4D-Var data assimilation using an adjoint model of a neural network surrogate model

Seiya Nishizawa^{1,1}

¹RIKEN Center for Computational Science

November 30, 2022

Abstract

Four-dimensional variational (4D-Var) data assimilation is an effective method for obtaining physically consistent time-varying states. In this study, a method using a neural network surrogate model obtained by machine learning is proposed to solve one of the most serious challenges in 4D-Var: to construct an adjoint model. The feasibility of the proposed method was demonstrated by a 4D-Var experiment using a surrogate model for the Lorenz 96 model. In the method, several effective procedures have been proposed to obtain an accurate surrogate model and the assimilated initial conditions, including two-stage learning (i.e., single- and multi-step learning) of neural networks, limiting the target states of the surrogate model to a small subspace of the state phase space, and updating the surrogate model during 4D-Var iterations.

¹ 4D-Var data assimilation using an adjoint

² model of a neural network surrogate model

Seiya Nishizawa

3

5

⁴ RIKEN Center for Computational Science, Kobe, Japan

March 4, 2022

Corresponding author: Seiya Nishizawa, RIKEN Center for Computational Science, 7-1-26 Minatojima-minami-machi, Chuo-ku, Kobe, Hyogo 650-0047, Japan. E-mail: s-nishizawa@riken.jp

Abstract

Four-dimensional variational (4D-Var) data assimilation is an effective method 7 for obtaining physically consistent time-varying states. In this study, a 8 method using a neural network surrogate model obtained by machine learn-9 ing is proposed to solve one of the most serious challenges in 4D-Var: to 10 construct an adjoint model. The feasibility of the proposed method was 11 demonstrated by a 4D-Var experiment using a surrogate model for the 12 Lorenz 96 model. In the method, several effective procedures have been 13 proposed to obtain an accurate surrogate model and the assimilated initial 14 conditions, including two-stage learning (i.e., single- and multi-step learn-15 ing) of neural networks, limiting the target states of the surrogate model 16 to a small subspace of the state phase space, and updating the surrogate 17 model during 4D-Var iterations. 18

6

¹⁹ Keywords Data assimilation; surrogate model; machine learning.

20 1. Introduction

Optimal initial conditions are crucial for accurate deterministic numeri-21 cal simulations. Data assimilation is widely used to obtain initial conditions, 22 e.g., it is an essential component of numerical weather forecasting systems. 23 Four-dimensional variational (4D-Var) data assimilation has the advantage 24 of obtaining a time evolution consistent with model physics. This is im-25 portant, especially when obtaining four-dimensional analysis data of target 26 phenomena to determine its mechanism. The 4D-Var method requires an 27 adjoint model of the simulation model for the backward calculation of a 28 cost function's gradient with respect to the initial conditions. Building the 29 adjoint model and updating it with the simulation model is costly, which is 30 an important challenge of the 4D-Var method. Despite this disadvantage, 31 the method has been employed in several operational numerical weather 32 forecasting systems. Research-purpose simulation models tend to have a 33 shorter lifetime than operational models, and usually have multiple simula-34 tion paths (several different schemes) for individual physical processes, from 35 which users can choose according to their objectives. Therefore, developing 36 and managing adjoint models of research models may require more effort 37 than operational models. 38

Machine learning techniques have developed rapidly and are used in an 39 increasing range of domains. Data assimilation and machine learning have 40 some similarities (Geer 2021); both minimize error (the cost or loss function) 41 by optimizing target quantities, such as the state vector (data assimilation) 42 and network parameters (machine learning). In neural network training, 43 network parameters are updated according to their loss function's gradient. 44 To obtain the gradient, a backward propagation algorithm is generally used. 45 Recently, excellent machine learning frameworks, such as Pytorch (Paszke 46 et al. 2019) and TensorFlow (Abadi et al. 2015), have been developed, 47 which can easily compute gradients. Note that the procedure for learning 48 network parameters is the same as that for updating the initial conditions 49 with the adjoint model in 4D-Var. Therefore, once the forward simulation 50 model is constructed, it is not necessary to manually build its adjoint model; 51 the backward calculation of the gradient with respect to the initial condi-52 tions can be performed using the functionality of the framework. However, 53 physics-based simulation models built using the framework generally require 54 more computational resources, such as CPU time and memory usage, than 55 conventional models written in C or Fortran, which may not be practical. 56 Replicating physics-based simulations with a neural network surrogate 57

