 surrogate model. This shows that our online learning method has been able to extract all the information  
 608 from the observations, and that we cannot expect better results with this surrogate model.

### 609 5.3. Results with TC-CNN-c as surrogate model

610 We now apply the method to TC-CNN-c. For this experiment, the algorithmic parameters, once again  
 611 chosen on empirical grounds, are slightly different:

$$612 \quad b_x = 0.26 + 0.20 \times \exp(-t/256), \quad (37a)$$

$$613 \quad \widehat{b}_p = 0.01 + 0.05 \times \exp(-t/3072), \quad (37b)$$

$$614 \quad b_p = \min \left[ 0.05, \widehat{b}_p \right], \quad (37c)$$

616 with  $t$  being once again the time measured in number of DAWs after the 1024 preliminary cycles.

617 Figure 9 shows the time series of sRMSE and tMSE throughout the experiment, after the 1024 preliminary  
 618 cycles. The success of the experiment is as clear as with TC-CNN-b: the algorithm steadily improves the  
 619 model, which can be seen both in the sRMSE and in the tMSE. After 128 to 256 cycles the model becomes  
 620 more accurate than if trained offline with the analysis. However, even though the total number of DA cycles  
 621 increases, the accuracy has not yet converged at the end of the experiment. One must keep in mind that  
 622 the correction provided by CNN-c is here nonlinear, contrary to the previous experiment with TC-CNN-b  
 623 where the correction provided by CNN-b is linear. Therefore, we think that the present experiment is a good  
 624 illustration of the increased complexity of training nonlinear NNs. Nevertheless, the accuracy of the model  
 625 after 8192 is remarkable and in particular the sRMSE is lower than when using the true model! We also think  
 626 that, if we were to extend the experiment with an appropriate tuning for  $b_p$ , the model would in the end be  
 627 almost as accurate as if trained offline with the truth, just as in the previous experiment with TC-CNN-b.

#### 628 5.4. Additional remarks on the online experiments

629 While preparing the online experiments, we have found that tuning the algorithmic parameter  $b_p$  is very  
 630 important. If  $b_p$  is too small, the algorithm gives too much weight to the background on model parameters  
 631  $\mathbf{p}_k^b$ . Even if this does not stop the learning process, it tapers the model parameter update, which makes the  
 632 convergence slower<sup>2</sup>. On the other hand if  $b_p$  is too large, the algorithm overfits the model parameters to the  
 633 observation window, which can yield divergence. As mentioned in section 5.2, what makes the tuning of  $b_p$   
 634 really complex here is that, as the model steadily improves, the optimal value of  $b_p$  decrease. The values  
 635 selected for the experiments, eqs. (36) and (37), have been chosen by *trial and error*. Even though they  
 636 yield good performances, they have not been optimally tuned. This means that a faster convergence could  
 637 have been most likely obtained with other values. Finally, note that the algorithmic parameter  $b_p$  here is  
 638 very similar to the tapering parameter introduced by Bocquet et al. (2020b); Malartic et al. (2021) in their  
 639 variants of the ensemble Kalman filter.

640 The online experiments use the zero/random NN initialisation method described in section 3.4.3. Therefore,  
 641 at the start of the WC 4D-Var cycles the surrogate model is equivalent to the physical model. Other  
 642 initialisation methods are possible. For example, one can use the set of parameters obtained after the offline  
 643 learning method of section 3.4. We have implemented this method (note illustrated here) and checked that,  
 644 with an appropriate tuning of  $b_p$ , the final scores are the same than with the zero/random initialisation but  
 645 the convergence is somewhat faster.

646 Finally, the online experiments use the same DAW length as with the physical model,  $L = 6$ . However,  
 647 since the accuracy of the model increases during the experiment, increasing  $L$  is a reasonable option. We have

---

<sup>2</sup>The convergence speed is measured here in number of DA cycles before convergence and not in wall-clock execution time.

performed the online experiments with  $L = 10$  (not illustrated here). The evolution of the tMSE (forecast score) is very similar to that with  $L = 6$ , but, as expected from fig. 6, the sRMSE (DA score) is close to 0.2, significantly lower than with  $L = 6$ .

