# Journal of the Meteorological Society of Japan

## EARLY ONLINE RELEASE

# Nonlinear Data Assimilation by Deep Learning Embedded in an Ensemble Kalman Filter

**Tadashi TSUYUKI**[1]

*Meteorological Research Institute*
*Japan Meteorological Agency, Tsukuba, Japan*

**and**

**Ryosuke TAMURA**

*Research Institute for Sustainable Humanosphere*
*Kyoto University, Kyoto, Japan*

------------------------------------

1) Corresponding author: Tadashi Tsuyuki, Observation and Data Assimilation Research Department, Meteorological Research Institute, 1-1 Nagamine, Tsukuba, 305-0052 JAPAN.
Email: ttuyuki@mri-jma.go.jp
Tel: +81-29-853-8642
Fax: +81-29-853-8649

**Abstract**

Recent progress in the particle filter has made it possible to use it for nonlinear or non-Gaussian data assimilation in high-dimensional systems, but a relatively large ensemble is still needed to outperform the ensemble Kalman filter (EnKF) in terms of accuracy. An alternative ensemble data assimilation method based on deep learning is presented, in which deep neural networks are locally embedded in the EnKF. This method is named the deep learning-ensemble Kalman filter (DL-EnKF). The DL-EnKF analysis ensemble is generated from the DL-EnKF analysis and the EnKF analysis deviation ensemble. The performance of the DL-EnKF is investigated through data assimilation experiments in both perfect and imperfect model scenarios using three versions of the Lorenz 96 model and a deterministic EnKF with an ensemble size of 10. Nonlinearity in data assimilation is controlled by changing the time interval between observations. Results demonstrate that despite such a small ensemble the DL-EnKF is superior to the EnKF in terms of accuracy in strongly nonlinear regimes and that the DL-EnKF analysis is more accurate than the output of deep learning due to positive feedback in assimilation cycles. Even if the target of training is an EnKF analysis with a large ensemble or a simulation by an imperfect model, the improvement introduced by the DL-EnKF is not very different from the case where the target of training is the true state.

## 1. Introduction

Data assimilation in nonlinear or non-Gaussian systems has been a challenge in meteorology and other geosciences (Bocquet et al., 2010). For instance, it is well known that cumulus convection exhibits strong non-Gaussianity in data assimilation (e.g., Kondo and Miyoshi, 2019; Kawabata and Ueno, 2020). The ensemble Kalman filter (EnKF) is formulated under the Gaussian assumption and is close to optimal in weakly nonlinear regimes, but it does not work well if nonlinearity is strong.   On the other hand, the particle filter (PF) does not need the Gaussian assumption, but the weight degeneracy had been preventing the use of the PF for high-dimensional data assimilation (Snyder et al., 2008; van Leeuwen 2009). However, this limitation is disappearing due to recent developments in the PF, including the use of localization and hybrids with the EnKF (Farchi and Bosquet, 2018; van Leeuwen et al., 2019). Despite this progress, a relatively large ensemble is still needed for the PF to outperform the EnKF (e.g., Penny and Miyoshi, 2016).   This may be plausible since non-Gaussian data assimilation needs some information on higher-order moments of probability density functions (PDFs). As for the 4-dimensional variational method (4D-Var), Tsuyuki (2014) showed that the 4D-Var with a conventional cost function implicitly used a non-Gaussian prior PDF that evolved according to the Liouville equation (Ehrendorfer, 1994) if a certain condition was satisfied, and that the difficulty caused by multiple minima could be alleviated by combining with the EnKF. The iterative ensemble Kalman filter/smoother (IEnKF/IEnKS) have been shown to be the missing link between the PF and the EnKF and

3

73  4D-Var, and can work very well with mild nonlinearity and generate a much better analysis

74  than the above data assimilation methods (Sakov et al., 2012; Bocquet and Sakov, 2013,

75  2014; Bocquet, 2016). However, the IEnKF/IEnKS need much larger computational cost due

76  to the iterative application of the EnKF/EnKS and the use of a long assimilation window.

77  Recent developments in machine learning, in particular in deep learning (Le cum et al.,

78  2015), have demonstrated impressive skills in various fields. Data-driven modeling,

79  including data-driven parametrizations, based on machine learning has been extensively

80  explored for improving simulations and predictions of nonlinear dynamical systems. Dueben

81  and Bauer (2018) discussed the question of whether models that were based on deep

82  learning and trained on atmospheric data could compete with weather and climate models

83  that were based on physical principles. Reichstein et al. (2019) advocated a hybrid modeling

84  approach in which physical process models were coupled with machine learning to further

85  improve understanding and predictive ability in earth system science. Abarbanel et al. (2018)

86  and Geer (2020) showed an equivalence in formulation between data assimilation and deep

87  learning. Lists of literature of recent studies are available in Reichstein et al. (2019) and

88  Chattopadhyay et al. (2020), for instance. Quite recently the combination of data assimilation

89  and machine learning has been explored to address sparse and noisy observations in data-

90  driven modeling (Brajard et al., 2020a; Bocquet et al., 2020; Tomizawa and Sawada, 2020;

91  Gottwald and Reich, 2021; Wikner et al., 2021), data-driven parametrizations (Brajard et al.,

92  2020b), and model error correction (Farchi et al., 2021).

93  Research on the application of deep learning to data assimilation itself has also started. Arcucci et al. (2021) proposed a method for integrating variational data assimilation with deep learning, in which a recurrent neural network is trained on the state of a dynamical model and the result of data assimilation. Silva et al. (2021) proposed the use of a generative adversarial network to make prediction and to assimilate observations by using a low-dimensional space for the spatial distribution of the simulated state. However, it is difficult to directly apply those methods to data assimilation in high-dimensional systems such as atmospheric models for numerical weather prediction.

101  In this study, we present an ensemble data assimilation method combining the EnKF and deep learning as an alternative to the PF for high-dimensional systems. The additional computational cost to assimilate observations is a very small fraction of that of the EnKF. Since a deep neural network (DNN) can learn a data assimilation method for a specific dynamical system and a specific observing system by training, we could expect this method to work with a relatively small ensemble size even in strongly nonlinear regimes. However, data assimilation in meteorology is generally a large-scale problem, and the background error covariance and the distribution of radar and satellite data change with the analysis time. The EnKF, as well as the PF and 4D-Var, can properly deal with this nonstationarity in data assimilation. On the other hand, deep learning is based on the minimization of the sum of errors over many samples. In addition, it would be difficult to provide sufficient information on the forecast error covariance to a DNN, because the feasible size of a DNN is limited,

113   where we define the size of a DNN as the total number of weights including bias parameters

114   to be optimized by training. If the output of a DNN is not well optimized for each analysis

115   time, the analysis accuracy may deteriorate in assimilation cycles. However, since the EnKF

116   does not work very well in strongly nonlinear regimes, we could expect data assimilation by

117   deep learning to outperform the EnKF in such regimes.

118        The purpose of this study is to propose a nonlinear data assimilation method based on

119   deep learning that is locally embedded in an EnKF and to investigate its performance

120   through data assimilation experiments in both perfect and imperfect model scenarios using

121   toy models. By applying deep learning in combination with an EnKF, we can reduce the size

122   of a DNN and address the nonstationarity in data assimilation. This method is named the

123   deep learning-ensemble Kalman filter (DL-EnKF).

124        The remainder of this paper is organized as follows. Section 2 introduces the method of

125   DL-EnKF. Section3 describes the design of experiments in both perfect and imperfect model

126   scenarios. Section 4 presents the results of these experiments. Summary and discussion

127   are mentioned in Section 5.

128

129   **2. Method**

130        Since data assimilation is generally a large-scale problem, it is desirable to keep the

131   size of a DNN as small as possible. For instance, the size of a feedforward neural network

132   with $m$ layers with $n$ nodes per layer is about $n^2(m\text{-}1)$ and a greater number of samples

133    would be required for training. If we directly apply deep learning to data assimilation, the

134    number of input nodes is at least the sum of the number of observations and the degrees of

135    freedom of a numerical model, and the number of output nodes is the degrees of freedom

136    of the model, while the number of nodes of a hidden layer is usually required to be larger

137    than the number of input or output nodes. For high-dimensional systems such as

138    atmospheric models, the size of a DNN would become too large to be stored in the memory

139    of a computer and to prepare sufficient training samples. To apply deep learning to data

140    assimilation for atmospheric models, we need to introduce a localization procedure and to

141    train a DNN to have some versatility so that it is applicable to each grid point in a certain

142    range of geographical areas.

143        Figure 1a shows the workflow of the DL-EnKF, in which deep learning is locally     Fig. 1

144    embedded in an EnKF. The "EnKF" box in this figure represents the analysis step of the

145    EnKF, and "Deep Learning" box consists of an ensemble of several DNNs. The inputs of

146    DNNs to create the DL-EnKF analysis at a grid point are the EnKF analysis, forecast,

147    observations, availability of observations in binary, and pseudo-observations that

148    supplement missing observations. The EnKF analysis and forecast are the ensemble means

149    of each ensemble. We do not explicitly use the information contained in the forecast

150    ensemble other than the ensemble mean to reduce the size of DNNs. Since observational

151    data for which the DNN has input nodes may be sometimes missing, it is necessary to

152    provide the information on the availability of observations to DNNs. The pseudo-

153   observations are prepared by using the EnKF analysis and the observation operators of the

154   missing observations. Those input data are extracted from a small domain centered at the

155   analysis grid point. The radius of this domain is hereafter referred to as the input radius, and

156   it is assumed that this value is smaller than the covariance localization radius of the EnKF.

157   According to Hsieh and Tang (1998), the DL-EnKF analysis is the average of outputs from

158   the ensemble of DNNs. The individual outputs from DNNs would be scattered in phase

159   space due to multiple minima of a loss function of deep learning, and we would likely obtain

160   a better DL-EnKF analysis by averaging those individual outputs.

161   The analysis ensemble $\{x_m^a\}_{m=1}^M$, where $M$ is the ensemble size, is created by

162   modifying the EnKF analysis ensemble $\{x_{\mathrm{EnKF},\,m}^a\}_{m=1}^M$ such that its ensemble mean is equal

163   to the DL-EnKF analysis $x_{\mathrm{DL-EnKF}}^a$ as follows:

164   $$x_m^a = x_{\mathrm{DL-EnKF}}^a + \alpha\big(x_{\mathrm{EnKF},\,m}^a - x_{\mathrm{EnKF}}^a\big), \tag{1}$$

165   where $x_{\mathrm{EnKF}}^a$ is the EnKF analysis and $\alpha$ is a parameter for adjusting the spread of the

166   analysis ensemble. If adaptive covariance inflation is used in the EnKF, we can set $\alpha$ to 1

167   since the effect of tuning $\alpha$ is almost canceled by this procedure. However, if we conduct

168   ensemble forecasts using the analysis ensemble, we may need to adjust the value of $\alpha$.

169   The members of the analysis ensemble thus generated are evolved by the time integration

170   of a numerical model to prepare the forecast ensemble for the next analysis time.

171   For the training of a DNN, we use the EnKF analysis and forecast provided by an EnKF

172   run as shown in Fig. 1b. The weights including bias parameters are optimized by reducing

173   a loss function that measures a difference between the output of the DNN and the target of

174   training. We prepare several DNNs by randomly initializing the weights before the training.

175       One of the reasons for including the EnKF analysis in the inputs of DNNs is that this

176   analysis at a grid point contains some information on the forecast, observations, and

177   forecast error covariance in a domain within the covariance localization radius, so that we

178   can reduce the input radius of DNNs and implicitly utilize some information of the forecast

179   error covariance. In addition, even if DNNs cannot deal with some observational data

180   because the input nodes for these observations are absent, they are assimilated by the

181   EnKF part of DL-EnKF and their information is partly provided to the deep learning part

182   through the EnKF analysis.

183       We can prepare pseudo-observations by other methods. However, it is easily shown by

184   the sequential assimilation method (e.g., Houtekamer and Mitchell, 2001) that if observation

185   errors are independent of each other the additional assimilation of pseudo-observations

186   does not change the EnKF analysis. Therefore, it can be considered that the pseudo-

187   observations thus created are assimilated in the EnKF part of DL-EnKF along with the real

188   observations, and that the same observations including the pseudo-observations are

189   provided to the deep learning part. In this sense, the method adopted in this study may be

190   a natural choice, although it may not be optimal and spurious correlations and biases will be

191   generated. We could expect that DNNs will learn to properly deal with this problem by

192   training.