⁵⁸ model is a possible solution. Surrogate models are not based on physical
⁵⁹ laws (e.g., governing equations), but on statistical relationships between the

initial conditions and simulation results (Grzeszczuk et al. 1998; Dueben 60 and Bauer 2018). Surrogate models are built by machine learning on inputs 61 and outputs of physics-based simulations and can be designed to be compu-62 tationally less expensive. The functionality of machine learning framework 63 can be used to calculate the cost function's gradient of the neural network 64 surrogate model with respect to the initial conditions. Even without the 65 functionality, building a adjoint model of the neural network model man-66 ually is much easier than building a physics-based adjoint model because 67 neural networks generally consist of a limited number of simple operations, 68 such as weighted sums, and only a few nonlinear activation functions. Us-69 ing a surrogate model's adjoint model may make 4D-Var data assimilation 70 easier. 71

There are two major concerns with using surrogate models in 4D-Var 72 data assimilation. (1) Can a surrogate model be obtained that provides 73 sufficiently accurate simulation results? For systems with many degrees of 74 freedom (e.g., atmospheric system), surrogate models must also have suffi-75 cient degrees of freedom (Dueben and Bauer 2018). The greater the degrees 76 of freedom, the more difficult it is to build a surrogate model. Limiting the 77 target space of the surrogate model in the phase space to a small subspace 78 around the target state, building a surrogate model is expected to be more 79 easier than when targeting the entire space. (2) Can the gradients be ac-80

curate enough to improve the initial conditions? Even if a surrogate model 81 providing accurate forward computations can be obtained, its Jacobian may 82 not be accurate (Chevallier and Mahfouf 2001; Aires et al. 2004). For exam-83 ple, if the resulting network overfits the training data, the gradients may be 84 unrealistic, even if the results of the forward simulation appear reasonable. 85 Several studies have proposed a similar concept using machine learning 86 for data assimilation. Brajard et al. (2020) combined data assimilation and 87 machine learning without a physics-based model; the amount of training 88 data was capped because they were limited to the observation data. This 89 limitation can make program overfitting more serious. Hatfield et al. (2021) 90 used an adjoint model obtained by machine learning for 4D-Var data assimi-91 lation. Training data were generated using simulations with a physics-based 92 model and there was no limit to the amount of training data in principle. 93 They demonstrated 4D-Var by replacing a parameterization scheme in the 94 general circulation model with a neural network. The scheme was only 95 a part of the model, and most parts of the adjoint model were derived 96 manually, as in the conventional method. Nonnenmacher and Greenberg 97 (2021) replaced a whole physics-based model with a neural network model, 98 which output tendencies of the prognostic variables and was trained with aa the output tendencies of the physics-based model. For time integration, a 100 conventional method, such as the Runge-Kutta method, was used. This 101

scheme has advantages; tendencies, which can help understand the mecha-102 nisms of phenomena, and states after the arbitrary integration period can 103 be obtained. However, even if the tendencies only have a tiny error, the 104 error may accumulate (grow) during time integration, since the network is 105 trained by instantaneous data. Alternatively, the neural network model can 106 be designed to output a state after a certain time integration period. In this 107 case, the model could be trained so that the error grown in a finite time 108 integration period is reduced. 100

This study investigates the feasibility of using a neural network surrogate 110 model to improve the initial conditions in a 4D-Var data assimilation, where 111 the surrogate model is trained by simulation results from a physics-based 112 model and outputs a state after certain time period. A simple dynamical 113 system was used to study the feasibility. Several physics-based simulations 114 were performed, then a neural network surrogate model was built using 115 the simulations' output data. Efficient ways to train the network were also 116 investigated. Using this surrogate model, a 4D-Var data assimilation exper-117 iment was conducted, and a promising method was proposed to efficiently 118 improve the initial conditions during the assimilation iteration. 119

¹²⁰ 2. Model and Methodologies

$_{121}$ 2.1 Lorenz 96 model

The Lorenz 96 model is a dynamical system model (Lorenz 1996):

$$\frac{dx_i}{dt} = (x_{i+1} - x_{i-2})x_{i-1} - x_i + F, i = 1, 2, \cdots I,$$
(1)

where I is the number of grid points. The first, second, and last terms on the right-hand side correspond to the advection, diffusion, and forcing terms, respectively. The system exhibits chaotic behavior for a range of Fvalues (Lorenz and Emanuel 1998).