### 5.5. Tendency correction in a realistic model

Both the offline and online learning experiments illustrate the efficiency of the TC method. Including a TC into a complex realistic model is not immediate, depending on the structure of the code, because the correction term must be added to the code computing the tendencies. In order to use the variational methods illustrated in the present work to train a TC, both offline and online, we must be able to compute the gradient of the resolvent of the total model (physical model with correction) with respect to the model state and to the model parameters. As we have shown in section 2.3, these gradients depend on the gradients of the correction term, which can be easily obtained with a ML library, but also on the TL operator of the physical model. The present work, with a simple model, has been largely facilitated by the fact that the code of the physical model has been written entirely with a ML library, which implies that we have benefited from automatic differentiation. For realistic models, which are inherently more complex, it is essential to have efficient differentiation methods.

Beyond these technical aspects, the number of trainable model parameters is a potential source of concern. In the present work, we have tried to keep it as small as possible, and ended up with 113 parameters (for both CNN-b and CNN-c). In preliminary experiments, a much larger NN (a fully-connected NN with a total of  $36 \times 64 + 64 + 64 \times 36 + 36 = 4708$  parameters) has been tested and we found similar performances, although with more training data. Even though 113 parameters (or even 4708) will most likely not be enough to correct a complex, realistic model, we have the hope that with smart ML models, we will be able to keep the number of trainable parameters under control (Bonavita and Laloyaux, 2020).

### 5.6. Generalisation to non-autonomous dynamics

We conclude this test series by briefly discussing the possibility to extend the present work to non-autonomous dynamics. With an offline learning method, a non-autonomous dynamics can only be learnt if (i) the time-dependency of the model is parametrised (which implies that time should be among the set of predictors) and (ii) the training dataset is large enough to infer the parametrisation. With an online learning method, parametrisation of the time-dependency is also an option, but such parametrisation is not mandatory if the time evolution of the dynamics is slow. For example, the online experiments of sections 5.2 and 5.3 would probably also work if the dynamics evolution is no faster than a few hundred DA cycles.

## 678 6. Conclusions

679 Combining DA and ML to emulate a dynamical model has been originally proposed by Brajard et al.  
680 (2020). The use of DA is essential here to assimilate sparse and noisy observations, which cannot be rigorously  
681 treated with ML methods alone. The same strategy can be used for model error correction instead of full  
682 model emulation. This has many advantages as it makes the inference problem easier (Jia et al., 2019;  
683 Watson, 2019). In practice, a correction term can be included either in the model resolvent (*i.e.* as an  
684 integrated term between two forecast times) or directly in the model tendencies (Bocquet et al., 2019; Farchi  
685 et al., 2021).

686 In the present article, we have compared the two methods. The first method, which we have called  
687 RC (resolvent correction), is easy to implement and has already been illustrated using low-order models  
688 by (Brajard et al., 2021; Farchi et al., 2021). The second method, which we have called TC (tendency  
689 correction), is more technical to implement, in particular because it requires the adjoint of the physical  
690 model to correct, but it has the advantage of being more flexible, in particular for short-term forecasts. Both  
691 methods have been tested using the two-scale Lorenz system. In this case, model error primarily comes  
692 from unresolved small-scales processes (the fast variables). We have used noisy observations of the slow  
693 variables to build three different surrogate models. All three surrogate models are hybrid, with a physical  
694 part, the non-corrected model, and a statistical part, the correction term, built using simple NNs. The first  
695 surrogate model uses the RC method, while the other two use the TC method. The surrogate models are  
696 trained using the offline DA-ML method of Brajard et al. (2020) and then evaluated in forecast and DA  
697 experiments. The results show that with the TC method, the surrogate models benefit from the interaction  
698 between the physical model and the NN, in such a way that it is possible to use much smaller NNs and much  
699 fewer training data to get similar results. The accuracy in forecast experiments is somewhat similar between  
700 the models using the RC and the TC method. By contrast, the models using the TC method significantly  
701 outperform the models using the RC method in DA experiments. This can be explained by the violation of  
702 the assumption of linear error growth in time, necessary with the RC method for the short-term forecasts  
703 within each DA experiment.