193     Lawson and Hansen (2004) showed that an analysis ensemble generated by a

194     deterministic EnKF tends to retain multi-modality that may appear in a forecast ensemble,

195     while this is not the case for a stochastic EnKF. Therefore, a stochastic EnKF is better than

196     a deterministic EnKF for generating an analysis ensemble of the DL-EnKF. However, it is

197     well known that if the ensemble size is relatively small, a stochastic EnKF is inferior in terms

198     of the accuracy of the ensemble mean due to random perturbations that are added to

199     observations (e.g., Sakov and Oke, 2008; Bowler et al., 2013), so that we adopt a

200     deterministic EnKF for the EnKF part in the present paper.

201

202     **3. Design of experiments**

203     *3.1. Outline*

204     The performance of the DL-EnKF is investigated through both perfect and imperfect

205     model experiments using three versions of the 40-variable Lorenz 96 models (Lorenz, 1996)

206     and the serial ensemble square root filter (EnSRF; Whitaker and Hamill, 2002), which is one

207     of the deterministic EnKFs. The ensemble size of the serial EnSRF is set to 10, because we

208     are interested in the performance of the DL-EnKF with a relatively small ensemble. In this

209     and the next sections, the EnKF means the serial EnSRF unless otherwise stated. The

210     purpose of the perfect model experiments is to clarify the basic performance of the DL-EnKF,

211     while that of the imperfect model experiments is to gain insight into the performance of the

212     DL-EnKF when applied to data assimilation in the real atmosphere.

213    The experiments consist of two phases: the training phase of DNNs and the test phase

214    using data assimilation experiments. In the training phase, we run the models and the EnKF

215    to prepare training and validation datasets, which are used to train DNNs and to verify the

216    accuracy of the output of DNNs, respectively. The length of period and time interval of these

217    datasets are 1 000 and 1, respectively. This large time interval is taken to ensure that each

218    data is almost independent of each other. In the test phase, we run the models to prepare a

219    test dataset for the data assimilation experiments. The length of period of this dataset is also

220    1 000. The accuracy of the DL-EnKF analysis is compared with the deep learning and EnKF

221    analyses to evaluate the performance of the DL-EnKF. The workflow to create the deep

222    learning analysis is the same as that of the DL-EnKF analysis except for the absence of

223    feedback from the deep learning part to the EnKF part (Fig. 2). The analysis accuracy is | Fig. 2

224    evaluated by the RMSE that is the square root of the squared error averaged over the grid

225    points and the period of the test dataset at a time interval of 1.

226    In the perfect model experiments, we use the original 40-variable Lorenz 96 model and

227    conduct two types of experiments, Exp-PA and Exp-PB, in which the targets of training are

228    different. The target in Exp-PA is the true state generated by the model, while the target in

229    Exp-PB is an analysis by the stochastic EnKF (Evensen, 1994; Burgers et al., 1998) with an

230    ensemble size of 1 000. This analysis is hereafter referred to as the EnKF1000 analysis.

231    Although this ensemble size may be unrealistic for the 40-variable model, the purpose of

232    Exp-PB is to examine the performance of the DL-EnKF when an analysis with a high

233 accuracy is used as a target.

234    In the imperfect model experiments, the two-scale Lorenz 96 model with 40 large-scale

235 variables and 400 small-scale variables is used as a substitute of the real atmosphere, while

236 a parameterized Lorenz 96 model with a parameterization of large-scale forcing by small-

237 scale variables is used as a substitute of a numerical model of the real atmosphere. We

238 conduct two types of experiments, Exp-IA and Exp-IB. In Exp-IA, we train DNNs using the

239 simulation data generated by the parameterized model, and conduct the data assimilation

240 experiment using observations generated by the two-scale model. The idea behind Exp-IA

241 is that if a dynamical system and a observing system that are used for the taining of a DNN

242 resemble the real-world systems, we could expect that a data assimilation method the DNN

243 has learned by training also works in the real-world applications. In Exp-IB, the two-scale

244 Lorenz 96 model is used for the training and data assimilation experiments in a perfect model

245 scenario for comparison to Exp-IA.

246    Table 1 summarizes the models used in the training and test phases of the experiments. | Table 1 |

247 The following subsections describe further details of the experimental design.

248

249 *3.2. Models*

250    The governing equations of the Lorenz 96 model for the perfect model experiments are

251 $$\frac{dX_k}{dt} = -X_{k-1}(X_{k-2} - X_{k+1}) - X_k + F,$$ (2)

252 for $k = 1, \cdots, K$, satisfying periodic boundary conditions: $X_{-1} = X_{K-1}$, $X_0 = X_K$, and $X_1 =$

12

253     $X_{K+1}$. The number of variables $K$ and the forcing parameter $F$ are set to 40 and 8,

254     respectively. Note that since the number of positive Lyapunov exponents of the model is 13

255     for those parameter values (Lorenz and Emanuel, 1998), the ensemble size of 10 is not very

256     small. The leading Lyapunov exponent corresponds to a doubling time of 0.42 (Lorenz and

257     Emanuel, 1998). When the nonlinearity in data assimilation is controlled by changing the

258     time interval between observations as in the present study, this value can be used as a

259     reference for estimating the degree of nonlinearity. The fourth-order Runge-Kutta scheme is

260     adopted for the time integration of the model with a time step 0.01.

261       The governing equations of the two-scale Lorenz 96 model for the imperfect model

262     experiments are

263
$$\frac{dX_k}{dt} = -X_{k-1}(X_{k-2} - X_{k+1}) - X_k + F - \frac{hc}{b}\sum_{j=1}^{J} Y_{j,k}, \tag{3}$$

264
$$\frac{dY_{j,k}}{dt} = -cbY_{j+1,k}(Y_{j+2,k} - Y_{j-1,k}) - cY_{j,k} + \frac{hc}{b}X_k, \tag{4}$$

265     for $k = 1, \cdots, K$ and $j = 1, \cdots, J$, where $\{X_k\}$ are large-scale variables and $\{Y_{j,k}\}$ are small-

266     scale variables, satisfying periodic boundary conditions: $X_{-1} = X_{K-1}$, $X_0 = X_K$, $X_1 = X_{K+1}$,

267     $Y_{0,1} = Y_{J,K}$, $Y_{J+1,K} = Y_{1,1}$, and $Y_{J+2,K} = Y_{2,1}$. To make Eq. (4) meaningful, we further define

268     $Y_{0,k} = Y_{J,k-1}$, $Y_{J+1,k} = Y_{1,k+1}$, and $Y_{J+2,k} = Y_{2,k+1}$. Large- and small-scale variables interact

269     with each other through the last terms on the right-hand side of Eqs. (3) and (4). We set the

270     parameters as follows: $K = 40$, $J = 10$, $F = 10$, $h = 1$, $c = 10$, and $b = 10$. These values

271     are the same as the ones used by Lorenz (1996) except for $K$. Note that the forcing

272     parameter $F$ is larger than in the perfect model experiments. The fourth-order Runge-Kutta

273    scheme is adopted for the time integration of the model with a time step 0.005.

274        We also need the parameterized Lorenz 96 model in the imperfect model experiments.

275    Although advanced parametrization methods such as stochastic parametrization (e.g., Wilks,

276    2005) and machine learning-based parametrization (e.g., Schneider et al., 2017) are

277    available, a simple function fitting is adopted in the present study; the last term on the right-

278    hand side of Eq. (3) is approximated by a linear function of $X_k$. The reason we adopt such

279    a simple approach is that we intend to demonstrate that even an unsophisticated imperfect

280    model works well for the training of a DNN. Then the governing equations of the

281    parameterized Lorenz 96 model are

282    $$\frac{dX_k}{dt} = -X_{k-1}(X_{k-2} - X_{k+1}) - X_k + F + (a_1 X_k + a_0), \tag{5}$$

283    for $k = 1, \cdots, K$, satisfying periodic boundary conditions: $X_{-1} = X_{K-1}$, $X_0 = X_K$, and $X_1 =$

284    $X_{K+1}$. The number of variables $K$ and the forcing parameter $F$ were set to 40 and 10,

285    respectively, to be consistent with the two-scale Lorenz 96 model. The constants $a_1$ and

286    $a_0$ are to be determined by the function fitting. The fourth-order Runge-Kutta scheme is

287    adopted for the time integration of the model with a time step 0.01.

288        The three models are integrated from $t = 0$ to $t = 2\,050$ for preparing the training and

289    validation datasets. The initial condition at each grid point is $F$ plus an independent random

290    number drawn from the normal distribution with the mean 0 and the variance 1, except that

291    the small-scale variables of the two-scale Lorenz 96 model are set to 0 at the initial time.

292    The data from $t = 51$ to $t = 1\,050$ are used for preparing the training dataset, and those

293      from $t = 1\,051$ to $t = 2\,050$ for the validation dataset. The other time integration of the

294      models from $t = 0$ to $t = 1\,050$ with initial conditions generated by using another random

295      number sequence is conducted for preparing the test dataset, and the state variables from

296      $t = 51$ to $t = 1\,050$ are used as the true state (target) for computing the analysis error.

297

298      *3.3. Observations*

299      Observations are generated by adding random errors to the results of the time

300      integration of the models. The observation errors are independent random draws from the

301      normal distribution with the mean 0 and the variance 1, so that the standard deviation of

302      observation errors is 1. Observations used in the imperfect experiments are of large-scale

303      variables of the two-scale Lorenz 96 model except for the training phase of Exp-IA, in which

304      observations are prepared by the parameterized Lorenz 96 model.

305      Nonlinearity in data assimilation is controlled by changing the time interval between

306      observations $\Delta t$. All experiments are performed for three values of $\Delta t$: 0.05, 0.20, and 0.50.

307      The case of $\Delta t = 0.05$ corresponds to a weakly nonlinear case, and that of $\Delta t = 0.50$

308      corresponds to a strongly nonlinear one. Note that the latter value is close to the doubling

309      time 0.42 mentioned in Subsection 3.2, and that Penny and Miyoshi (2016) used $\Delta t = 0.50$

310      for their experiments of a local PF. All observations are prepared such that observations at

311      the same analysis time are the same regardless of the time interval between observations.

312      For the spatial distribution of observations, two cases are examined. In one case,

313    observations are available at all grid points, and the number of observations is always 40.

314    In other words, observations are available at each grid point with a probability of 1. In the

315    other case, observations are available at each grid point with a probability of 1/2. It is

316    assumed that events that an observation exists are independent of each other in space and

317    time so that the spatial distribution of observations randomly changes at every observation

318    time. The average number of observations is 20, and the standard deviation of the number

319    of observations is $\sqrt{40 \cdot (1/2)^2} \approx 3.16$. Hence, the number of pseudo-observations used by

320    deep learning is about the same as that of observations. The probability of observations is

321    hereafter denoted by $p$.

322

323    *3.4. Data assimilation by EnKF*

324        Covariance localization and covariance inflation are needed to optimize the performance

325    of the EnKF. The correlation function defined by Eq. (4.10) of Gaspari and Cohn (1999) is

326    taken for covariance localization. The parameter $c$ in this equation is regarded as the

327    localization radius $r_L$ (unit: grid intervals) in the present study, at which radius the

328    correlation coefficient decreases to 5/24. An adaptive inflation method proposed by Li et al.

329    (2009) is used for multiplicative covariance inflation. This method is based on the innovation

330    statistics by Desroziers et al. (2005). Li et al. (2009) imposed lower and upper limits in the

331    "observed" inflation factor $\tilde{\Delta}^o$ before applying a smoothing procedure: $0.9 \leq \tilde{\Delta}^o \leq 1.2$. Since

332    we conduct data assimilation over a much wider range of the time interval between

16

333  observations $\Delta t$, we optimize the upper limit of $\tilde{\Delta}^o$ for each set of parameters ($r_L$, $\Delta t$, $p$)

334  leaving the lower limit at 0.9. The candidates of the upper limit are 1.2, 1.3, 1.4, 1.5, 2.0, 3.0,

335  5.0, and no limit. In addition, although Li et al. (2009) set the error growth parameter $\kappa$ to

336  1.03, we adopt a larger value $\kappa$ = 1.1 because this value leads to a better analysis in the

337  present study. A set of values of $r_L$ and the upper limit of $\tilde{\Delta}^o$ with the best analysis accuracy

338  is hereafter referred to as the optimal parameters. We determine the optimal parameters for

339  each pair of ($\Delta t$, $p$) in Exp-PA and Exp-IA by data assimilation experiments using the target

340  and observations in each training dataset. The optimal parameters thus determined are also

341  used in Exp-PB and Exp-IB, respectively, unless otherwise stated.

342      In Exp-PB, the target of training is the EnKF1000 analysis that is yielded by the

343  stochastic EnKF with an ensemble size of 1 000, as mentioned in Subsection 3.1. The

344  reason we adopt the stochastic EnKF is that when an ensemble size is very large the

345  accuracy of the serial EnSRF tends to deteriorate and to become less accurate than the

346  stochastic EnKF. We can avoid this problem with the serial EnSRF by applying the mean-

347  preserving random rotation of an analysis ensemble (Sakov and Oke, 2008), but an

348  additional computational cost is very large for the random rotation of a 1000-member

349  ensemble. Figure 3 compares the analysis accuracy of the two EnKFs for ensemble sizes | Fig. 3 |

350  of 10 (in cold colors) and 1 000 (in warm colors), plotting the RMSEs averaged over the

351  period from $t = 51$ to $t = 1050$. This result is obtained by using the target and observations

352  in the training datasets of Exp-PA. The localization radius and the upper limit of $\tilde{\Delta}^o$ are

17

353    optimized in the case of the ensemble size 10, while no covariance localization is applied

354    and the upper limits of $\tilde{\Delta}^o$ is set to 1.2 in the case of the ensemble size 1 000. It is found

355    from Fig. 3 that in the latter case the RMSE of the stochastic EnKF is smaller than that of

356    the serial EnSRF for all values of the time interval between observations $\Delta t$ and probability

357    of observations $p$. Note that the serial EnSRF with an ensemble size 10 outperforms the

358    serial EnSRF with an ensemble size 1 000 in the three cases: $(\Delta t,\; p) = (0.05, 1),\; (0.05, 1/2),$

359    and $(0.20, 1)$.

360    In Exp-IB, the two-scale Lorenz 96 model is used to assimilate observations of large-

361    scale variables. As noted by Tsuyuki (2019), when the ensemble size is relatively small,

362    forecast correlations between large- and small-scale variables are not reliable. Hence, these

363    forecast correlations are neglected in the EnKF, and the analysis ensemble of small-scale

364    variables is left unchanged from the forecast ensemble at each analysis time.

365

366    *3.5. Deep learning*

367    A simple feedforward neural network with the same number of nodes for all hidden layers

368    is adopted for the deep learning part of the DL-EnKF. As we assume that the input radius of

369    DNNs is relatively small, a convolutional neural network would not be needed. We could use

370    a recurrent neural network instead of the feedforward neural network to utilize the

371    information obtained by the previous processing in deep learning, but it is not adopted in the

372    present study for simplicity.

373     The inputs of a DNN are the EnKF analysis using the optimal parameters, forecast,

374     observations, availability of observations, and pseudo-observations in a small domain

375     centered on an analysis grid point within the input radius $r_I$ (unit: grid intervals). The

376     availability of observation at a grid point is set to 1 if the observation is available and set to

377     -1 if not available. The DNN assumes that observations are always available at all grid points,

378     and we supplement missing observations with pseudo-observations. Since the availability

379     of observations is not necessary in the case of $p = 1$, the input layer of the DNN has $3(2r_I +$

380     $1)$ nodes for $p = 1$, and $4(2r_I + 1)$ nodes for $p = 1/2$. The number of hidden layers is set

381     to 5 or 10 and the number of nodes per hidden layer is optimized as will be mentioned later.

382     Since the balance of analysis (Kalnay, 2003) is not a serious issue in Lorenz 96 models, we

383     let the output of the DNN be the analysis value at the analysis grid point only. For Exp-IB in

384     which the two-scale Lorenz 96 model is used in the EnKF part of DL-EnKF, the input and

385     output of the DNN are of large-scale variables only.

386     Table 2 summarizes the architecture and training of the DNN. All input and output data | Table 2 |

387     except for the availability of observations are normalized by using the mean and standard

388     deviation of the target state in the training dataset of Exp-PA for the perfect model

389     experiments and of Exp-IB for the imperfect model experiments. Since the statistical

390     behavior of the models does not depend on the location of a grid point, the data at all grid

391     points are used to prepare the training and validation datasets. Hence, the number of

392     samples of each dataset is $40 \times 1\,000 = 40\,000$. The training dataset is split into small

19

393    batches called mini-batches that are used to compute the loss function and update the

394    weights of the DNN. Learning rate decay is adopted in the training to avoid the situation in

395    which the DNN converges towards minima in a noisy manner and ends up oscillating far

396    away from actual minima. The number of epochs is the number of times each element in the

397    training dataset is used by the DNN for optimizing the weights. For most of the cases,

398    iterations almost converge within 10 epochs. We use PyTorch (Paszke et al., 2019) as the

399    deep learning software.

400         To determine the optimal number of nodes per hidden layer, we train two DNNs with 5

401    and 10 hidden layers by changing the number of nodes using the training and validation

402    datasets of Exp-PA. Figure 4 plots the RMSEs of the two DNNs against the input radius $r_I$     Fig. 4

403    for the case of $\Delta t = 0.50$ and $p = 1$. The RMSEs are computed by using the validation

404    datasets. For the DNN with 10 hidden layers and 5 nodes per hidden layer, the training fails

405    for six values of the input radius, so that the RMSE of this case is not plotted in Fig. 4b.

406    Since the RMSE is not very different between 5 and 10 hidden layers, we adopt the DNN

407    with 5 hidden layers. It is found from Fig. 4 that when the input radius and number of nodes

408    are large to some extent, the RMSE tends to increase due to the generalization error of deep

409    learning. Since the DNN with 20 nodes per hidden layer has the smallest RMSE for most of

410    the input radii, we set the optimal number of nodes to 20 for the case of $\Delta t = 0.50$ and $p =$

411    1.

412         The optimal numbers of nodes of the DNN with 5 hidden layers are summarized in Fig.     Fig. 5

413   5 by blue bars for $p = 1$ and by cyan bars for $p = 1/2$. They tend to increase as the time

414   interval between observations increases because the estimation of state variables becomes

415   more difficult as nonlinearity increases. Although this result is obtained for Exp-PA, those

416   number of nodes are used in all experiments. We also compute the optimal numbers of

417   nodes for the case where the EnKF analysis is not included in the inputs of the DNN. The

418   RMSEs for this case are plotted in Fig. 5 by red bars for $p = 1$ and by orange bars for $p =$

419   $1/2$. It is found that the inclusion of the EnKF analysis tends to reduce the optimal number

420   of nodes. This is probably because the EnKF analysis plays the role of a first guess and

421   makes it easier to estimate state variables.

422       The appendix discusses the impacts of the increase in the ensemble size of EnKF and

423   the sample size for training on the accuracy of output of a DNN. The result of experiments

424   shows that when the ensemble size of EnKF is increased to 40 the improvement by deep

425   learning is considerably reduced, and that we need to increase the sample size much larger

426   to obtain a larger improvement.

427

428   *3.6 Data assimilation by DL-EnKF*

429       In the data assimilation experiments with the DL-EnKF, the EnKF part of DL-EnKF

430   adopts the optimal parameters, and the DL-EnKF analysis is the average of outputs of 5 or

431   10 DNNs. Since the adaptive covariance inflation is used in the EnKF part, the parameter

432   $\alpha$ in Eq. (1) is set to 1. In the test phase of Exp-IB, the deep learning part receives only

433  large-scale variables from the EnKF part and generates the DL-EnKF analysis of large-scale

434  variables. The analysis ensemble of large-scale variables is modified by using this analysis,

435  while the analysis ensemble of small-scale variables is left unchanged from the one

436  generated by the EnKF part. The RMSEs of the DL-EnKF, deep learning, and EnKF analysis

437  in Exp-IB are computed by using large-scale variables only.

438

439  **4. Results**

440  *4.1 perfect model experiments*

441  The first issue to be clarified is whether deep learning can outperform the EnKF in terms

442  of analysis accuracy. The EnKF is close to optimal in weakly nonlinear regimes, and the

443  deep learning part of DL-EnKF does not explicitly utilize the forecast error covariance. Figure  | Fig. 6 |

444  6a compares the analysis accuracy between deep learning and the EnKF in Exp-PA for all

445  values of $\Delta t$ and $p$, in which the RMSEs are plotted against the RMSE of EnKF for the

446  input radius of 2 grid intervals. The dots indicate the RMSEs of a single DNN, and the short

447  horizontal bars indicate the RMSE of the average of outputs from 5 DNNs. It is found from

448  this figure that all RMSEs of deep learning analysis are the same for each case and that

449  deep learning outperforms the EnKF when $\Delta t = 0.50$. Note that since the EnKF analysis is

450  included in the inputs of DNNs, the accuracy of deep learning analysis does not become

451  worse than that of the EnKF analysis if sufficient training samples are available.

452  The second issue is whether the accuracy of the DL-EnKF analysis is better than that

22

of the EnKF and deep learning analyses. As mentioned in the introduction, the analysis by deep learning is based on the minimization of the sum of errors over many samples and not optimized for each analysis time. Hence, the analysis accuracy may deteriorate during assimilation cycles by the DL-EnKF. Figure 6b compares the analysis accuracy between the DL-EnKF and EnKF for Exp-PA. The dots indicate the RMSEs of DL-EnKF when the output of a single DNN is used as the DL-EnKF analysis, and the horizontal bars indicate the ones when the average of outputs from 5 DNNs is used as the DL-EnKF analysis. It is found that when $\Delta t = 0.05$ the RMSEs of DL-EnKF based on a single DNN are scattered and larger than that of EnKF. In other words, the accuracy of the deep learning analysis shown in Fig. 6a is not maintained during assimilation cycles in a weakly nonlinear case. Taking the average over 5 DNNs does not improve the accuracy very well. When $\Delta t = 0.20$ and $0.50$, on the other hand, the RMSEs of DL-EnKF based on a single DNN become almost the same for each case and taking the average over 5 DNNs slightly improves the accuracy in the case of $\Delta t = 0.50$. In addition, a comparison of Fig. 6a and 6b shows that the RMSE of DL-EnKF is smaller than that of deep learning when $\Delta t = 0.50$ due to positive feedback in assimilation cycles.

We conduct additional experiments in which the ensemble size of DNNs is increased to 10 in Exp-PA. The initial conditions of weights used for the training are different from the ones used in the case of 5 DNNs. Results are presented in Figs. 6c and 6d, and the former figure looks the same as Fig. 6a. The benefit of taking the average over 10 DNNs for the

473    DL-EnKF is evident when $\Delta t = 0.05$, although its analysis accuracy is still lower than that of

474    EnKF. When $\Delta t = 0.20$ and $\Delta t = 0.50$, the RMSEs of DL-EnKF are almost the same as in

475    the case of 5 DNNs. This suggests that an ensemble size of 5 is sufficient except for a

476    weakly nonlinear case. Then, all the results shown below are based on the average of

477    outputs from 5 DNNs, because our interest is primarily in the performance of the DL-EnKF

478    in strongly nonlinear regimes.

479       Figure 7 compares the time sequences of RMSEs of the EnKF (red line) and the DL-    | Fig. 7 |

480    EnKF (green line) in the case of $p = 1$. When $\Delta t = 0.05$ (Fig. 7a), the DL-EnKF is

481    outperformed by the EnKF during the whole period. When $\Delta t = 0.20$ (Fig. 7b), the analysis

482    accuracy of the two methods is close; the correlation coefficient between the two RMSEs

483    computed for the period from $t = 51$ to $t = 1\,050$ is 0.761. When $\Delta t = 0.50$ (Fig. 7c), the

484    EnKF sometimes exhibits a significant deterioration of accuracy, but the DL-EnKF does not

485    show such a tendency. This result demonstrates an excellent performance of the DL-EnKF

486    in strongly nonlinear regimes.

487       The third issue is whether the optimal input radius of deep learning is smaller than the

488    optimal localization radius of the EnKF. Figure 8 plots the RMSEs of EnKF (orange lines),    | Fig. 8 |

489    deep learning (green lines), and DL-EnKF (blue lines) analysis against the input radius for

490    all cases of Exp-PA (solid lines) and Exp-PB (broken lines). The RMSE of EnKF in Exp-PB

491    is the same as the one in Exp-PA. An orange broken line indicates the RMSE of the

492    EnKF1000 analysis that is used for the training in Exp-PB. The optimal localization radius is

24

493    plotted by a red arrow, except for the case of $(\Delta t,\ p) = (0.05, 1/2)$ where the optimal

494    localization radius is 11 grid intervals. The RMSEs of EnKF and deep learning overlap in

495    Figs. 8a and 8b, and the RMSEs except for the EnKF1000 analysis almost overlap in Figs.

496    8c and 8d.

497       When $\Delta t = 0.05$ (Figs. 8a and 8b), the DL-EnKF is outperformed by the EnKF.

498    Reflecting that an ensemble of 5 DNNs is not sufficient (see Fig. 5b), the graphs of the DL-

499    EnKF are not smooth due to large sampling errors. When $\Delta t = 0.20$ (Fig. 8c and 8d) the

500    two data assimilation methods exhibit almost the same accuracy while when $\Delta t = 0.50$

501    (Figs. 8e and 8f) the DL-EnKF outperforms the EnKF irrespective of the input radius. The

502    accuracy of the DL-EnKF analysis is higher than that of the deep learning analysis for the

503    latter case due to positive feedback in assimilation cycles. We can conclude that the input

504    radius of 2 grid intervals is sufficient to attain the best accuracy of the DL-EnKF analysis.

505    This value is smaller than the optimal localization radii for both $p = 1$ and $p = 1/2$. Even if

506    the input radius is further increased, the accuracy of the DL-EnKF and deep learning

507    analysis remains almost the same, although slight degradations are seen due to the

508    generalization error of deep learning. This small sensitivity of RMSEs on the input radius

509    indicates that the information at distant grid points contributes little to the DL-EnKF analysis

510    even within the localization radius of the EnKF. The inclusion of the EnKF analysis in the

511    inputs of DNNs also contributes to this insensitivity.

512       Another point to be noted in Figs. 8e and 8f is that even if DNNs are trained on the

513    EnKF1000 analysis, the accuracy of the DL-EnKF analysis is not very different from the one

514    trained on the true state. Given the large errors of the EnKF1000 analysis shown in Figs. 8e

515    and 8f, this result may look surprising. That is probably because this analysis well represents

516    the basic dynamics of the Lorenz 96 model despite the large errors. If the difference in the

517    ensemble size between the DL-EnKF and the target analysis is decreased, the accuracy of

518    the DL-EnKF analysis in Exp-PB is more deteriorated. For instance, according to an

519    additional experiment in which the ensemble size of the DL-EnKF is set to 40 for the case

520    of $\Delta t = 0.50$ and $p = 1$ (see the appendix), the RMSEs of the EnKF analysis and DL-EnKF

521    analyses in Exp-PA and Exp-PB are 0.682, 0.617, and 0.638, respectively, for the input

522    radius of 2 grid intervals. The corresponding values for the ensemble size of 10 are 0.798,

523    0.675, and 0.689 (see Fig. 8e), so that the deterioration of accuracy in Exp-PB is still not

524    very large.

525        Finally, we examine the impact of including the EnKF analysis in the inputs of DNNs on

526    the accuracy of the deep learning analysis using the test datasets of Exp-PA. Figure 9 plots    Fig. 9

527    the RMSE of deep learning in which the EnKF analysis is not included (cyan line) and the

528    one in which the EnKF analysis is included (green line) against the input radius. The green

529    lines are the same as in Fig. 8, and the two RMSEs overlap in Fig. 9f. For comparison, the

530    RMSE of EnKF of which localization radius is not optimized is also plotted by an orange line

531    against the localization radius with the upper limit of $\tilde{\Delta}^o$ optimized for each localization

532    radius. Note that the RMSE of EnKF does not always attain the minimum at the optimal

26

533  localization radius indicated by a red arrow, because the values of the optimal parameters

534  are determined by using the training datasets.

535  We can see from Fig. 9 that the accuracy of the deep learning analysis is improved by

536  including the EnKF analysis in the inputs of DNNs except for Fig. 9f, in which the EnKF

537  analysis is too inaccurate to be useful. It is also found in Figs. 9c and 9e for $p = 1$ that this

538  procedure reduces the dependence of the analysis accuracy on the input radius. This is

539  because the EnKF analysis contains some information on the forecast ensemble and

540  observations in a domain within the localization radius. Such a reduction in the dependence

541  brought about by including the EnKF analysis is not clearly seen in Figs. 9d and 9f for $p =$

542  1/2, since deep learning partly utilizes the EnKF analysis through pseudo-observations.

543

544  *4.2 Imperfect model experiments*

545  The parametrization procedure for the parameterized Lorenz 96 model is described in

546  Subsection 3.2. Figure 10 is the scatter plot between the large-scale variables and the  Fig. 10

547  forcing. The initial condition is the same as that used for preparing the training dataset of

548  Exp-IB. The number of samples is 40 000 and the result of linear function fitting is plotted by

549  a straight line. The values of constants in Eq. (5) are $a_1 = -0.320$ and $a_0 = -0.165$. Since

550  the slope of this line is negative, the forcing acts on large-scale variables as negative

551  feedback. Figure 11 compares the Hovmöller diagrams of the Lorenz 96 model,  Fig. 11

552  parameterized Lorenz 96 model, and large-scale variables of the two-scale Lorenz 96 model.

553  The initial condition is the same as that used in Fig. 10. Note that the forcing parameter $F$

554  is larger than in the perfect model experiments. A comparison of the three panels in Fig. 11

555  shows that the parameterization works well, although the parametrized Lorenz 96 model

556  evolves a little more regularly than the two-scale Lorenz 96 model. Stochastic

557  parameterizations could remedy this defect (Wilks, 2005).

558      Figure 12 plots the RMSEs of EnKF (orange lines), deep learning (green lines), and DL-     <span style="border:1px solid">Fig. 12</span>

559  EnKF (blue lines) analysis against the input radius for all cases of Exp-IA (solid lines) and

560  Exp-IB (broken lines). Note that Exp-IA is conducted in an imperfect model scenario, while

561  Exp-IB is conducted in a perfect model scenario for comparison. Unlike Fig. 8, an orange

562  broken line indicates the RMSE of EnKF using the two-scale Lorenz 96 model. The optimal

563  localization radius of the EnKF is indicated by a red arrow. The RMSEs of EnKF and deep

564  learning for each case overlap in Figs. 12a-12d. It is found that the RMSE of EnKF using

565  the two-scale Lorenz 96 model is smaller than the one using the parameterized Lorenz 96

566  model. We can confirm that the basic performance of the DL-EnKF is the same as in the

567  perfect model experiments; the DL-EnKF is inferior to the EnKF in a weakly nonlinear case

568  (Figs. 12a and 12b), while the opposite is true in a strongly nonlinear case (Figs. 12e and

569  12f), in which the optimum input radius is smaller than the optimum localization length. A

570  difference from the perfect model experiments is that when $\Delta t = 0.20$ (Figs. 12c and 12d),

571  the accuracy of the DL-EnKF analysis is a little worse than that of the EnKF analysis for $p =$

572  1 and a little better for small values of the input radius for $p = 1/2$.

<div align="center">28</div>

573　　　　An important point to be noted in Figs. 12e and 12f is that even if DNNs are trained on

574　　the training dataset prepared by the parameterized Lorenz 96 model, the improvement in

575　　analysis accuracy introduced by the DL-EnKF is not very different from the case where the

576　　training dataset is prepared by the two-scale Lorenz 96 model. The former model is run in

577　　the data assimilation experiments in the test phase of Exp-IA without any trouble, implying

578　　that this model well represents the basic dynamics of large-scale variables of the two-scale

579　　Lorenz 96 model. When we use the Lorenz 96 model with $F = 10$, of which evolution is

580　　shown in Fig.11a, in the above data assimilation experiments, we often experience failures.

581

582　**5. Summary and discussion**

583　　　　An ensemble data assimilation method based on deep learning was presented, in which

584　　an ensemble of DNNs was locally embedded in an EnKF. This method was named the DL-

585　　EnKF. The inputs of a DNN were the EnKF analysis, forecast, observations, availability of

586　　observations, and pseudo-observations in a small domain centered on an analysis grid point.

587　　Missing observations were supplemented with the pseudo-observations created from the

588　　EnKF analysis. The DL-EnKF analysis was the average of outputs from an ensemble of

589　　DNNs. The DL-EnKF analysis ensemble was generated from the DL-EnKF analysis and the

590　　EnKF analysis deviation ensemble. The members of the DL-EnKF analysis ensemble thus

591　　generated were evolved by the time integration of a numerical model to prepare the forecast

592　　ensemble for the next analysis time.

The performance of the DL-EnKF was investigated through data assimilation experiments in both perfect and imperfect model scenarios using three versions of the Lorenz 96 model and the serial EnSRF with an ensemble size of 10. The target of training in the perfect model experiments was the true state generated by the Lorenz 96 model or the EnKF1000 analysis generated by the stochastic EnKF with an ensemble size of 1000. In the imperfect model experiments, the true state and observations were provided by the two-scale Lorenz 96 model, while the training dataset was prepared by using the parameterized Lorenz 96 model. Nonlinearity in data assimilation was controlled by changing the time interval between observations.

The DL-EnKF was outperformed by the serial EnSRF in a weakly nonlinear case, but it was superior to the serial EnSRF in terms of analysis accuracy in a strongly nonlinear case despite such a small ensemble size. The DL-EnKF analysis was more accurate than the output of deep learning due to positive feedback in assimilation cycles in the latter case. Even if the target of training was the EnKF1000 analysis or the simulation by the parametrized Lorenz 96 model, the improvement introduced by the DL-EnKF was not very different from the case where the target of training was the true state. The inclusion of EnKF analysis in the inputs of DNNs not only improved the accuracy of the deep learning analysis but also reduced the optimal number of nodes per hidden layer and the dependence of the accuracy on the input radius.

Although the above results were obtained from experiments using toy models, they

suggest that the DL-EnKF may be a promising methods for data assimilation in strongly

nonlinear regimes. The DL-EnKF works with a relatively small ensemble size compared to

the PF, and we can prepare a training dataset for deep learning from simulation data by a

numerical model used in data assimilation. Observational data and EnKF analysis data

generated with a large ensemble could be used for this purpose, but a huge computational

cost may be needed to obtain sufficient samples and a period when observations are

available is limited.

The DL-EnKF may be suitable for data assimilation in cloudy or convective regions in

the atmosphere to assimilate radar and satellite observations. We need to extend the inputs

and output of DNNs in the vertical to assimilate satellite radiance data, since they are

nonlocal observations. As for radial wind data by a Doppler radar, the direction and distance

of a radar site differ depending on a grid point. However, if the radar site position relative to

the grid point is included in the inputs of a DNN, we can train the DNN collectively regardless

of the grid point as in the present study.

There are a couple of issues to be addressed before applying the DL-EnKF to data

assimilation in the atmosphere. In the data assimilation experiments using Lorenz 96 models,

the analysis value at a single grid point is sufficient for the output of a DNN, but we need to

take the balance of analysis into account for atmospheric data assimilation. One of the

methods for ensuring the balance is to extend the output of a DNN to include analysis values

at surrounding grid points. Then the target of training consists of the target state in a small

633    domain centered on an analysis grid point. Since the target state is usually well balanced,

634    the DNN could learn the balance. Adding a penalty term for suppressing imbalance to a loss

635    function of the DNN may help enhance the balance. In addition, taking a weighting average

636    of the outputs in adjacent domains may be effective in improving analysis accuracy.

637    This study demonstrates that the DL-EnKF is inferior to the EnKF in a weakly nonlinear

638    case. It is found that an increase in the ensemble size of DNNs can mitigate this problem,

639    but it would be difficult to increase the ensemble size sufficiently, given the computational

640    cost needed for the training of DNNs. We may need a criterion for replacing the EnKF

641    analysis with the corresponding DL-EnKF analysis in DL-EnKF assimilation cycles. An

642    advanced DNN such as a recurrent neural network would be useful for improving the

643    performance of DL-EnKF in weakly nonlinear regimes as well as in strongly nonlinear ones.

644

645    **Data Availability Statement**

646    The Python programs used for Exp-PA in this study are available on J-STAGE Data.

647

653 Fugaku (Large Ensemble Atmospheric and Environmental Prediction for Disaster

654 Prevention and Mitigation; hp200128, hp210166), and by the Japan Society for the

655 Promotion of Science through KAKENHI (Study on Uncertainty of Cumulonimbus Initiation

656 and Development Using Particle Filter; JP17H02962).

657

658 **Appendix**

659     In this study, we perform the experiments using the serial EnSRF with 10 members and

660 the 40 000 training samples. It may be of interest to examine how the accuracy of output of

661 a DNN changes when the ensemble size and the sample size are increased. This appendix

662 presents some results of additional experiments in which the ensemble size of EnKF is set

663 to 10 and 40 and the sample size is set to 40 000, 160 000 and 640 000. These experiments

664 correspond to Exp-PA for $\Delta t = 0.50$. The ensemble size of 40 is the same as the degrees

665 of freedom of the Lorenz 96 model and, according to Fig. 5 of Penny and Miyoshi (2016),

666 an EnKF still outperforms a local PF with this ensemble size. The periods of time integration

667 of the model for preparing the training and validation datasets are 2 050, 8 050, and 32 050

668 with the first periods of 50 in length are discarded.

669     Figure A1 shows the optimum numbers of nodes per hidden layer of a DNN, obtained  | Fig. A1 |

670 by using the validation datasets. The maximum number of nodes is limited to 100. Although

671 we choose the number of nodes that performs the best for the various input radius, the

672 determination of the optimal number becomes difficult with the increase of the sample size.

33

When the number of training samples is increased, the generalization error of deep learning

tends to reduce and, therefore, the optimal number of nodes per hidden layer tends to

increase.

Figure A2 compares RMSE between the serial EnSRF and the output of a DNN obtained [Fig. A2]

by using the test datasets of Exp-PA. Note that they are not the average over 5 DNNs, so

that the RMSEs shown by green lines in Figs. A2a and A2b are different from those in Figs.

8e and 8f, respectively. The optimal localization radius of the serial EnSRF in Fig. A2c is 12

grid intervals. It is found from this figure that the improvement by deep learning is

considerably reduced for the ensemble size of 40. When the sample size is increased, the

RMSE of the output of a DNN is reduced, but we need much more training samples to obtain

a large improvement.


**References**

Abarbanel, H. D., P. J. Rozdeba, and S. Shirman, 2018: Machine learning: deepest learning

    as statistical data assimilation problems. *Neural Computation,* **30**, 2025–2055.

Arcucci, R., J. Zhu, S. Hu and Y.-K. Guo, 2021:   Deep data assimilation: Integrating deep

    learning with data assimilation. *Appl. Sci.*, **11**, 1114-1134. doi:10.3390/app11031114.

Bocquet, M., 2016: Localization and the iterative ensemble Kalman smoother. *Quart. J. Roy.*

    *Meteor. Soc.*, **142**, 1075–1089. doi:10.1002/qj.2711.

Bocquet, M., J. Brajard, A. Carrassi and L. Bertino, 2020: Bayesian inference of chaotic dynamics by merging data assimilation, machine learning and expectation-maximization. *Foundations of Data Science*, **2**, 55–80.

Bocquet, M., C. A. Pires and L. Wu, 2010: Beyond Gaussian statistical modeling in geophysical data assimilation. *Mon. Wea. Rev*., **138**, 1997-3023. doi:10.1175/2010MWR3164.1.

Bocquet, M. and P. Sakov, 2013: Joint state and parameter estimation with an iterative ensemble Kalman smoother. *Nonlin. Processes Geophys.,* **20**, 803–818. doi:10.5194/npg-20-803-2013.

Bocquet, M. and P. Sakov, 2014: An iterative ensemble Kalman smoother. *Quart. J. Roy. Meteor. Soc.,* **140**, 1521–1535. doi:10.1002/qj.2236.

Bowler, N. E., J. Flowerdew, and S. R. Pring, 2013: Tests of different flavours of EnKF on a simple model. *Quart. J. Roy. Meteor. Soc.*, **139**, 1505-1519. doi:10.1002/qj.2055.

Brajard, J., A. Carrassi, M. Bocquet and L. Bertino, 2020a: Combining data assimilation and machine learning to emulate a dynamical model from sparse and noisy observations: a case study with the Lorenz 96 model. *J. Comp. Sci.*, **44**, 101171.

Brajard, J., A. Carrassi, M. Bocquet and L. Bertino, 2020b: Combining data assimilation and machine learning to infer unresolved scale parametrisation. *Philos. Trans. Roy. Soc. London, Ser. A*, **379**, 2020.0086. doi:10.1098/rsta.2020.0086.

Burgers, G., P. J. van Leeuwen, and G. Evensen, 1998: Analysis schemes in the ensemble

713       Kalman filter. *Mon. Wea. Rev.*, **126**, 1719-1724.

714    Chattopadhyay, A., P. Hassanzadeh and D. Subramanian, 2020: Data-driven prediction of a

715       multi-scale Lorenz 96 chaotic system using deep learning methods: Reservoir

716       computing, artificial neutral network, and long short-term memory network. *Nonlinear*

717       *Processes Geophys.*, **27**, 373–389. doi:10.5194/npg-27-373-2020.

718    Desroziers, G., L. Berre, B. Chapnik, and P. Poli, 2005: Diagnosis of observation,

719       background and analysis-error statistics in observation space. *Quart. J. Roy. Meteor.*

720       *Soc.*, **131**, 3385–3396.

721    Dueben, P. D. and P. Bauer, 2018: Challenges and design choices for global weather and

722       climate models based on machine learning, *Geosci. Model Dev*., **11**, 3999–4009.

723       doi:10.5194/gmd-11-3999-2018.

724    Ehrendorfer, M., 1994: The Liouville equation and its potential usefulness for the prediction

725       of forecast skill. Part I: Theory. *Mon. Wea. Rev.*, **122**, 703-713.

726    Evensen, G., 1994: Sequential data assimilation with a nonlinear quasigeostrophic model

727       using Monte Carlo methods to forecast error statistics. *J. Geophys., Res.*, **99**, 10143-

728       10162.

729    Farchi, A. and M. Bocquet, 2018: Review article: Comparison of local particle filters and new

730       implementations, *Nonlinear Processes Geophys.*, **25**, 765-807. doi:10.5194/npg-25-

731       765-2018.

732    Farchi, A. M. Laloyaux, M. Bonavita and M. Bocquet, 2021: Using machine leaning to correct

733     model error in data assimilation and forecast applications. *Quart. J. Roy. Meteor. Soc.*,

734     **147**, 3067-3084. doi:10.1002/qj.4116.

735 Gaspari, G. and S. E. Cohn, 1999: Construction of correlation functions in two and three

736     dimensions. *Quart. J. Roy. Meteor. Soc.*, **125**, 723-757.

737 Geer, A. J., 2020: Learning earth system models from observations: machine learning or

738     data assimilation? *Technical Memorandom 863*, ECMWF, 23pp.

739 Gottwald, G. A. and Reich, S. (2021). Supervised learning from noisy observations:

740     Combining machine-learning techniques with data assimilation. *Physica D: Nonlinear*

741     *Phenomena*, **423**, 132911. doi:10.1016/j.physd.2021.132911.

742 Houtekamer, P. L. and H. L. Mitchell, 2001: A sequential ensemble Kalman filter for

743     atmospheric data assimilation. *Mon. Wea. Rev.*, **129**, 123-137.

744 Hsieh, W. W. and B. Tang, 1998: Applying neural network models to prediction and data

745     analysis in meteorology and oceanography. *Bull. Amer. Meteor. Soc.*, **79**, 1855-1870.

746 Kalnay, E., 2003: *Atmospheric Modeling, Data Assimilation and Predictability*. Cambridge

747     University Press, Cambridge, 364 pp.

748 Kawabata, T. and G. Ueno, 2020: Non-Gaussian probability densities of convection initiation

749     and development investigated using a particle filter with a storm-scale numerical

750     weather prediction model. *Mon. Wea. Rev.*, **148**, 3–20. doi:10.1175/MWR-D-18-0367.1.

751 Kingma, D. and J. Ba, 2015: Adam: A method for stochastic optimization. arXiv:1412.6980.

752 Kondo, K. and T. Miyoshi, 2019: Non-Gaussian statistics in global atmospheric dynamics: a

753    study with a 10 240-member ensemble Kalman filter using an intermediate atmospheric

754    general circulation model. *Nonlinear Processes Geophys.*, **26**, 211–225.

755    doi:10.5194/npg-26-211-2019.

756  Lawson, G. W. and J. A. Hansen, 2004: Implications of stochastic and deterministic filters

757    as ensemble-based data assimilation methods in varying regimes of error growth. *Mon.*

758    *Wea. Rev.*, **132**, 1966-1981.

759  Le Cum, Y., Y. Bengio and G. Hinton, 2016: Deep learning. *Nature*, **521**, 436-444.

760    doi:10.1038/nature14539.

761  Li, H., E. Kalnay and T. Miyoshi, 2009: Simultaneous estimation of covariance inflation and

762    observation errors within an ensemble Kalman filter. *Quart. J. Roy. Meteor. Soc.*, **135**,

763    523-533. doi:10.1002/qj.371.

764  Lorenz, E. D., 1996: Predictability – A problem partly solved, *Proceedings of the ECMWF*

765    *Seminar on Predictability (4-9 September 1995, Reading, UK)*, ECMWF, 1-18.

766  Lorenz, E. D. and K. A. Emanuel, 1998: Optimal sites for supplementary weather

767    observations: Simulation with a small model. *J. Atmos. Sci.*, **55**, 399-414.

768  Paszke, A., S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N.

769    Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A.

770    Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai and S. Chintala, 2019: PyTorch: An

771    imperative style, high-performance deep learning library. In H. Wallach et al.,

772    eds. *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc.,

773    8024–8035. Available at: http://papers.neurips.cc/paper/9015-pytorch-an-imperative-

774    style-high-performance-deep-learning-library.pdf.

775    Penny, S. G. and T. Miyoshi, 2016: A local particle filter for high-dimensional geophysical

776    systems, *Nonlinear Processes Geophys.*, **23**, 391-405. doi:10.5194/npg-23-391-2016.

777    Reichstein, M., G. Camps-Vallis, B. Stevens, M. Jung, J. Denzler and N. Carvalhais and

778    Prabhat, 2019: Deep learning and process understanding for data-driven Earth system

779    science, *Nature*, **566**, 195–204. doi:10.1038/s41586-019-0912-1.

780    Sakov, P. and P. R. Oke, 2008: Implications of the form of the ensemble transformation in

781    the ensemble square root filters. *Mon. Wea. Rev.*, **136**, 1042–1053.

782    doi:10.1175/2007MWR2021.1.

783    Sakov, P., D. S. Oliver and L. Bertino, 2012: An iterative EnKF for strongly nonlinear systems.

784    *Mon. Wea. Rev.*, **140**, 1988–2004. doi:10.1175/MWR-D-11-00176.1.

785    Schneider, T., S. Lan, A. Stuart and J. Teixeira, 2017: Earth system modeling 2.0: A blueprint

786    for models that learn from observations and targeted high-resolution simulations.

787    *Geophys. Res. Lett.*, **44**, 12396–12417. doi:10.1002/2017GL076101.

788    Silva, V. L., C. E. Heaney, Y. Li and C. C. Pain, 2021: Data Assimilation Predictive GAN

789    (DAPredGAN): Applied to determine the spread of COVID-19. arXiv:2105.07729.

790    Snyder, C, T. Bengtsson, P. Bickel and J. Anderson, 2008: Obstacles to high dimensional

791    particle filtering. *Mon. Wea. Rev.* **136**, 4629–4640. doi:10.1175/2008MWR2529.1

792    Tomizawa, F. and Y. Sawada, 2020: Combining ensemble Kalman filter and reservoir

computing to predict spatio-temporal chaotic systems from imperfect observations and models. *Geosci. Model Dev. Discuss.* doi:10.5194/gmd-2020-211.

Tsuyuki, T., 2014: Deterministic predictability of the most probable state and reformulation of variational data assimilation. *J. Meteor. Soc. Japan*, **92**, 599-622. doi:10.2151/jmsj.2014-606.

Tsuyuki, T., 2019: Ensemble Kalman filtering based on potential vorticity for atmospheric multi-scale data assimilation. *J. Meteor. Soc. Japan*, **97**, 1191-1210. doi:10.2151/jmsj.2019-067.

Whitaker, J. S. and T. M. Hamill, 2002: Ensemble data assimilation without perturbed observations. *Mon. Wea. Rev.*, **130**, 1913-1924.

Wikner, A., J. Pathak, B. R. Hunt, I. Szunyogh, M. Girvan and E. Ott, 2021: Using data assimilation to train a hybrid forecast system that combines machine-learning and knowledge-based components. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, **31**, 053114.

Wilks, D. J., 2005: Effects of stochastic parametrizations in the Lorenz '96 system. *Quart. J. Roy. Meteor. Soc.*, **131**, 389-407. doi:10.1256/qj.04.03.

van Leeuwen, P. J., 2009: Particle filtering in geophysical systems. *Mon. Wea. Rev.*, **137**, 4089‑4114. doi:10.1175/2009MWR2835.1.

van Leeuwen, P. J., H. R. Kunsch, L. Nerger, R. Potthast and S. Reich, 2019: Particle filters for high-dimensional geoscience applications: A review. *Quart. J. Roy. Meteor. Soc.*, **145**,

813     2335-2365. doi:10.1002/qj.3551.
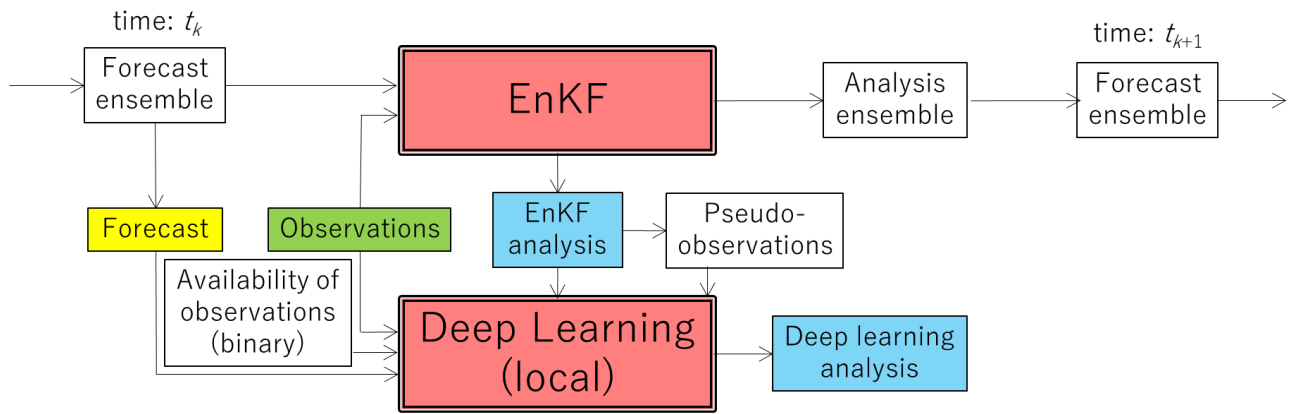
814

815

816

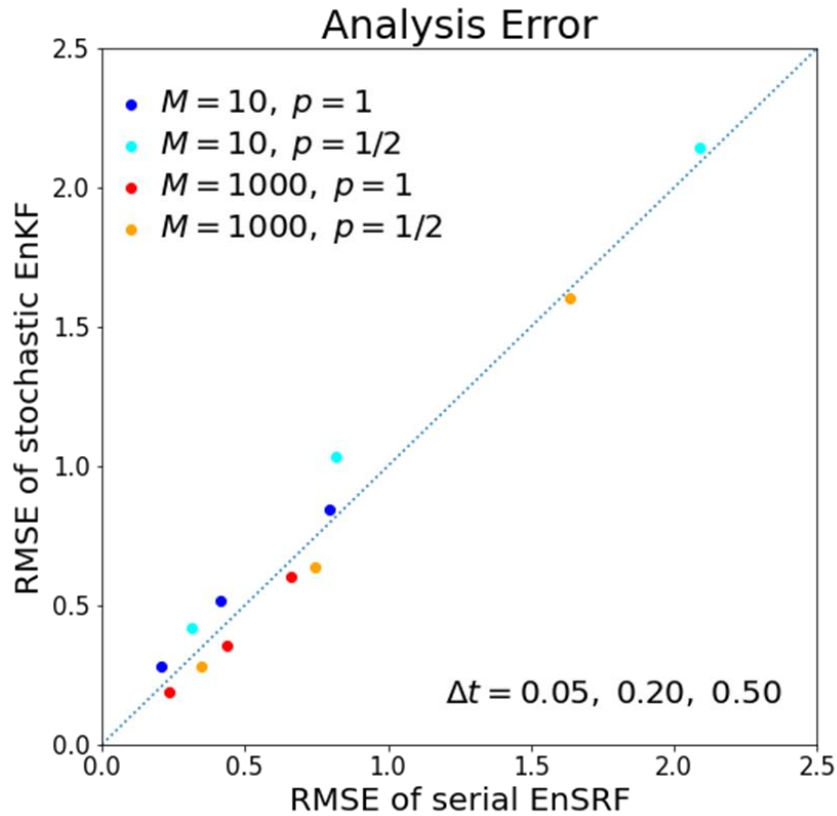Fig. 1    (a) Workflow of DL-EnKF and (b) workflow to train a DNN for DL-EnKF. See text

for details.
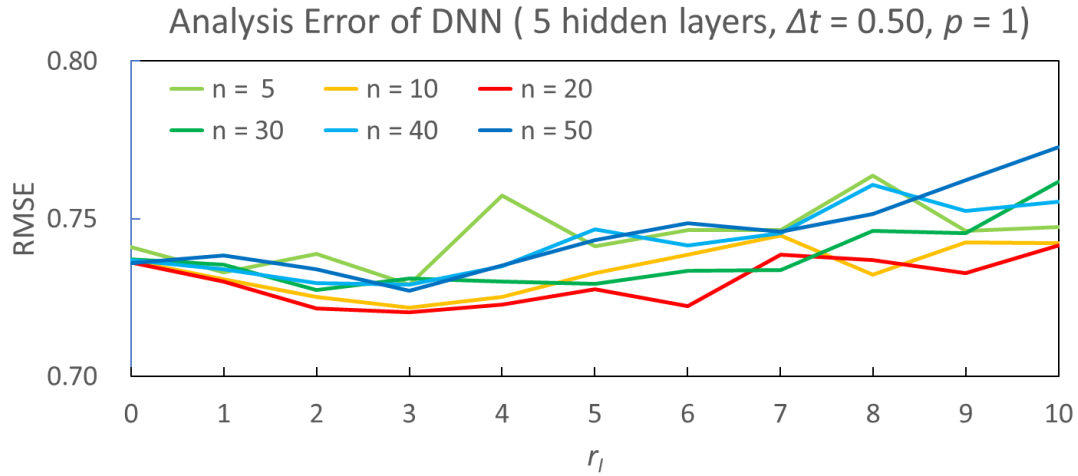
819
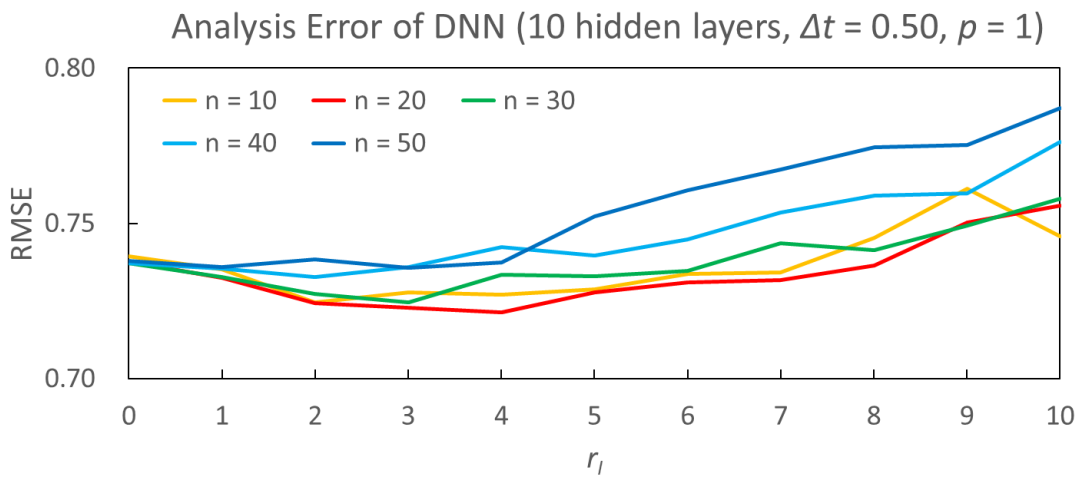
Fig. 2    Workflow to generate deep learning analysis.

821

Fig. 3    Comparison of RMSEs between the serial EnSRF (abscissa) and stochastic EnKF

(ordinate) for the training dataset of Exp-PA. Dots in cold colors are for an ensemble size

10 and dots in warm colors are for an ensemble size 1 000. Dots in dark colors are for

the probability of observations 1 and dots in light colors are for the probability of

observations 1/2. The three dots in the same color correspond to the observation time

intervals 0.05, 0.20, and 0.50 from left to right.

828
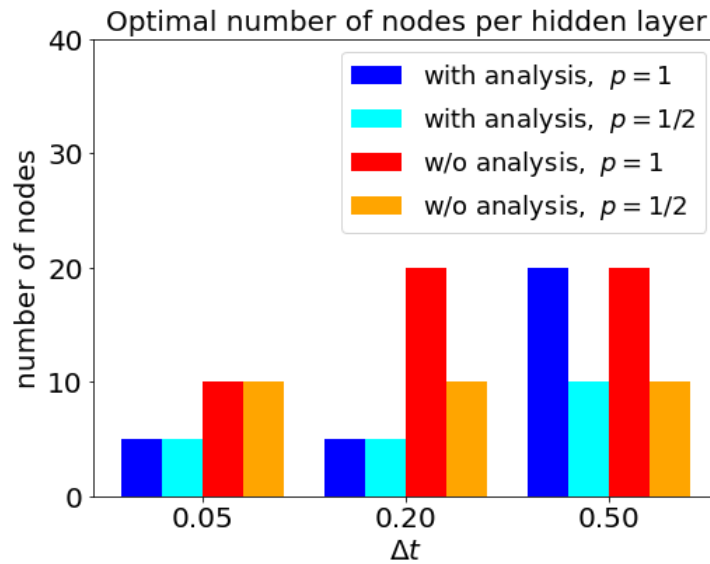
Fig. 4    Comparison of RMSEs of a single DNN with 5 (light green), 10 (orange), 20 (red),

30 (green), 40 (cyan), and 50 (blue) nodes per hidden layer for the validation dataset of

Exp-PA. The RMSEs are plotted against the input radius. The number of hidden layers is

(a) 5 and (b) 10. The observation time interval is 0.50 and the probability of observations
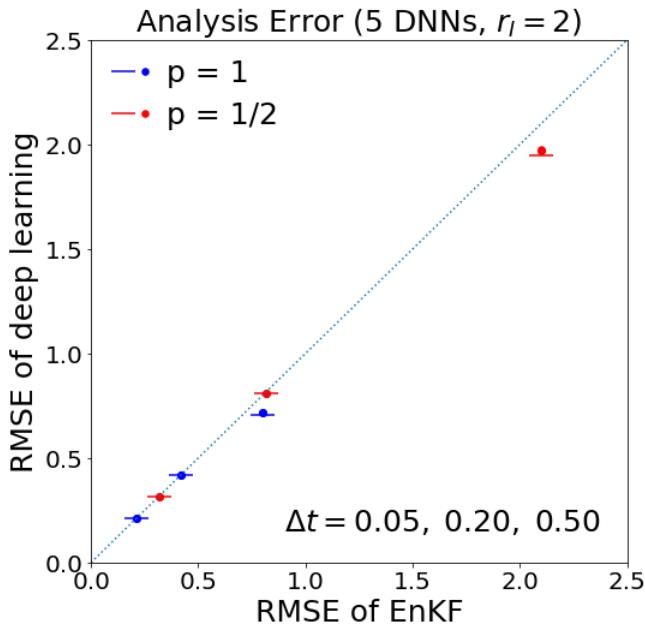
is 1.

834

835    Fig. 5    The optimal number of nodes per hidden layer of a DNN with 5 hidden layers. The

836    abscissa is the observation time interval. Blue and cyan bars are for the case of including

837    EnKF analysis in input for the probability of observations 1 and 1/2, respectively. Red and

838    orange bars are for the case of not including EnKF analysis in input for the probability of

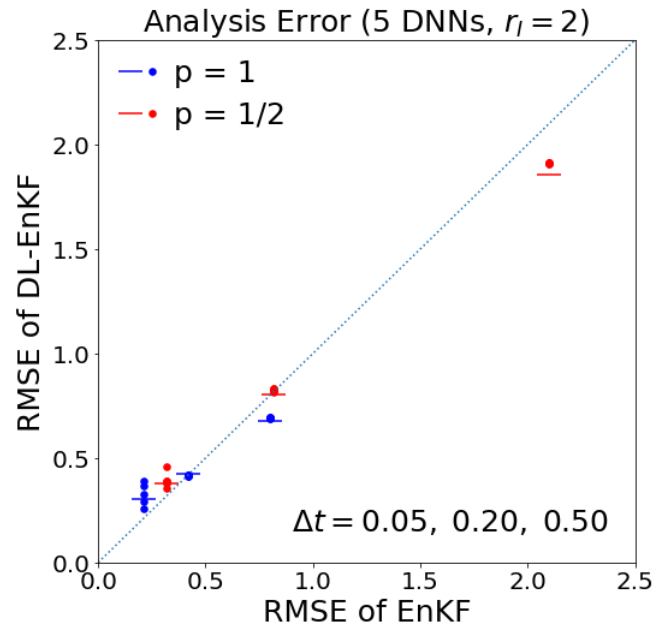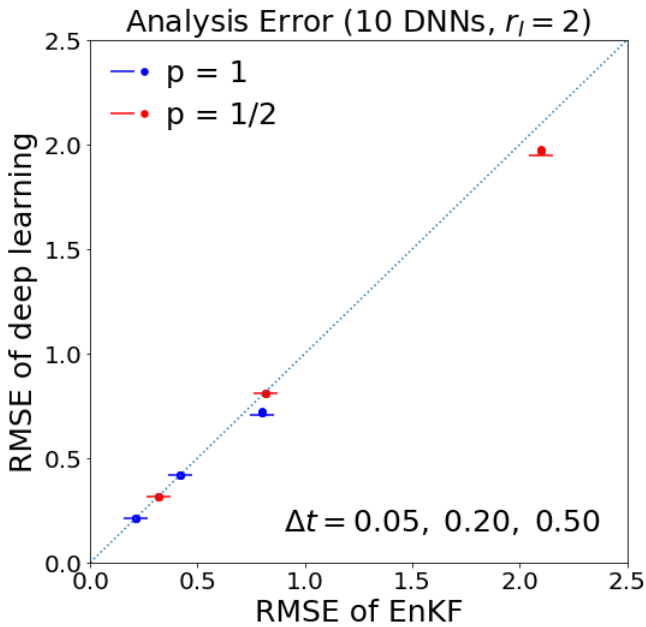839    observations 1 and 1/2, respectively.

840

841    Fig. 6    Comparison of RMSEs between (a) EnKF (abscissa) and deep learning with 5

842    DNNs (ordinate), (b) EnKF and DL-EnKF with 5 DNNs, (c) EnKF and deep learning with

843    10 DNNs, and (d) EnKF and deep learning with 10 DNNs for Exp-PA. The probability of

844    observations 1 is in blue and 1/2 in red, and the input radius is 2 grid intervals. Dots

845    indicate RMSEs based on a single DNN, and short horizontal bars indicate RMSEs based

846    on an ensemble of DNNs. The three groups of dots and a horizontal bar in the same color

847    correspond to the observation time intervals 0.05, 0.20, and 0.50 from left to right.

848

849    Fig. 7    Time sequences of RMSEs of EnKF (red lines) and DL-EnKF (blue lines) for

850    observation time interval (a) 0.05, (b) 0.20, and (c) 0.50 for Exp-PA. The probability of

851    observations is 1 and the input radius is 2 grid intervals.

852

853    Fig. 8    Comparison of RMSEs of EnKF (orange lines), deep learning (green lines), and

854    DL-EnKF (blue lines) for Exp-PA (solid lines) and Exp-PB (broken lines). An orange

855      broken line indicates the RMSE of EnKF1000 analysis used for training in Exp-PB. The

856      RMSEs are plotted against the input radius, and a red arrow indicates the optimal

857      localization radius of EnKF. The observation time interval and the probability of

858      observations are (a) 0.05 and 1, (b) 0.05 and 1/2, (c) 0.20 and 1, (d) 0.20 and 1/2, (e)

859      0.50 and 1, and (f) 0.50 and 1/2, respectively.

860

861 Fig. 9      Comparison of RMSEs of EnKF (orange line), deep learning not including EnKF

862      analysis in input (cyan line), and deep learning including the EnKF analysis in input

863      (green line) for Exp-PA. The RMSE of EnKF is computed for each localization radius. The

864      abscissa is the input radius for deep learning and the localization radius for the EnKF. A

865      red arrow indicates the optimal localization radius. The observation time interval and the

866      probability of observations are (a) 0.05 and 1, (b) 0.05 and 1/2, (c) 0.20 and 1, (d) 0.20

867      and 1/2, (e) 0.50 and 1, (f) 0.50 and 1/2, respectively.

868

869 Fig. 10      Scatter plot between large-scale variables (abscissa) and large-scale forcing by

870      small-scale variables (ordinate) of the two-scale Lorenz 96 model. A solid line is the result

871      of linear function fitting.

872

873 Fig. 11      Hovmöller diagrams of (a) the Lorenz 96 model, (b) the parametrized Lorenz 96

874      model, and (c) large-scale variables of the two-scale Lorenz 96 model.

875

Fig. 12    Same as Fig. 8 except for Exp-IA (solid lines) and Exp-IB (broken lines) and that

an orange broken line indicates the RMSE of EnKF using the two-scale Lorenz 96 model.

878

Fig. A1    The optimal number of nodes per hidden layer of a DNN with 5 hidden layers for

the time interval between observations of 0.50. The abscissa is the number of samples.

Blue and cyan bars are for the EnKF ensemble size of 10 for the probability of

observations 1 and 1/2, respectively. Red and orange bars are for the EnKF ensemble

size of 40 for the probability of observations 1 and 1/2, respectively.

884

Fig. A2    Comparison of RMSE between EnKF (orange lines) and the output of a DNN with

the number of samples of 40 000 (green lines), 160 000 (blue), and 640 000 (cyan) for

the observation time interval of 0.50. The ensemble size of EnKF and the probability of

observations are (a) 10 and 1, (b) 10 and 1/2, (c) 40 and 1, (d) 40 and 1/2, respectively.

The RMSEs are plotted against the input radius, and a red arrow indicates the optimal

localization radius of EnKF.

891

892

893 (a)

894

895 (b)

896

897 Fig. 1    (a) Workflow of DL-EnKF and (b) workflow to train a DNN for DL-EnKF. See text

898    for details.

899

900

901    Fig. 2    Workflow to generate deep learning analysis.

902

903



904

905     Fig. 3     Comparison of RMSEs between the serial EnSRF (abscissa) and stochastic EnKF

906        (ordinate) for the training dataset of Exp-PA. Dots in cold colors are for an ensemble size

907        10 and dots in warm colors are for an ensemble size 1 000. Dots in dark colors are for

908        the probability of observations 1 and dots in light colors are for the probability of

909        observations 1/2. The three dots in the same color correspond to the observation time

910        intervals 0.05, 0.20, and 0.50 from left to right.

911

912 (a)



913 (b)



914

915 Fig. 4    Comparison of RMSEs of a single DNN with 5 (light green), 10 (orange), 20 (red),

916     30 (green), 40 (cyan), and 50 (blue) nodes per hidden layer for the validation dataset of

917     Exp-PA. The RMSEs are plotted against the input radius. The number of hidden layers is

918     (a) 5 and (b) 10. The observation time interval is 0.50 and the probability of observations

919     is 1.

920

Fig. 5    The optimal number of nodes per hidden layer of a DNN with 5 hidden layers. The abscissa is the observation time interval. Blue and cyan bars are for the case of including EnKF analysis in input for the probability of observations 1 and 1/2, respectively. Red and orange bars are for the case of not including EnKF analysis in input for the probability of observations 1 and 1/2, respectively.

Fig. 6 Comparison of RMSEs between (a) EnKF (abscissa) and deep learning with 5 DNNs (ordinate), (b) EnKF and DL-EnKF with 5 DNNs, (c) EnKF and deep learning with 10 DNNs, and (d) EnKF and deep learning with 10 DNNs for Exp-PA. The probability of observations 1 is in blue and 1/2 in red, and the input radius is 2 grid intervals. Dots indicate RMSEs based on a single DNN, and short horizontal bars indicate RMSEs based on an ensemble of DNNs. The three groups of dots and a horizontal bar in the same color correspond to the observation time intervals 0.05, 0.20, and 0.50 from left to right.

940　(a)



Analysis Error ( $\Delta t = 0.05$, $p = 1$, $r_l = 2$ )

941　(b)



Analysis Error ( $\Delta t = 0.20$, $p = 1$, $r_l = 2$ )

942　(c)



Analysis Error ( $\Delta t = 0.50$, $p = 1$, $r_l = 2$ )

943

944　Fig. 7　　Time sequences of RMSEs of EnKF (red lines) and DL-EnKF (blue lines) for

945　　　　observation time interval (a) 0.05, (b) 0.20, and (c) 0.50 for Exp-PA. The probability of

946　　　　observations is 1 and the input radius is 2 grid intervals.

947

Fig. 8    Comparison of RMSEs of EnKF (orange lines), deep learning (green lines), and DL-EnKF (blue lines) for Exp-PA (solid lines) and Exp-PB (broken lines). An orange broken line indicates the RMSE of EnKF1000 analysis used for training in Exp-PB. The RMSEs are plotted against the input radius, and a red arrow indicates the optimal localization radius of EnKF. The observation time interval and the probability of observations are (a) 0.05 and 1, (b) 0.05 and 1/2, (c) 0.20 and 1, (d) 0.20 and 1/2, (e) 0.50 and 1, and (f) 0.50 and 1/2, respectively.

959 

960

961

54

Fig. 10    Scatter plot between large-scale variables (abscissa) and large-scale forcing by small-scale variables (ordinate) of the two-scale Lorenz 96 model. A solid line is the result of linear function fitting.

(a)

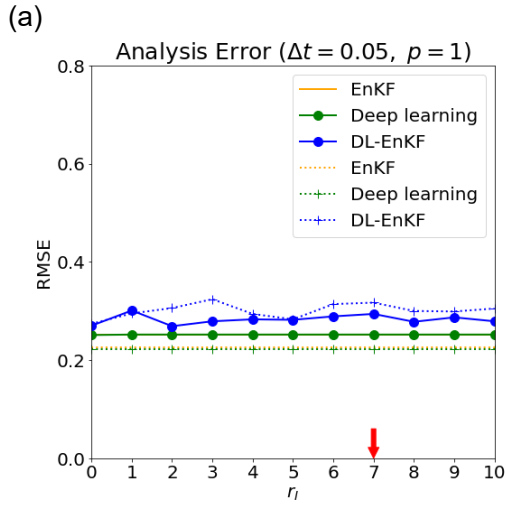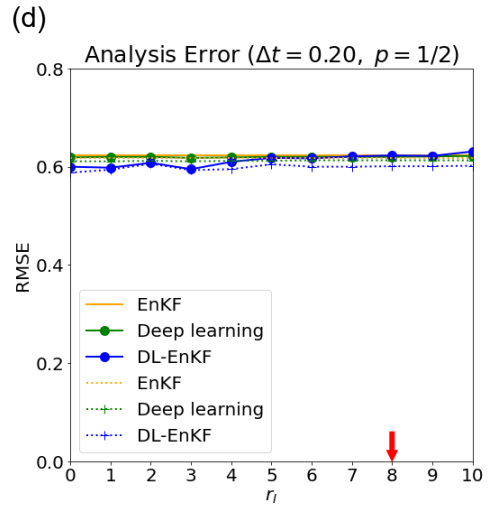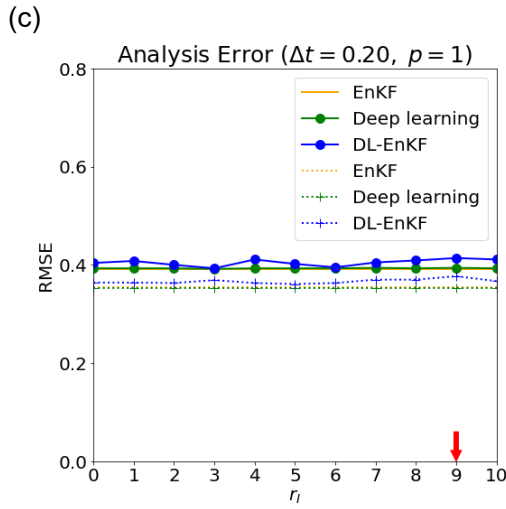Lorenz-96 model w/o parameterization

(b)

Lorenz-96 model with parameterization

(c)

Two-scale Lorenz-96 model (large scale)

Fig. 11    Hovmöller diagrams of (a) the Lorenz 96 model, (b) the parametrized Lorenz 96

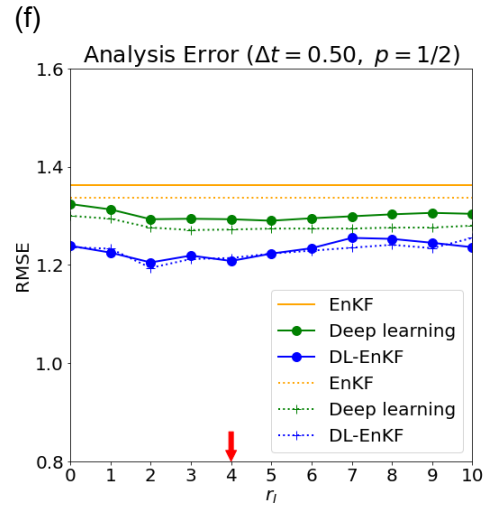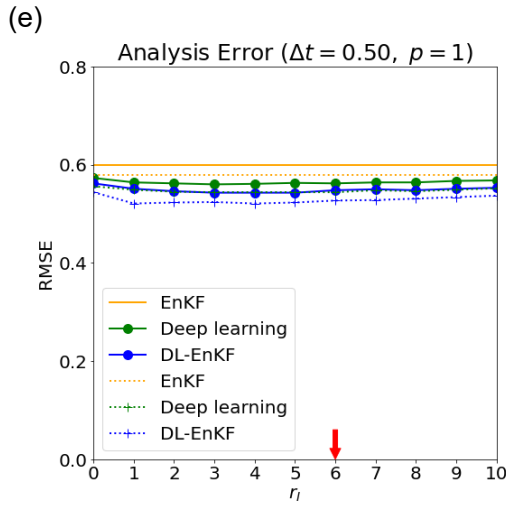model, and (c) large-scale variables of the two-scale Lorenz 96 model.

981

(a) Analysis Error ($\Delta t = 0.05$, $p = 1$)

(b) Analysis Error ($\Delta t = 0.05$, $p = 1/2$)

982

(c) Analysis Error ($\Delta t = 0.20$, $p = 1$)

(d) Analysis Error ($\Delta t = 0.20$, $p = 1/2$)

983

(e) Analysis Error ($\Delta t = 0.50$, $p = 1$)

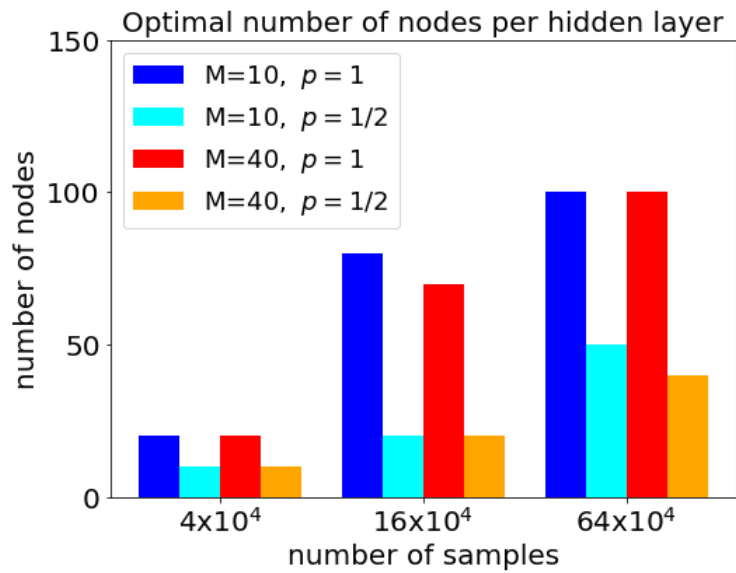(f) Analysis Error ($\Delta t = 0.50$, $p = 1/2$)

984

985    Fig. 12    Same as Fig. 8 except for Exp-IA (solid lines) and Exp-IB (broken lines) and that

986        an orange broken line indicates the RMSE of EnKF using the two-scale Lorenz 96 model.
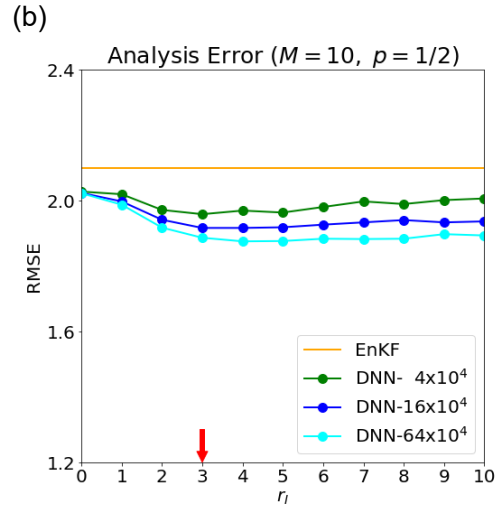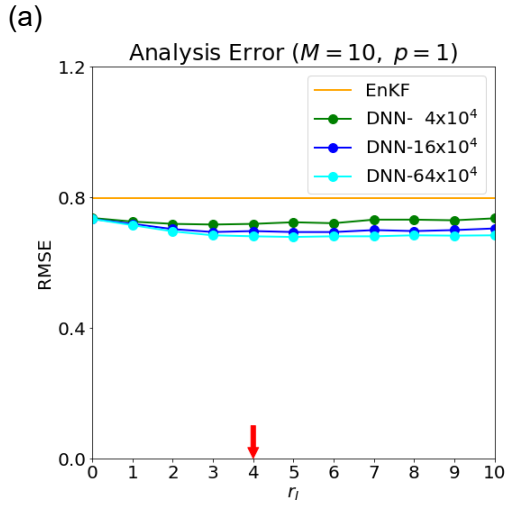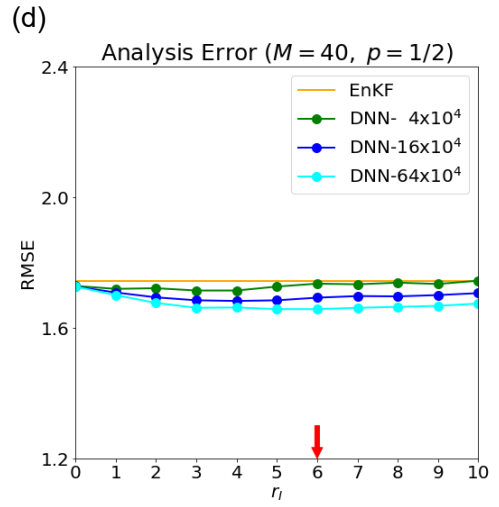
987

Optimal number of nodes per hidden layer

Legend:
- M=10, $p = 1$ (blue)
- M=10, $p = 1/2$ (cyan)
- M=40, $p = 1$ (red)
- M=40, $p = 1/2$ (orange)

y-axis: number of nodes

x-axis: number of samples ($4 \times 10^4$, $16 \times 10^4$, $64 \times 10^4$)
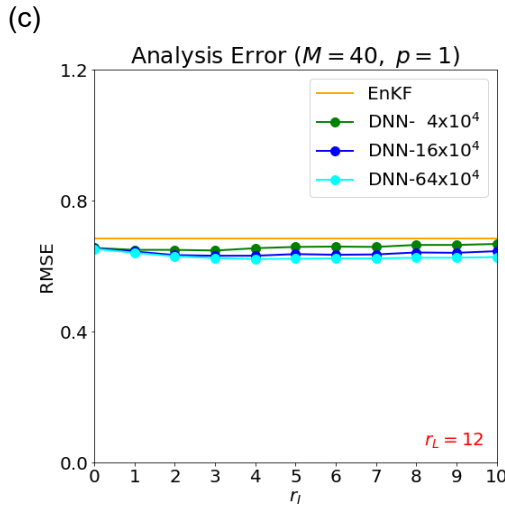
988

Fig. A1    The optimal number of nodes per hidden layer of a DNN with 5 hidden layers for the time interval between observations of 0.50. The abscissa is the number of samples. Blue and cyan bars are for the EnKF ensemble size of 10 for the probability of observations 1 and 1/2, respectively. Red and orange bars are for the EnKF ensemble size of 40 for the probability of observations 1 and 1/2, respectively.

995     (a)

996     (c)

997

998    Fig. A2     Comparison of RMSE between EnKF (orange lines) and the output of a DNN with

999     the number of samples of 40 000 (green lines), 160 000 (blue), and 640 000 (cyan) for

1000     the observation time interval of 0.50. The ensemble size of EnKF and the probability of

1001     observations are (a) 10 and 1, (b) 10 and 1/2, (c) 40 and 1, (d) 40 and 1/2, respectively.

1002     The RMSEs are plotted against the input radius, and a red arrow indicates the optimal

1003     localization radius of EnKF.

1004

# List of Tables

Table 1　Models used in the training and test phases in the experiments.

(a) Perfect model experiments

|  | Exp-PA | | Exp-PB | |
|---|---|---|---|---|
|  | Training | Test | Training | Test |
| Target | L | L | Analysis | L |
| Observations | L | L | L | L |
| Forecast ensemble | L | L | L | L |

L: Lorenz 96 model

(b) Imperfect model experiments

|  | Exp-IA | | Exp-IB | |
|---|---|---|---|---|
|  | Training | Test | Training | Test |
| Target | P | T | T | T |
| Observations | P | T | T | T |
| Forecast ensemble | P | P | T | T |

P: parameterized Lorenz 96 model,　T: two-scale Lorenz 96 model

Table 2    Architecture and training of the feedforward neural network.

| | |
|---|---|
| No. of nodes of input layer | $3(2r_I + 1)$ for $p = 1$ |
| | $4(2r_I + 1)$ for $p = 1/2$ |
| No. of hidden layers | 5 |
| No. of nodes per hidden layer | 5, 10, or 20 (optimized) |
| No. of nodes of output layer | 1 |
| Activation function | ReLU |
| Loss function | Sum of squared error |
| Gradient descent method | Adam* |
| Learning rate | 0.01 to 0.0001 (linear decay) |
| No. of samples | 40 000 |
| Mini-batch size | 100 |
| No. of epochs | 100 |

*: Kingma and Ba (2014)