Physics-based simulations were performed using the fourth-order Runge-126 Kutta scheme with a time step of $\Delta t = 0.01$. I and F were set to 40 and 8.0, 127 respectively, at which the system was chaotic. Periodic boundary conditions 128 were employed, and the initial conditions were $x_i = F + \epsilon_i$, where ϵ is a small 129 random normally distributed perturbation with a standard deviation of 0.01. 130 After a spin-up of 5,100 integration time steps, time integration of 100 131 steps from t = 0 to 1 was performed (hereafter referred to as the reference). 132 Then, the ensembles were generated by adding random normally distributed 133 perturbations, with a standard deviation of 0.1, to the reference state after 134 the first spin-up 5,000 steps. After 100 integration steps for the second 135 spin-up, 100 time-integration steps were performed for each ensemble. 136

¹³⁷ The ensemble average of the mean squared error (MSE) from the refer-

ence grows exponentially, with a growth rate (Lyapunov exponent) of approximately 2.14. The states were output every five steps ($\Delta_{output} = 0.05$), and there were 21 outputs of 40-dimensional vectors \boldsymbol{x} , including the initial state, for each ensemble. These were used to train the neural network surrogate model.

The states of these ensembles lie within a limited region (a subspace) 143 around the reference in the state phase space (hereafter referred to as the 144 localized ensemble set). To examine the effect of extent of the training 145 data's state in the phase space on the surrogate model trained from the 146 data, another ensemble set was generated with a second spin-up of 1,000 147 steps (hereafter referred to as the spread ensemble set); the second spin-up 148 was 100 steps for the localized ensemble set. The longer spin-up resulted in 149 a wider spread; the states of the spread ensemble set are widely spread in 150 the phase space with a large variance that is comparable in magnitude to 151 the variance of a very long time series. The MSE of the spread ensemble 152 set is approximately 24–29 throughout the integration period, whereas the 153 MSE of the localized ensemble set is approximately 0.27 and 2.30 at the 154 beginning and end of the integration period, respectively. 155

156 2.2 Surrogate model

Using the state vectors \boldsymbol{x} of the physics-based simulation as both input and target data for training, a neural network surrogate model replicating the physics-based simulation was built. The calculations in the neural network were conducted using Pytorch.

¹⁶¹ *a.* Network architecture

In the physics-based model, the state at the next time step depends 162 only on that of the previous step. To emulate this behavior, the network 163 was designed as a recurrent neural network; an identical network module 164 is connected recurrently, and each module corresponds to a time interval 165 of 0.05 (Figure 1a). Each module consists of a stacked hourglass network 166 (Newell et al. 2016) of two stages: down-sampling and up-sampling (Fig-167 ure 1b). Through these stages, multiple horizontal scales are considered. 168 In down-sampling, the grid size is halved at each step by the max-pooling 169 layer; there are four steps and the grid size at each step is 40, 20, 10, or 170 5. In up-sampling, the grid size is doubled at each step, from 5 to 40, by 171 the max-unpooling layer. The output of each down-sampling step is added 172 to the input of the corresponding up-sampling step via skip connections. 173 Each step consists of convolution layers, batch normalization layers, and 174 rectified linear unit activation layers. In the convolution layers, the values 175

of neighboring grid points interact, and the convolution and following activation layers are expected to replicate advection. The convolution kernel size for the hourglass network is three, and periodic boundary padding is applied before the convolution. Before the hourglass network, the number of channels is increased from one to four by convolution with a kernel size of one, and after the network, the number of channels is reduced from four to one.