704 In the second part of this paper, we explored the possibility to train the surrogate models in an online  
705 fashion. With sparse and noisy observations, this means that we would have to learn both model state and  
706 model parameters at the same time. To address this problem, we introduce a new DA algorithm, which  
707 can be seen as a new formulation of WC 4D-Var. The algorithm has been implemented and tested using  
708 the same model error setup with the two-scale Lorenz system. The results show that online learning works  
709 as expected: the surrogate model is steadily improved, which can be seen both in the DA score and the  
710 forecast score; it quickly becomes more accurate than if trained offline with the analysis. At the end of the  
711 experiment, the model is almost as accurate as if trained offline with the truth. These results show that with

712 online learning, it is possible to extract all the information from sparse and noisy observations.

713 Even though the results of the online experiments are promising, the NN formulation of WC 4D-Var  
 714 derived in the present work could still benefit from methodological developments. In our experiments, we  
 715 have found that tuning the algorithmic parameter  $b_p$  is a difficult but critical task. Ideally, the tuning method  
 716 should be adaptive, with the objective of making the convergence faster and more accurate (Bocquet et al.,  
 717 2020b). Furthermore, the method implicitly assumes that background errors for model state and model  
 718 parameters are uncorrelated. Including cross-correlations between model state and model parameters is  
 719 possible; it would be interesting to check whether this could make a difference in the accuracy of the analysis.  
 720 More generally, the conclusions of the present work, and in particular the advantages of the TC method over  
 721 the RC method, must be confirmed using higher-dimensional models.

## 722 Acknowledgements

723 The authors are grateful to an anonymous reviewer for insightful comments which helped improve the  
 724 manuscript. CEREa is a member of Institut Pierre-Simon Laplace.

## 725 References

- 726 Abarbanel, H.D.I., Rozdeba, P.J., Shirman, S., 2018. Machine learning: Deepest learning as statistical data assimilation  
 727 problems. *Neural Computation* 30, 2025–2055. doi:{10.1162/neco\\_a\\_01094}.
- 728 Arcomano, T., Szunyogh, I., Pathak, J., Wikner, A., Hunt, B.R., Ott, E., 2020. A machine learning-based global atmospheric  
 729 forecast model. *Geophysical Research Letters* 47. doi:10.1029/2020GL087776.
- 730 Arcucci, R., Zhu, J., Hu, S., Guo, Y.K., 2021. Deep data assimilation: Integrating deep learning with data assimilation. *Applied*  
 731 *Sciences* 11, 1114. doi:10.3390/app11031114.
- 732 Bachlechner, T., Majumder, B.P., Mao, H.H., Cottrell, G.W., McAuley, J., 2020. Rezero is all you need: Fast convergence at  
 733 large depth. *arXiv:2003.04887 [cs, stat]* arXiv:2003.04887.
- 734 Bocquet, M., Brajard, J., Carrassi, A., Bertino, L., 2019. Data assimilation as a learning tool to infer ordinary differential  
 735 equation representations of dynamical models. *Nonlinear Processes in Geophysics* 26, 143–162. doi:10.5194/npg-26-143-2019.
- 736 Bocquet, M., Brajard, J., Carrassi, A., Bertino, L., 2020a. Bayesian inference of chaotic dynamics by merging data assimilation,  
 737 machine learning and expectation-maximization. *Foundations of Data Science* 2, 55–80. doi:10.3934/fods.2020004.
- 738 Bocquet, M., Farchi, A., Malartic, Q., 2020b. Online learning of both state and dynamics using ensemble kalman filters.  
 739 *Foundations of Data Science* 0, 0–0. doi:10.3934/fods.2020015.
- 740 Bolton, T., Zanna, L., 2019. Applications of deep learning to ocean data inference and subgrid parameterization. *Journal of*  
 741 *Advances in Modeling Earth Systems* 11, 376–399. doi:10.1029/2018MS001472.
- 742 Bonavita, M., Laloyaux, P., 2020. Machine learning for model error inference and correction. *Journal of Advances in Modeling*  
 743 *Earth Systems* 12. doi:10.1029/2020MS002232.
- 744 Brajard, J., Carrassi, A., Bocquet, M., Bertino, L., 2020. Combining data assimilation and machine learning to emulate a  
 745 dynamical model from sparse and noisy observations: A case study with the lorenz 96 model. *Journal of Computational*  
 746 *Science* 44, 101171. doi:10.1016/j.jocs.2020.101171.

- 747 Brajard, J., Carrassi, A., Bocquet, M., Bertino, L., 2021. Combining data assimilation and machine learning to infer unresolved  
748 scale parametrization. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*  
749 379, 20200086. doi:10.1098/rsta.2020.0086.
- 750 Brunton, S.L., Proctor, J.L., Kutz, J.N., 2016. Discovering governing equations from data by sparse identification of nonlinear  
751 dynamical systems. *Proceedings of the National Academy of Sciences* 113, 3932–3937. doi:10.1073/pnas.1517384113.
- 752 Byrd, R.H., Lu, P., Nocedal, J., Zhu, C., 1995. A limited memory algorithm for bound constrained optimization. *SIAM Journal*  
753 *on Scientific Computing* 16, 1190–1208. doi:10.1137/0916069.
- 754 Chollet, F., 2018. *Deep learning with Python*. Manning Publications Co, Shelter Island, New York.
- 755 Dueben, P.D., Bauer, P., 2018. Challenges and design choices for global weather and climate models based on machine learning.  
756 *Geoscientific Model Development* 11, 3999–4009. doi:10.5194/gmd-11-3999-2018.
- 757 Fablet, R., Oualla, S., Herzet, C., 2018. Bilinear residual neural network for the identification and forecasting of geophysical  
758 dynamics, in: *2018 26th European Signal Processing Conference (EUSIPCO)*, IEEE, Rome. pp. 1477–1481. doi:10.23919/  
759 EUSIPCO.2018.8553492.
- 760 Farchi, A., Laloyaux, P., Bonavita, M., Bocquet, M., 2021. Using machine learning to correct model error in data assimilation  
761 and forecast applications. *Quarterly Journal of the Royal Meteorological Society* 0, 0. Accepted for publication.
- 762 Fillion, A., Bocquet, M., Gratton, S., 2018. Quasi static ensemble variational data assimilation: a theoretical and numerical  
763 study with the iterative ensemble Kalman smoother. *Nonlin. Processes Geophys.* 25, 315–334. doi:10.5194/npg-25-315-2018.
- 764 Gagne, D.J., Christensen, H.M., Subramanian, A.C., Monahan, A.H., 2020. Machine learning for stochastic parameterization:  
765 Generative adversarial networks in the lorenz '96 model. *Journal of Advances in Modeling Earth Systems* 12. doi:10.1029/  
766 2019MS001896.
- 767 Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep learning. Adaptive computation and machine learning*, The MIT Press,  
768 Cambridge, Massachusetts.
- 769 Gottwald, G.A., Reich, S., 2021. Supervised learning from noisy observations: Combining machine-learning techniques with  
770 data assimilation. *Physica D: Nonlinear Phenomena* 423, 132911. doi:10.1016/j.physd.2021.132911.
- 771 Hamilton, F., Berry, T., Sauer, T., 2016. Ensemble kalman filtering without a model. *Physical Review X* 6, 011021.  
772 doi:10.1103/PhysRevX.6.011021.
- 773 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers,  
774 D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., Chiara,  
775 G., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L.,  
776 Healy, S., Hogan, R.J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., Rosnay, P.,  
777 Rozum, I., Vamborg, F., Villaume, S., Thépaut, J.N., 2020. The era5 global reanalysis. *Quarterly Journal of the Royal*  
778 *Meteorological Society* 146, 1999–2049. doi:10.1002/qj.3803.
- 779 Hsieh, W.W., Tang, B., 1998. Applying neural network models to prediction and data analysis in meteorology and oceanography.  
780 *Bulletin of the American Meteorological Society* 79, 1855–1870. doi:10.1175/1520-0477(1998)079<1855:ANNMTP>2.0.CO;2.
- 781 Jazwinski, A.H., 1970. *Stochastic processes and filtering theory*. Number 64 in *Mathematics in science and engineering*, Acad.  
782 Press, San Diego.
- 783 Jia, X., Willard, J., Karpatne, A., Read, J., Zwart, J., Steinbach, M., Kumar, V., 2019. Physics guided RNNs for modeling  
784 dynamical systems: A case study in simulating lake temperature profiles, in: *Proceedings of the 2019 SIAM International*  
785 *Conference on Data Mining, Society for Industrial and Applied Mathematics, Philadelphia, PA*. pp. 558, 566. doi:10.1137/1.  
786 9781611975673.
- 787 Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization, in: Bengio, Y., LeCun, Y. (Eds.), *3rd International*  
788 *Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*,  
789 San Diego, CA, USA.