Fig. 1

183 b. Training and evaluation

Training the neural network was divided into two stages: single-step and multi-step learning. In single-step learning, a non-recurrent single network module was trained. The training data input were $\boldsymbol{x}(t)$ and the target data were $\boldsymbol{x}(t+0.05)$, where $t = 0, 0.05, \dots, 0.95$. The loss function l_1 is defined as

$$l_1 = \frac{1}{I} \left| f(\boldsymbol{x}(t)) - \boldsymbol{x}(t+0.05) \right|^2, \qquad (2)$$

where f is the operator corresponding to the single network module. With the 21-step output dataset obtained in each ensemble run, 20 training inputoutput pairs were available; from M ensemble runs, a training dataset of 20M pairs could be used.

In the multi-step learning process, the input data were $\boldsymbol{x}(t)$ and the target data were $(\boldsymbol{x}(t+0.05), \boldsymbol{x}(t+0.1), \cdots, \boldsymbol{x}(t+0.5))$, where $t = 0, 0.05, \cdots, 0.5$.

Each ensemble run contained 11 training data pairs, and a total of 11M training data pairs could be used from M ensemble runs. The loss function l_{10} is defined as

$$l_{10} = \frac{1}{10I} \sum_{n=1}^{10} |f^n(\boldsymbol{x}(t)) - \boldsymbol{x}(t+0.05n)|^2.$$
(3)

The network parameters obtained by single-step learning were used as the initial parameters for multi-step learning. In the multi-step learning process, the dropout layers were disabled, and the mean and standard deviations in the batch normalization layers were fixed to the values obtained in singlestep learning.

For both the single- and multi-step learning processes, the error of the 193 trained network was evaluated by the ensemble average of their loss func-194 tions, calculated using the evaluation data of another ensemble dataset of 195 100 runs, which was generated in the same manner as the training datasets, 196 with the same number of spin-up steps. The batch size was swept and de-197 termined such that the error of the network was minimized. The size of the 198 training dataset, which is proportional to the ensemble size, was also var-199 ied for sensitivity testing. The Adam optimization algorithm (Kingma and 200 Ba 2014) was used to update the network parameters. The initial learning 201 rates were set to 10^{-5} and 10^{-6} times batch size for single- and multi-step 202 learning, respectively, with a decay every 1,000 epochs by a factor of 0.99. 203

204 2.3 4D-Var data assimilation

In the 4D-Var data assimilation experiment, the neural network surrogate model and its adjoint model were used for the forward and backward computations, respectively. The time window for assimilation was 0.5, and the observed data were assimilated at intervals of 0.05, for 10 time steps. Since the non-dimensional time unit in this system is roughly equivalent to 5 days (Lorenz 1996), the time window corresponds to 2–3 days. The cost function J_s was defined as follows:

$$J_{s}(\boldsymbol{x}(0)) = \frac{1}{2} (\boldsymbol{x} - \boldsymbol{x}_{b})^{T} B^{-1} (\boldsymbol{x} - \boldsymbol{x}_{b}) + \frac{1}{2} \sum_{n=1}^{10} \left[H(f^{n}(\boldsymbol{x}(0))) - \boldsymbol{y}(0.05n) \right]^{T} R^{-1} \left[H(f^{n}(\boldsymbol{x}(0))) - \boldsymbol{y}(0.05n) \right],$$
(4)

where \boldsymbol{x}_b is the first guess for $\boldsymbol{x}(0), \boldsymbol{y}(t)$ is the observed state at time t, B is 205 the background covariance matrix, R is the observation covariance matrix, 206 and H is the observation operator. The first guess was the initial state of 207 one of the ensemble runs in Section 2.1. B was calculated from the 1,000 208 ensembles, and the mean of the diagonal components was 0.27. Observa-209 tions were generated to have a normally distributed error with a standard 210 deviation of 0.1; the matrix R is diagonal and its diagonal components were 211 0.01. The observed data were located at 10 random grids, chosen from the 212 40 grids: i = 3, 12, 17, 20, 25, 26, 27, 31, 32, and 33. The operator H reveals 213

²¹⁴ the data at these locations.