- 790 Laloyaux, P., Bonavita, M., Chrust, M., Gürol, S., 2020a. Exploring the potential and limitations of weak-constraint 4d-var.  
791 Quarterly Journal of the Royal Meteorological Society 146, 4067–4082. doi:10.1002/qj.3891.
- 792 Laloyaux, P., Bonavita, M., Dahoui, M., Farnan, J., Healy, S., Hólm, E., Lang, S.T.K., 2020b. Towards an unbiased stratospheric  
793 analysis. Quarterly Journal of the Royal Meteorological Society 146, 2392–2409. doi:10.1002/qj.3798.
- 794 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521, 436–444. doi:10.1038/nature14539.
- 795 Lguensat, R., Tandeo, P., Ailliot, P., Pulido, M., Fablet, R., 2017. The analog data assimilation. Monthly Weather Review 145,  
796 4093–4107. doi:10.1175/MWR-D-16-0441.1.
- 797 Lorenz, E.N., 2005. Designing chaotic models. Journal of the Atmospheric Sciences 62, 1574–1587. doi:10.1175/JAS3430.1.
- 798 Lorenz, E.N., Emanuel, K.A., 1998. Optimal sites for supplementary weather observations: Simulation with a small model.  
799 Journal of the Atmospheric Sciences 55, 399–414. doi:10.1175/1520-0469(1998)055<0399:OSFSWO>2.0.CO;2.
- 800 Malartic, Q., Farchi, A., Bocquet, M., 2021. State, global and local parameter estimation using local ensemble kalman filters:  
801 applications to online machine learning of chaotic dynamics. SIAM/ASA Journal on Uncertainty Quantification 0, 0.  
802 Submitted.
- 803 Mitchell, L., Carrassi, A., 2015. Accounting for model error due to unresolved scales within ensemble kalman filtering. Quarterly  
804 Journal of the Royal Meteorological Society 141, 1417–1428. doi:10.1002/qj.2451.
- 805 Pathak, J., Hunt, B., Girvan, M., Lu, Z., Ott, E., 2018. Model-free prediction of large spatiotemporally chaotic systems from  
806 data: A reservoir computing approach. Physical Review Letters 120, 024102. doi:10.1103/PhysRevLett.120.024102.
- 807 Pires, C., Vautard, R., Talagrand, O., 1996. On extending the limits of variational assimilation in nonlinear chaotic systems.  
808 Tellus A 48, 96–121. doi:10.1034/j.1600-0870.1996.00006.x.
- 809 Pulido, M., Tandeo, P., Bocquet, M., Carrassi, A., Lucini, M., 2018. Stochastic parameterization identification using ensemble  
810 kalman filtering combined with maximum likelihood methods. Tellus A: Dynamic Meteorology and Oceanography 70, 1–17.  
811 doi:10.1080/16000870.2018.1442099.
- 812 Rasp, S., Pritchard, M.S., Gentine, P., 2018. Deep learning to represent subgrid processes in climate models. Proceedings of the  
813 National Academy of Sciences 115, 9684–9689. doi:10.1073/pnas.1810286115.
- 814 Ruiz, J.J., Pulido, M., Miyoshi, T., 2013. Estimating model parameters with ensemble-based data assimilation: A review.  
815 Journal of the Meteorological Society of Japan. Ser. II 91, 79–99. doi:10.2151/jmsj.2013-201.
- 816 Scher, S., Messori, G., 2019. Generalization properties of feed-forward neural networks trained on lorenz systems. Nonlinear  
817 Processes in Geophysics 26, 381–399. doi:10.5194/npg-26-381-2019.
- 818 Trémolet, Y., 2006. Accounting for an imperfect model in 4d-var. Quarterly Journal of the Royal Meteorological Society 132,  
819 2483–2504. doi:10.1256/qj.05.224.
- 820 Watson, P.A.G., 2019. Applying machine learning to improve simulations of a chaotic dynamical system using empirical error  
821 correction. Journal of Advances in Modeling Earth Systems 11, 1402–1417. doi:10.1029/2018MS001597.
- 822 Weyn, J.A., Durran, D.R., Caruana, R., 2019. Can machines learn to predict weather? using deep learning to predict gridded  
823 500-hpa geopotential height from historical weather data. Journal of Advances in Modeling Earth Systems 11, 2680–2693.  
824 doi:10.1029/2019MS001705.
- 825 Wikner, A., Pathak, J., Hunt, B., Girvan, M., Arcomano, T., Szunyogh, I., Pomerance, A., Ott, E., 2020. Combining machine  
826 learning with knowledge-based modeling for scalable forecasting and subgrid-scale closure of large, complex, spatiotemporal  
827 systems. Chaos: An Interdisciplinary Journal of Nonlinear Science 30, 053111. doi:10.1063/5.0005541.
- 828 Wilks, D.S., 2005. Effects of stochastic parametrizations in the lorenz '96 system. Quarterly Journal of the Royal Meteorological  
829 Society 131, 389–407. doi:10.1256/qj.04.03.

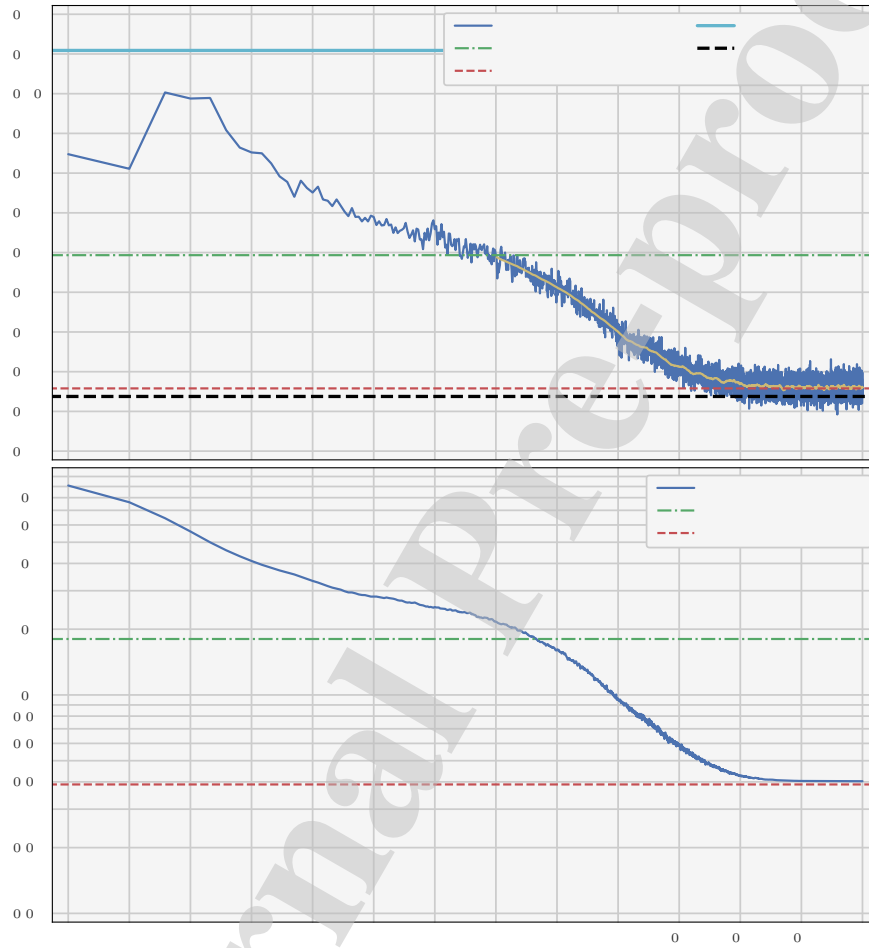


Figure 8: Time series of sRMSE (top panel) and tMSE (bottom panel) for the online experiment with TC-CNN-b (in blue). The tMSE is computed using a test dataset of size  $N_t = 1024$  and then normalised by the tMSE of the physical model. Both sRMSE and tMSE are averaged over 512 repetitions of the experiment. In addition, the yellow line