The gradient of J_s for $\boldsymbol{x}(0)$ was obtained using the surrogate model's 215 adjoint model; the gradient was calculated automatically by Pytorch. The 216 gradient was then used to update $\boldsymbol{x}(0)$. A one-dimensional golden-section 217 search procedure (Kiefer 1953) was used for the update. $\boldsymbol{x}(0)$ was updated 218 iteratively to reduce the cost function, and the maximum iteration count 219 was 1,000. At each of the K iterations (hereafter referred to as the update 220 interval), the physics-based simulation was performed using the latest $\boldsymbol{x}(0)$, 221 and the network parameters of the surrogate model were updated using the 222 simulation results in the same way as multi-step learning. 223

To evaluate the assimilated initial state, an extended forecast simulation was conducted from the initial state using the physics-based model and the root mean squared error (RMSE) of the forecast state from the reference state at t = 1 was calculated. The learning rate in the surrogate model update swept between 10^{-5} , 3×10^{-5} , and 10^{-4} , and was determined such that the RMSE improvement was maximized.

The 4D-Var experiment was conducted with ten different first guesses for each parameter. In the analysis, the RMSE was averaged over the ten results.

233 **3.** Results

²³⁴ 3.1 Obtaining the surrogate model

First, the neural network was trained using single-step learning. The 235 error of the network depends on the size of the training data; the larger 236 the size of the dataset, the smaller the error is (Figure 2a). The error is 237 $O(10^{-4})-O(10^{-2})$, which is much smaller than the O(10) background vari-238 ance of \boldsymbol{x} . The batch sizes with the smallest errors were 125, 250, 2000, 239 4000, and 4000 for ensemble sizes of 50, 100, 200, 400, and 800, respec-240 tively. Note that the ensemble size is proportional to the training dataset 241 size. Then, using the surrogate model, a time integration experiment for 242 t = 0-1 was conducted, i.e., the network obtained above was repeated 20 243 times. This time integration was calculated from 1,000 different generated 244 initial states, as in Section 2.1. The accuracy of the surrogate model was 245 evaluated by MSE from the physics-based model solution from the same 246 initial states. Note that the MSE is only due to model error, since the ini-247 tial conditions have no error. Figure 2c shows the temporal evolution of the 248 ensemble average of the surrogate model's MSEs obtained from an ensemble 249 set of M = 800. The MSE grows over time, with the growth rate initially 250 gradually decreasing and then remaining nearly constant. Even later, the 251 growth rate is still larger than the growth rate of the physical growth mode 252

 $_{253}$ in Section 2.1.

Next, multi-step learning was performed. Figure 2b shows the multi-254 step learning network error. Because the discrepancy between the surrogate 255 model and the physics-based simulations tends to increase with time, the 256 magnitude of the network error is larger than that in single-step learning. 257 The error depends on the ensemble size, as in the case of single-step learning, 258 but the dependency on batch size is smaller than in single-step learning. Us-250 ing the multi-step learning surrogate model, time integration was performed 260 as for single-step learning. The early growth rate of the error was improved 261 compared to that of the single-step learning model (Figure 2c). As a result, 262 the error at the end of the timespan is approximately 60% of that obtained 263 by the single-step learning model. Conversely, the error after the first step 264 (t = 0.05) is larger than that of the single-step learning model. This can be 265 explained as follows. In single-step learning, the network learns such that 266 the error after single-step integration is small, whereas in the multi-step 267 learning process, the network learns such that the average error of 10 steps 268 is small. This means that unstable modes with large Jacobian eigenvalues 269 become smaller in single-step learning, and unstable modes with large sin-270 gular values become smaller in the multi-step learning process. The fastest 271 growing mode in terms of instantaneous temporal difference is represented 272 by the eigenvector and the mode over a finite-time interval is represented 273

by the singular vector. This is consistent with the expectation of Brenowitz and Bretherton (2018), that a multiple-time-step loss function penalizes a growing unstable mode.

To evaluate the efficacy of the two-stage learning process (single- and 277 multi-step learning), one-stage learning (multi-step learning only) was also 278 conducted. Randomly generated initial network parameters were used for 279 multi-step learning. In this case, the network did not learn well; the loss 280 function did not decrease significantly during the epoch iteration and sat-281 urated at the level of O(1). As a result, the error of the surrogate model 282 obtained by one-stage learning was much larger than that obtained by two-283 stage learning. This is due to the total number of layers in the network 284 being too deep. This may be solved by a more appropriate network design. 285 Regardless of the case, the network was successfully trained by two-stage 286 learning. This suggests that two-stage learning is an efficient way to build 287 surrogate models. 288

To investigate the effect on the accuracy of the surrogate model by limiting the training data's state to a subspace of the phase space, the same training was performed using the spread ensemble set generated with the longer second spin-up of 1,000 steps, instead of the localized ensemble set of the 100-step spin-up. The errors of the networks obtained using the spread ensemble set were 0.40, 0.20, 0.088, 0.044, and 0.018 for ensemble

sizes of M = 50, 100, 200, 400, and 800, respectively. These errors were 295 approximately 6.5 to 10 times larger than the corresponding errors with 296 the localized ensemble set (0.045, 0.024, 0.010, 0.0041, and 0.0028, respec-297 tively). This shows that limiting the state of the target space is effective in 298 increasing the model's accuracy. A surrogate model targeting wider space 299 may require larger training data size and/or more complex network archi-300 tecture. This suggests that the difficulty of building a surrogate model can 301 be reduced by limiting the target states to a small phase subspace. 302

303 3.2 4D-Var experiment

A 4D-Var data assimilation experiment was conducted using the neu-304 ral network surrogate model. Figures 3a and 3b show the evolution of the 305 cost function J_s with the number of iterations. As the number of iterations 306 increases, the cost generally decreases. The same cost function using the 307 physics-based model was also calculated (J_p) . J_s and J_p were calculated 308 from the time series integrated from the same initial conditions. J_p is due 309 to the errors in the initial state and observation, while J_s is due to not 310 only these errors, but also the model error. J_p is generally larger than J_s , 311 since the initial conditions have been updated so that J_s decreases. Nev-312 ertheless, we see that J_p decreases with an increasing number of iterations. 313 This indicates that 4D-Var data assimilation using an adjoint model of a 314

neural network surrogate model is effective. The cost is smaller for a larger 315 ensemble size, which corresponds to a smaller error of the surrogate model, 316 suggesting that the accuracy of the surrogate model affects the accuracy of 317 the assimilation. Additionally, updating the surrogate model during 4D-318 Var iterations improved the cost. We observed large improvements in the 319 cost when the network was updated, e.g., at 100 iterations for the case of 320 M = 200 and K = 100. We also observed that the smaller the update 321 interval, the smaller the cost. 322

To evaluate the accuracy of the assimilated initial conditions, the ex-323 tended forecasts' RMSE at t = 1 from the assimilated initial states after 324 1,000 iterations was examined. Figure 4a shows the RMSE averaged over 325 10 samples; it was approximately 0.42 for large ensemble-size cases, and 326 1.59 for cases from the first guess. This indicates that assimilation using 327 the surrogate model improved the accuracy of the initial state. The depen-328 dency of the RMSE on the ensemble size and update interval shows similar 329 characteristics to those of the cost; the RMSE is likely to be small for large 330 ensemble sizes and smaller update intervals. As a reference, a 4D-Var ex-331 periment with a manually constructed adjoint model of the physics-based 332 model was also conducted. The forecast's RMSE from the assimilated ini-333 tial state was 0.39. This shows that assimilation using the surrogate model 334 can achieve similar accuracy to that using conventionally manually obtained 335

Fig. 3

³³⁶ adjoint model of the physics-based model.

Fig. 4

337 4. Conclusions

A 4D-Var assimilation method was proposed using an adjoint model of a neural network surrogate model. Additionally, several procedures were proposed to efficiently obtain an accurate surrogate model and assimilated initial conditions. As a feasibility study, a surrogate model was constructed and a 4D-Var assimilation experiment was conducted using the Lorenz 96 model.

Two-stage learning was efficient for obtaining an accurate surrogate 344 model. In the first stage, the network was trained from a training dataset, 345 with target data one step forward from the input data, obtained using 346 the physics-based model (single-step learning). In the next stage, the net-347 work was trained using time-series data of multiple steps as the target data 348 (multi-step learning). In this stage, the network parameter obtained in the 340 first stage was used as the initial value. It is found that the neural network 350 model trained by time-sequence data with a longer time period has better 351 accuracy than that with shorter period. The neural network model which 352 output tendencies, as proposed by Nonnenmacher and Greenberg (2021), is 353 thought to have the same problem with models trained by shorter period 354 data. We found that limiting the target states of the surrogate model to 355

a state phase subspace around the target case was efficient for building an
accurate surrogate model.

The 4D-Var assimilation experiment showed that the initial conditions 358 were improved by assimilation by using the adjoint model of the surrogate 359 model. It was also found that updating the surrogate model during the 360 4D-Var iterations was effective in improving the accuracy of the initial con-361 ditions. Even if the accuracy of the initial surrogate model is not very high 362 (e.g., M = 50), accurate initial conditions could be obtained with frequent 363 updates (small K) of the surrogate model during 4D-Var iterations. In gen-364 eral, more accurate data contribute to better training in machine learning. 365 Therefore, learning during 4D-Var iterations is likely to be more efficient 366 than the earlier two-stage learning because the training data used for up-367 dating the network are more accurate due to better initial conditions. On 368 the other hand, frequent updates require large computational resources; up-369 dating the network requires physics-based simulation and training with the 370 simulation data. The optimal values of the ensemble size and update in-371 terval must be determined by balancing the computational costs for each 372 stage of training and assimilation. 373

Assimilation has an affinity for limiting states in the phase space to build are a surrogate model. The states in a finite assimilation window generally occupy only a small subspace and, therefore, a surrogate model that covers all possible states is not needed; we can focus on the subspace around the target state to be assimilated. However, the surrogate model needs to be rebuilt for different cases. The learning speed of the network in other cases can be significantly improved by using the network obtained in one case as the initial parameters, i.e., transfer learning.

As simulation models become more sophisticated, they become more complex, which requires more effort from researchers. The effective use of data science techniques will become increasingly important for various aspects of simulation research.

Data Availability Statements

The source code of the programs used in this study are available from the Zenodo repository at https://doi.org/10.5281/zenodo.6319009. The datasets generated by the experiment are available from the Zenodo repository at https://doi.org/10.5281/zenodo.6318869.

391

386

Acknowledgments

This work was supported by JSPS KAKENHI (Grant Number JP19H01974) and JST AIP (Grant Number JPMJCR19U2). Some results were obtained

using the Wisteria/BDEC-01 system at the Information Technology Cen-394 ter, The University of Tokyo. The calculations associated with the neural 395 network were conducted using Pytorch (https://pytorch.org/). The dia-396 grams were drawn using tools developed by the GFD-Dennou Club (https: 397 //www.gfd-dennou.org/index.html.en). 398

399

References

400	Abadi, M., A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S.
401	Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfel-
402	low, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser,
403	M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray,
404	C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Tal-
405	war, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals,
406	P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, 2015:
407	TensorFlow: Large-scale machine learning on heterogeneous systems.
408	Software available from tensorflow.org.
409	Aires, F., C. Prigent, and W. B. Rossow, 2004: Neural network uncertainty

- assessment using Bayesian statistics with application to remote sens-410
- ing: 3. Network Jacobians. J. Geophys. Res., 109, D10305. 411

Brajard, J., A. Carrassi, M. Bocquet, and L. Bertino, 2020: Combining data 412

414	from sparse and noisy observations: A case study with the Lorenz
415	96 model. J. Comput. Sci., 44, 101171.
416	Brenowitz, N. D., and C. S. Bretherton, 2018: Prognostic validation of
417	a neural network unified physics parameterization. Geophys. Res.
418	<i>Lett.</i> , 45 , 6289–6298.
419	Chevallier, F., and JF. Mahfouf, 2001: Evaluation of the Jacobians of
420	infrared radiation models for variational data assimilation. $J.$ Appl.
421	Meteor., 40 , 1445–1461.
422	Dueben, P. D., and P. Bauer, 2018: Challenges and design choices for global
423	weather and climate models based on machine learning. Geosci.
424	Model Dev., 11 , 3999–4009.
425	Geer, A. J., 2021: Learning earth system models from observations: ma-
426	chine learning or data assimilation? Phil. Trans. R. Soc. A., 379,
427	20200089.
428	Grzeszczuk, R., D. Terzopoulos, and G. Hinton, 1998: Neuroanimator: Fast
429	neural network emulation and control of physics-based models. In
430	Proceedings of the 25th Annual Conference on Computer Graphics
431	and Interactive Techniques, SIGGRAPH '98, Association for Com-
432	puting Machinery, New York, NY, USA, 9–20.

assimilation and machine learning to emulate a dynamical model

413

432

23

433	Hatfield, S., M. Chantry, P. Dueben, P. Lopez, A. Geer, and T. Palmer,
434	2021: Building tangent-linear and adjoint models for data assim-
435	ilation with neural networks. J. Adv. Model. Earth Syst., 13,
436	e2021MS002521.

- ⁴³⁷ Kiefer, J., 1953: Sequential minimax search for a maximum. *Proc. American*⁴³⁸ *Math. Soc.*, 4, 502–506.
- Kingma, D. P., and J. Ba, 2014: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Lorenz, E. N., 1996: Predictability a problem partly solved. In Seminar
 on predictability, ECMWF, Reading, UK, 1–18.
- Lorenz, E. N., and K. A. Emanuel, 1998: Optimal sites for supplementary weather observations: Simulation with a small model. J. Atmos. Sci., 55, 399–414.
- Newell, A., K. Yang, and J. Deng, 2016: Stacked hourglass networks for
 human pose estimation. , *Computer vision ECCV 2016*, Springer
 International Publishing, Cham, 483–499.
- ⁴⁴⁹ Nonnenmacher, M., and D. S. Greenberg, 2021: Deep emulators for differen⁴⁵⁰ tiation, forecasting, and parametrization in Earth science simulators.
 ⁴⁵¹ J. Adv. Model. Earth Syst., 13, e2021MS002554.

452	Paszke, A., S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan,
453	T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf,
454	E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy,
455	B. Steiner, L. Fang, J. Bai, and S. Chintala, 2019: Pytorch:
456	An imperative style, high-performance deep learning library. Ad -
457	vances in Neural Information Processing Systems 32, H. Wallach,
458	H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox and R. Gar-
459	nett, Eds., Curran Associates, Inc., 8024–8035.

List of Figures

461	1	Surrogate model network architecture: (a) recurrent network	
462		modules and (b) details of the hourglass network. "Conv,"	
463		"BN," "ReLU," and "Dropout" indicate convolution, batch	
464		normalization, rectified linear unit, and dropout layers, re-	
465		spectively. The number following the convolution indicates	
466		the kernel size. The number above each box in (b) is the	
467		channel size times the grid (neuron) size.	27
468	2	The error of the network obtained by (a) the single- and (b)	
469		multi-step learning, and (c) the temporal evolution of the	
470		error of the surrogate model obtained with 800 ensembles by	
471		(red) single-step and (blue) multi-step learning. The symbols	
472		and colors in (a) and (b) represent the batch size normalized	
473		by the ensemble size. The broken line in (c) represents the	
474		error growth rate of the physical growth mode	28
475	3	The temporal evolution of the cost as a function of the iter-	
476		ation count in the 4D-Var assimilation with (a) the update	
477		interval $K = 10$ and (b) the ensemble size $M = 200$. The	
478		solid and broken lines represent J_p and J_s , respectively	29
479	4	RMSE dependency of the extended forecast at $t = 1$ on (a)	
480		the ensemble size and (b) the update interval. The symbols	
481		and colors represent the ensemble size and update interval,	
482		respectively.	30



Fig. 1. Surrogate model network architecture: (a) recurrent network modules and (b) details of the hourglass network. "Conv," "BN," "ReLU," and "Dropout" indicate convolution, batch normalization, rectified linear unit, and dropout layers, respectively. The number following the convolution indicates the kernel size. The number above each box in (b) is the channel size times the grid (neuron) size.



Fig. 2. The error of the network obtained by (a) the single- and (b) multistep learning, and (c) the temporal evolution of the error of the surrogate model obtained with 800 ensembles by (red) single-step and (blue) multi-step learning. The symbols and colors in (a) and (b) represent the batch size normalized by the ensemble size. The broken line in (c) represents the error growth rate of the physical growth mode.



Fig. 3. The temporal evolution of the cost as a function of the iteration count in the 4D-Var assimilation with (a) the update interval K = 10 and (b) the ensemble size M = 200. The solid and broken lines represent J_p and J_s , respectively.



Fig. 4. RMSE dependency of the extended forecast at t = 1 on (a) the ensemble size and (b) the update interval. The symbols and colors represent the ensemble size and update interval, respectively.