shows the moving-average of the sRMSE over 128 DAWs. For comparison, the horizontal lines show the scores for the physical model (in cyan), the true model (in black), TC-CNN-b trained offline with the analysis (in green) and trained offline with the truth (in red).



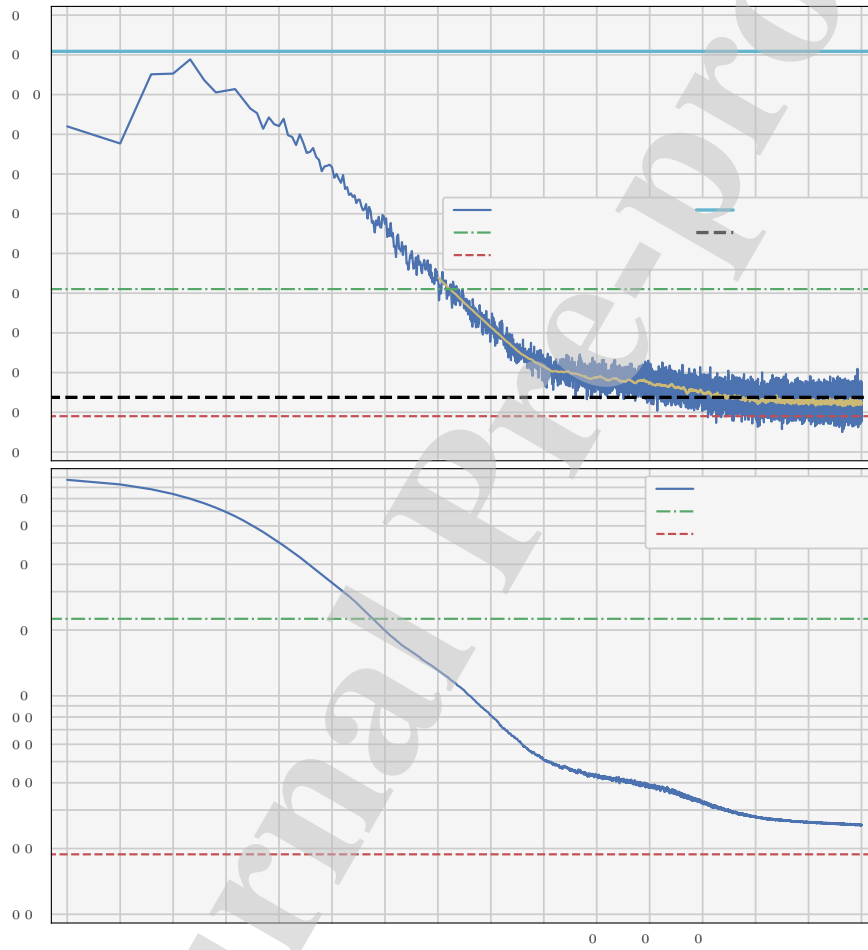


Figure 9: Same as fig. 8 for the online experiment with TC-CNN-c.

## **Highlights: A comparison of combined data assimilation and machine learning methods for offline and online model error correction**

Alban Farchi, Marc Bocquet,  
Patrick Laloyaux, Massimo Bonavita, Quentin Malartic

July 23, 2021

- Data assimilation and machine learning are combined to correct model error from sparse and noisy observations.
- The correction can be included in the model resolvent or in the model tendencies.
- With offline learning, the tendency correction is more accurate than the resolvent correction, although it is arguably less easy to implement.
- The tendency correction can be used for online learning, resulting in a new formulation of weak-constraint 4D-Var.
- The online correction is more accurate than the offline correction.

## **Highlights: A comparison of combined data assimilation and machine learning methods for offline and online model error correction**

Alban Farchi, Marc Bocquet,  
Patrick Laloyaux, Massimo Bonavita, Quentin Malartic

July 23, 2021

- Data assimilation and machine learning are combined to correct model error from sparse and noisy observations.
- The correction can be included in the model resolvent or in the model tendencies.
- With offline learning, the tendency correction is more accurate than the resolvent correction, although it is arguably less easy to implement.
- The tendency correction can be used for online learning, resulting in a new formulation of weak-constraint 4D-Var.
- The online correction is more accurate than the offline correction.

## Biography for all authors

Alban Farchi is a recently hired permanent researcher at CERE. He works in the field of data assimilation for the geosciences with application to atmospheric chemistry. In 2020, he has been a visitor scientist at ECMWF, working on the use of machine learning techniques to correct model error in data assimilation and forecast applications.

Marc Bocquet is Professor at École des Ponts (France) and deputy director of CERE laboratory. He works on the methods of data assimilation, environmental statistics and machine learning, with applications to dynamical systems, atmospheric chemistry and transport, and meteorology. He has published 107 peer-reviewed papers.

Trained as an applied mathematician, Patrick Laloyaux has expertise in optimisation methods for variational problems and in ensemble Kalman filters. He obtained his PhD from the University of Namur (Belgium) conducted in partnership with the European Centre for Research and Advanced Training in Scientific Computation (CERFACS, France). During his PhD, he developed new preconditioning and gradient-free techniques for solving variational problems. He was also a teaching assistant giving tutorials (220 hours/year) to undergraduate and graduate students in programming, scientific computation and numerical optimization.

Massimo Bonavita has a background in Physics. He has worked for 15 years in Italian Weather Service, first as a forecaster and then in data assimilation, where he developed ensemble and variational analysis systems for local area model initialisation. He joined ECMWF in March 2009 to work on the development of the variational data assimilation system and its ensemble extension. Since February 2016 Massimo leads the Data Assimilation Methodology Team at ECMWF, which works on the continuous improvement of the data assimilation algorithms used for the initialization of ECMWF forecasts.

Quentin Malartic is a PhD student at École des Ponts ParisTech and École Normale Supérieure. His research interest is data assimilation and machine learning in the context of chaotic dynamics. He holds a master's degree in both geosciences and civil engineering from Université Paris Saclay.

photo of author 4

[Click here to access/download;Author Photo;bonavita.png](#)



photo of author 3

[Click here to access/download;Author Photo;laloyaux.jpg](#) 



photo of author 5

[Click here to access/download;Author Photo;malartic.jpg](#)



photo of author 2

[Click here to access/download;Author Photo;bocquet.jpg](#) 





photo of author 1

[Click here to access/download;Author Photo;farchi.jpg](#) 



Author Statement

Alban Farchi: conceptualization; formal analysis; investigation; methodology; software; validation; visualization; writing original draft; writing review editing.

Marc Bocquet: conceptualization; methodology; supervision; writing review editing.

Patrick Laloyaux: conceptualization; methodology; supervision; writing review editing.

Massimo Bonavita: conceptualization; methodology; supervision; writing review editing.

Quentin Malartic: conceptualization; methodology; writing review editing.

**Declaration of interests**

X The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof