**Benefits of stochastic weight averaging in developing neural network radiation scheme for numerical weather prediction**

Hwan-Jin Song[1], Soonyoung Roh[1], Juho Lee[2], Giung Nam[2], Eunggu Yun[2], Jongmin Yoon[2], and Park Sa Kim[1]

[1]National Institute of Meteorological Sciences, Korea Meteorological Administration, Jeju-do, Republic of Korea

[2]Graduate School of Artificial Intelligence, Korea Advanced Institute of Science and Technology, Daejeon, Republic of Korea

**Key Points**

- The performance of the neural network radiation scheme was evaluated under a framework of ideal and real cases.

- Stochastic weight averaging is advantageous in generalization compared to the traditional neural network.

- Long-term forecast errors can be largely improved using stochastic weight averaging.

*\* Corresponding author's address*
Hwan-Jin Song
National Institute of Meteorological Sciences,
63568, Seogwipo-si, Jeju-do, Republic of Korea
E-mail: hwanjinsong@gmail.com

34    **Abstract**

35    Stochastic weight averaging (SWA) was applied to improve the radiation emulator based on a

36    sequential neural network (SNN) in a numerical weather prediction model over Korea. While

37    the SWA has advantages in terms of generalization such as the ensemble model, the

38    computational cost is maintained at the same level as that of a single model. The

39    performances of both emulators were evaluated under ideal and real case frameworks.

40    Various sensitivity experiments using different sampling ratios, activation functions, hidden

41    layers, and batch sizes were also conducted. The emulators showed a 60-fold speedup for the

42    radiation processes and 84–87% reduction of the total computation. In the ideal simulation,

43    compared to the infrequent radiation scheme by 60 times, SNN improved forecast errors by

44    5.8–14.1%, and SWA further increased these improvements by 18.2–26.9%. In the real case

45    simulation, SNN showed 8.8% and 4.7% improvements for longwave and shortwave fluxes

46    compared to the infrequent method; however, these improvements deceased significantly

47    after 5 days, resulting in 1.8% larger error for skin temperature. By contrast, SWA showed

48    stable one-week forecast features with 12.6%, 8.0%, and 4.4% improvements in longwave

49    and shortwave fluxes, and skin temperature, respectively. Although the use of two hidden

50    layers showed the best performance in this study, it was thought that the optimal number of

51    hidden layers could differ depending on the given problem. Compared to temperature and

52    precipitation observations, all experiments showed a variability of error within 1%, implying

53    that the operational use of the developed emulators is possible.

54    **Keywords:** neural network, stochastic weight averaging, emulator, speedup, WRF, RRTMG

55

56

**Plain Language Summary**

The NN emulators for radiation parameterization were developed to accelerate the computational speed of the numerical weather forecasting model. Although previous studies have demonstrated that the computational speed for radiation processes can be improved tens of times, guaranteeing stability in long-term forecasting has been recognized as imperative for the operational use of radiation emulator. In general, the multi-model ensemble approach is used to reduce the uncertainty of a single model. However, this approach induces a significant computation burden in proportion to ensemble members. The alternative method developed in this study uses a stochastic averaging technique for weight coefficients during the NN training process, allowing processing to be conducted at the same computational cost as the single model because the dimensions of the final weights are maintained. Application of the trained NN emulator to the numerical model has demonstrated the advantages of generalization in various test cases, while exhibiting significant improvements in accuracy in the latter part of the forecast. This method can therefore contribute to improving emulator studies that face problems related to generalization.

## 1. Introduction

Longwave (LW) and shortwave (SW) radiation physics are important for describing the exchange of energy between the Earth and the Sun. Radiation is a fundamental energy source that determines large-scale atmospheric circulation and consequent physical processes. Accurate calculation involving radiation physics using the line-by-line model (Clough et al., 1992; 2005) requires high computational burden, rendering it important to develop methods that allow rapid calculation of the radiation process. The recent rapid advances in machine learning techniques has led to the development of neural network (NN) emulators for radiation processes in the two main fields: the radiative transfer model (RTM) and radiation parameterization for numerical weather–climate model. An NN emulator that can be used in the RTM was developed some time ago (Chevallier et al., 1998) and was applied to the data assimilation system of the numerical weather prediction (NWP) model (Chevallier et al., 2000). Recent RTM studies based on clear sky simulations have shown a of 1.87–10.88-fold speedup (Liu et al., 2020) when used with the Rapid Radiative Transfer Model for GCMs (RRTMG; Iacono et al., 2008), and 1.8–3.5-fold (Ukkonen et al., 2020) and up to 4-fold (Veerman et al., 2021) for the RRTMGP scheme (Iacono et al., 2008; Pincus et al., 2019). Note that the results of Liu et al. (2020) should be interpreted differently because the measurements described were obtained under different parallelization conditions. Meanwhile, Meyer et al. (2021) showed that using an emulator to add 3D cloud radiative effects was less than 1% more expensive than the 1D scheme; this was a significant decrease in computational cost because the 3D scheme was usually five-times as expensive than the 1D scheme. These results demonstrate the effectiveness of emulating cloud processes in terms of computational cost.

It is difficult to develop an emulator for radiation parameterization within the global circulation model (GCM) and NWP because of complex interactions with various processes

98    within numerical models. However, the emulator for numerical models is more valuable

99    because it can provide important forecasting information that includes factors such as climate

100   change and rapid floods. Thus, the reduction in computational cost associated with the

101   development of an emulator for use with the numerical model would be advantageous in

102   many ways (such as producing national policy or saving lives). Krasnopolsky et al. (2010)

103   used a GCM model of the National Oceanic and Atmospheric Administration (NOAA) with

104   coarse horizontal (~ 100 km) and temporal resolutions, to show that the NN emulator can

105   improve the computational speed of the RRTMG radiation processes by approximately 30

106   times (an average of LW and SW) and reduce 20–25% computational cost for the total model.

107   Notably, the total reduction calculated can vary with the computational percentage used for

108   the radiation scheme to that used for the total model. The deep neural network (DNN)

109   emulator that was developed by Pal et al. (2019) showed 8–10 times speedup for radiation

110   parameterization; however, the total reduction achieved in terms of computational cost was

111   not elucidated. Song and Roh (2021), and Song et al. (2021) performed NWP studies with 5-

112   km spatial and 20-s temporal resolution from the Korea Meteorological Administration

113   (KMA) to show a 60-fold speedup in RRTMG-K radiation parameterization (Iacono et al.,

114   2008; Beak, 2017), with an 87% reduction in the time taken for total model computation. The

115   significant difference in the total computation reduction achieved in GCM and NWP studies

116   is because GCMs typically use an hourly scale radiation time step (radt), whereas the NWP

117   studies used the same time step for both the total model and the radiation process (i.e., 20 s),

118   leading to a more accurate result but a higher computational burden for the control run (i.e.,

119   more speedup for the emulator).

120   All these studies of radiation emulators have mainly been developed using the NN or

121   DNN techniques because these methods can be simply implemented into Fortran in both the

122   GCM and NWP. However, recent developments have been made in machine learning

123   techniques based on the Python code. Ott et al. (2020) recently developed the Fortran-Keras

124   Bridge to communicate between Fortran and Python, and it is actively used in emulator

125   studies. However, such efforts remain within the scope of the DNN, and other deep learning

126   techniques have not yet been attempted. Although Liu et al. (2020) applied a convolutional

127   neural network (CNN) to a single column model, it was based on the use of a Python wrapper

128   outside the numerical model. For real-case modeling such as the GCM or NWP, which are

129   based on large-scale Fortran codes, this approach is difficult to apply. Most NN emulators for

130   radiation parameterization in the GCM and NWP have been developed by the NOAA

131   (Krasnopolsky et al., 2005, 2008, 2010; Belochitski et al., 2011; Belochitski and

132   Krasnopolsky, 2021) or the KMA (Roh and Song, 2020; Song and Roh, 2021; Song et al.,

133   2021) using Fortran software (Krasnopolsky, 2014). However, this software does not support

134   other activation functions other than tangent hyperbolic (Tanh), DNN with multiple hidden

135   layers, and batch (or parallel) learning. Although functions other than Tanh (e.g., sigmoid,

136   softsign, arctan, and rectified linear unit (ReLU)-type functions) have been used in many

137   studies (Pal et al., 2019; Liu et al., 2020; Roh and Song, 2020; Ukkonen et al., 2020;

138   Veerman et al., 2020), the best activation function that is used for the radiation emulator is

139   still controversial. The development of DNN emulators has included several sensitivity

140   experiments investigating the number of neurons and hidden layers (Pal et al., 2019; Liu et al.,

141   2020; Veerman et al., 2020; Meyer et al. 2021); however, no attempt has yet been made to

142   investigate the radiation process at the same computational cost (or speedup the process). Pal

143   et al. (2019) compared the validation loss architecture of 32-32-32 (32 neurons and 3 hidden

144   layers) with 16-16-16 (16 neurons and 3 hidden layers), 32-32-32-32 (32 neurons and 4

145   hidden layers), and 64-64-64 (64 neurons and 3 hidden layers), but the computation costs of

146   the experiments differed because the numerical complexity is expressed as the total

147   dimension of the weight and bias coefficients. Furthermore, the use of a single hidden layer,

148    which can include the largest number of neurons at the same computational cost, was not

149    considered in Pal et al (2019). Belochitski and Krasnopolsky (2021) emphasized the risks of

150    using the DNN emulator in relation to increasing nonlinearity, and retained the use of a single

151    hidden layer in developing the NN emulator for radiation parameterization. However, no

152    practical evidence was provided (i.e., the DNN experiments were not performed), indicating

153    that the accuracy of NN (with a single hidden layer) and DNN (with multiple hidden layers)

154    emulators still requires comprehensive evaluation at the same computational cost and

155    numerical complexity. Sensitivity tests with different batch sizes have rarely been performed

156    in the field of radiation emulation, except for the speedup check that was reported in Liu et al.

157    (2020). In general, the use of an appropriate mini-batch is known to produce a more accurate

158    solution than the full batch (Li et al., 2014), while requiring more training (a small batch size

159    is equivalent to less parallelization). Thus, further consideration of batch size may contribute

160    to optimizing the performance of the radiation emulator.

161        Stochastic weight averaging (SWA), which was recently developed in the field of

162    machine learning, is aimed at increasing generalization in the NN training process (Izmailov

163    et al., 2018). In general, a multi-model ensemble approach is used to reduce the uncertainty in

164    a single model. However, this approach is not appropriate for use in emulators that are used

165    to speed up the GCM and NWP because the computational burden is directly proportional to

166    the number of ensemble members included. As an alternative approach in which the

167    computational cost can be minimized, SWA performs the averages for multiple points along

168    the trajectory of the stochastic gradient descent (SGD) (Bottou, 2012; Mandt et al., 2017)

169    under constant or cyclical learning rates. SWA tends to find a wide flat solution using this

170    method, whereas the SGD often converges to a sharp (or local) minimum that can cause

171    problems with generalization. Izmailov et al. (2018) noted that the use of SWA can improve

172    the accuracy of test sets with better generalization than conventional SGD in terms of several

173　benchmarks. To the best of our knowledge, SWA has never been used in climate and weather

174　models. In fact, as noted by Krasnopolsky et al. (2008), Belochitski and Krasnopolsky (2021),

175　and Song et al. (2021), emulator studies using the GCM and NWP face severe problems with

176　generalization because the errors that are accumulated during long-term integration by the

177　emulator can induce a blow-up of the entire numerical model. Because infinite training

178　datasets cannot be used, generalization is an important issue for developing universal

179　radiation emulator.

180　　This study therefore mainly examines the benefits of using SWA in developing a

181　radiation emulator for the NWP model under the frameworks of idealized squalline and real

182　case simulations. The ideal simulation will then serve as a testbed for various sensitivity

183　experiments. At the same computational cost, the results of SWA will be compared with NN

184　based on sequential training (SNN), which has been used in many previous studies

185　(Krasnopolsky et al., 2005, 2008, 2010; Belochitski et al., 2011; Roh and Song, 2020;

186　Belochitski and Krasnopolsky, 2021; Song and Roh, 2021; Song et al., 2021), and the

187　infrequent use of radiation scheme, which is a popular method in operational NWP fields

188　(Pauluis and Emanuel, 2004; Pincus et al., 2013). Sensitivity experiments investigating the

189　sampling ratio of training sets, activation functions, the number of hidden layers (at the same

190　speedup), and batch sizes (as well as learning rates) are also conducted. These all efforts will

191　contribute to reducing the forecast error of the NWP model using the NN radiation scheme

192　that can attain significant speedup.

193　**2. Data and Methods**

194　　This study considers two types of frameworks (i.e., ideal and real cases) to evaluate the

195　performance of a radiation emulator based on the Advanced Research Weather Research and

196　Forecasting (WRF-ARW) model (Skamarock et al., 2019). The ideal framework was based

197　on a two-dimensional squall-line simulation with 5-km resolution on 201 horizontal grids,

198  using 39 vertical layers up to 50 hPa and a 24-h integration period with a model time step (dt)

199  and radiation time step (radt) of 20 s serving as the control run for the ideal simulation.

200  Different horizontal resolution (0.25 km → 5 km), integration time (6 h → 24 h), and time

201  steps (3 s → 20 s) than those used in Roh and Song (2020) allowed consistency with the real

202  case experiment. Thus, this experiment can provide conceptual guidance for large-scale

203  datasets generated under real conditions. The use of small-scale data rendered it possible to

204  perform various sensitivity experiments. For the real case, this study used the horizontal

205  domain with 234×282 grids over the Korean peninsula, which is the same that utilized in the

206  Korea Local Analysis and Prediction System (KLAPS), one of the operational NWP models

207  used by the KMA. Note that the dynamics and physics processes of the KLAPS were based

208  on the WRF model. The radiation emulator used in both ideal and real case frameworks

209  targets the RRTMG-K radiation scheme (Baek, 2017), which calculates vertical heating rates

210  and LW fluxes with 256-g points in 16 bands and SW fluxes with 224-g points in 14 bands.

211  The WRF double moment 7-Class (WDM7) microphysics scheme (Bae et al., 2019) was used

212  in both simulations. The real case simulation further used the KIAPS Simplified Arakawa–

213  Schubert (SAS) cumulus (Kwon and Hong, 2017), the Shin and Hong planetary boundary

214  layer (Shin and Hong, 2015), the revised MM5 Monin–Obukhov surface layer (Jiménez et al.,

215  2012), and the Unified Noah land surface model (Tewari et al., 2004). The RRTMG-K

216  scheme accounted for 85.0% (for the ideal case) and 88.6% (for the real case) of the total

217  computational costs of using the WRF model under the same dt and radt (20 s). The ideal and

218  real case frameworks were initialized by default initial sounding in the WRF model (with

219  warm bubble forcing at low levels) and data from the European Center for Medium-Range

220  Weather Forecasts Reanalysis v5 (ERA5) (Hersbach et al., 2020) with 0.25° grid and 3-h

221  intervals, respectively.

222     The training sets for the ideal simulation were prepared through random sampling of the

223     full set (i.e., control run for 24 h) using sampling ratios from 10% to 90%. The representation

224     error was reduced under an increase in the sampling ratio. However, the ideal experiment is a

225     highly nonlinear system that is sensitive to small perturbations in the initial stage; therefore,

226     the radiation emulator was found to produce quite different results during the 24-h integration

227     over 201 grids (i.e., the emulator was applied 868,320 times) when it was applied to the WRF

228     model (i.e., via online prognostic testing). Thus, we did not expect a consistent trend with the

229     sampling ratio. The training sets were divided into LW clear, LW cloud, SW clear, and SW

230     cloud to maintain consistency with the input–output structure of the radiation emulator

231     developed by Song and Roh (2021). The training sets for the real case simulations were sub-

232     sampled from 10-min interval outputs from the period 2009–2019, with 48 days from the

233     period of 2009–2018 and the one-year period of 2019 used in Song and Roh (2021) evenly

234     considered (i.e., 50% of the 48 days and 50% in 2019). Note that the 48 days included events

235     on which the maximum and the second maximum precipitation occurred in each month

236     together with non-precipitating 24 days over the period of 2009–2018. To optimize the

237     hyperparameters used in the NN training, we further prepared independent validation sets

238     consisting of the days on which the third and fourth maximum precipitation occurred in each

239     month over the period of 2009–2018 along with other non-precipitating 24 days which were

240     not used in the training sets. Note that the validation sets were newly adopted in this study

241     because Song and Roh (2021) did not optimize the hyperparameters. The training and

242     validation sets were divided into 96 categories with 3 million cases in each, as in Song and

243     Roh (2021), who used a 96-categories approach (LW and SW, clear and cloud, land and

244     ocean, and 12 months) to effectively utilize as much data as possible to reduce the

245     representation error. The final evaluation of accuracy was performed for the year 2020 using

246     a one-week period and 3-h intervals (test sets), while the emulator was implemented in the

247    WRF model (i.e., online prognostic testing). Note that the one-week forecast period used in

248    this study was much extended compared to the one-day period used in Song and Roh (2021).

249        The inputs for the NN emulator for the ideal simulation consist of 187 variables,

250    including: pressure (39 profiles), temperature (39 profiles), water vapor (39 profiles), ozone

251    (39 profiles), and cloud fraction (30 profiles due to the removal of constant values), in

252    addition to skin temperature (LW) and the solar constant multiplied by the cosine zenith

253    angle (SW). The inputs were decreased by 157 variables in the clear case, because the cloud

254    fraction was not used. The inputs for the real case simulation further included surface

255    emissivity (LW), surface albedo (SW), and monthly variant cloud fraction (28 to 35 profiles).

256    Unlike Song and Roh (2021), topography (longitude, latitude, and elevation) was excluded in

257    this study. The outputs for both the ideal and real case simulations consist of 39 heating rate

258    profiles and three fluxes (upward fluxes at the top and bottom, and downward flux at the

259    bottom). Hereafter, the heating rate and flux in this study refer to the heating rates in the 39

260    layers and the three fluxes, respectively. The inputs and outputs are summarized in Table 1.

261    For given input–output pairs, two NN methods were applied: SNN (Krasnopolsky, 2014) and

262    SWA (Izmailov et al., 2018). Both are fully connected and feed-forward NN methods. Here,

263    the same min-max normalization and standardization were used for the inputs and outputs,

264    respectively. In addition, because the SNN provides the utility of early stopping, the

265    maximum number of epochs used in SWA was determined from the SNN. The SWA mode

266    was applied to the last 25% of the epochs, as in Izmailov et al. (2018), while the former 75%

267    of the epochs was trained by the common SGD. Under the ideal simulation, the mean and

268    standard deviation of epochs were 13,499±4697 for clear and 4,089±832 for cloud cases

269    with different sampling ratios of 10–90%. When the number of samples is large, the required

270    epoch tends to decrease. For the real case, the mean of 3,011 epochs was used for clear and

271  2,251 for cloud conditions; thus, approximately 3,000 and 2,200 epochs were used,

272  respectively.

273      After the NN training, the weight and bias coefficients were obtained and inserted into the

274  radiation emulator, replacing the RRTMG-K code (module_ra_rrtmg_swk.F) in the WRF

275  model. In the emulator code, the NN outputs were forced into the range between the

276  minimum and maximum values of the training sets to prevent extrapolation. Because the

277  numerical complexity in the NN is defined as the total sum of the dimensions of the weight

278  and bias coefficients, the use of 90 neurons in a single hidden layer for the radiation process

279  corresponds to a 60-fold speedup and an 87% reduction in the total computation time (Song

280  and Roh, 2021). We follow this methodology for the real case simulation. For the ideal case,

281  the mean computation time for the radiation process and the total model were measured using

282  the Intel Xeon E5-2690v3 central processing unit (CPU) with serial compilation condition.

283  As a result of averaging 10 experiments, a mean speedup of 60 times (3086 s ÷ 51.5 s) was

284  achieved for the radiation processes and the time taken to run the total model was 84% (3630

285  s vs. 593.5 s) lower. The small difference observed between the results obtained using the

286  SNN and SWA was thought to be due to different cloud conditions during integration. For the

287  situation in which there are the same number of neurons in the hidden layers, the numerical

288  complexity of the NN or DNN can be expressed as: $I \times N + N + (H–1) \times (N \times N + N) + N \times O + O$.

289  Here, I is the number of input variables, O is the number of output variables, N is the

290  number of neurons, and H is the number of hidden layers. For example, 68-68 (two hidden

291  layers), 58-58-58 (three hidden layers), 52-52-52-52 (four hidden layers), and 47-47-47-47-

292  47 (five hidden layers) neuron structures are comparable to 90 neurons with a single hidden

293  layer in terms of producing a 60-fold speedup. This is a fair approach in terms of

294  computational cost, unlike the sensitivity experiments in Pal et al. (2019), Liu et al. (2020),

295  Ukkonen et al. (2020), and Veerman et al. (2021). These comparisons can be used to obtain

296   an answer to the controversial argument raised by Belochitski and Krasnopolsky (2021), who

297   discussed the use of a single hidden layer (with a long history) and multiple hidden layers in

298   developing an NN emulator for radiation parameterization.

299       In conclusion, the idealized squalline simulation was used to perform the sensitivity

300   experiments using sampling ratios from 10–90% in generating a training set. Both SNN and

301   SWA methods were applied, and their accuracy was measured in terms of the root mean

302   square error (RMSE) by comparing with the control run over 24 h. As in a previous study

303   (Song and Roh, 2021), the 60-fold speedup (i.e., 90 neurons) emulator results were also

304   compared with the infrequent radiation scheme with a radt of 20 m (denoted as "WRF60" in

305   this study). Here, we did not adjust the time between the infrequent calls, as in Manners et al.

306   (2009) and Hogan and Bozzo (2015), because the treatment was not available in the

307   WRF model. To minimize the redundancy problem, a sampling ratio of 10% was selected

308   and then applied to subsequent experiments. For the second experiment, sensitivity tests were

309   conducted with 16 nonlinear activation functions (Tanh, Arctan, Tanhshrink, Sigmoid,

310   Logsigmoid, SiLU, Softsign, Softplus, Mish, Hardtanh, Hardsigmoid, Hardswish, ReLU,

311   LeakyReLU, ELU, and SELU) based on SWA. Detailed definitions of the activation

312   functions are presented in Table 2. The SWA results were evaluated with SNN using Tanh.

313   The third experiment involved sensitivity tests on the number of hidden layers (1–5). The

314   numerical complexity, and thereby speedup, for the radiation process was maintained by

315   reducing the number of neurons in a given hidden layer. Different speedup conditions of 15,

316   30, 45, 60, 90, and 120 times were considered in the ideal simulation. The best performance

317   for each speedup condition was selected from the mean RMSEs using five prediction

318   variables (LW/SW heating rates, LW/SW fluxes, and surface temperature) over 24 h. In real

319   case simulation, experiments for batch sizes and learning rates were performed for validation

320   sets. The experiment based on huge datasets (96×3 million data) was found to be extremely

321    time consuming compared to the ideal case. In fact, the SNN based on sequential training

322    with one batch size (Krasnopolsky, 2014) is fundamentally different from the batch learning

323    in SWA (or SGD). In addition, SNN was performed using adjustable learning rates ($10^{-3}$ to

324    $10^{-6}$) during the NN training and generally converged at optimal solutions of approximately

325    2,000 and 1,200 epochs with a learning rate of $10^{-4}$. Smith et al. (2018) insisted that batch

326    size and learning rate should be proportional to achieve similarly high performance among

327    the experiments. The empirical relationship observed between batch size and learning rate

328    under the SNN (1 and $10^{-4}$) was thus applied to the experiments investigating batch sizes

329    (100–9000) and initial learning rates (0.001–0.9) in the SWA. It should be noted that the

330    learning rate of the SWA mode was reduced by half of its initial value under cosine annealing.

331    The computation time taken for training all datasets (i.e., 96 sets) was 12 h using the NVIDIA

332    DGX A100 graphics processing unit (GPU) 16 units, in contrast to the 63 h taken by the SNN

333    using 96-node parallelization that was carried out with the Intel Xeon E5-2690v3 CPU. The

334    computation time taken by the GPU for training was based on a batch size of 500 (which will

335    be further discussed later). The learning rate of the SWA in the ideal simulation was

336    determined empirically by multiplying the full batch size (equal to the number of datasets) by

337    $2\times10^{-6}$ based on a learning rate of 0.92997, which is less than 1 for the maximum number of

338    datasets (464,985). Note that there were 316,322 LW clear, 464,985 LW cloud, 115,103 SW

339    clear, and 215,821 SW cloud datasets for the sampling ratio of 90%, and the numbers were

340    reduced proportionally to the sampling ratio. No further experiments were performed on

341    batch size or learning rate in the ideal simulation, although the use of mini-batches and a

342    proper learning rate may lead to better optimization. The SWA group with the highest

343    accuracy in the validation sets (2009–2018) was used in the final online testing for the year

344    2020. The RMSE evolutions during a one-week period were examined for LW/SW fluxes,

345    skin temperature, 2-m air temperature, and 3-h accumulated precipitation. The evaluation of

346    2-m temperature and precipitation was performed by comparing with surface observation in

347    South Korea, and the other variables were compared with the control run and WRF60. The

348    real case experiments on multiple hidden layers (2–5) were further examined in the final

349    evaluation step.

350    **3. Results and Discussion**

351    For the idealized squall line simulation, nine-type datasets with a sampling ratio ranging

352    from 10% to 90% were trained by the SNN and SWA methods. The two methods were based

353    on Tanh. The mean RMSEs for five variables are compared with the results of the control run,

354    which was executed over 24 h in 1-min intervals over the 1000-km domain in Fig. 1 (LW/SW

355    heating rates, LW/SW fluxes, and surface temperature). The emulator results were used 4,320

356    times temporally (number of time steps) and 201 times spatially (number of grids). Only

357    daytime variables were considered in the RMSE calculation of SW radiation. No apparent

358    dependency on the sampling ratio was observed in either SNN or SWA. Although the

359    representation error should decrease when the sampling ratio is increased, the strong

360    nonlinearity of the ideal simulation appears to have significantly influenced the results over

361    24 h. We can also suspect a strong correlation between training sets because 5-km and 20-s

362    interval data were used. In such a situation, finding an optimal sampling ratio for NN training

363    using advanced sampling techniques can be helpful and should be investigated in the future.

364    Compared to the SNN, improvements of 9.9% were observed in the mean RMSE for all

365    sampling ratios by using SWA, indicating that SWA can guarantee a better performance than

366    SNN, regardless of the datasets used. Because the NN approximation tends to be optimized to

367    reduce the total error, the improvements are not linear for all variables. On average, the SW

368    heating rate showed the largest improvement (20.7%) of the five variables, and can increase

369    the predictability during the daytime. Roh and Song (2020) also noted that the SW heating

370    rate is the most uncertain variable among radiation products. The uncertainty of the SW

371     heating rate is thought to be significantly reduced by using SWA. For a sampling ratio of

372     10%, the mean RMSE improvements generated by using SWA for the five variables were

373     13.2% higher than errors involved in using SNN (23.20% vs. 10.03%). The improvements in

374     the RMSE obtained by using SWA were relatively large for the SW outputs (12.2–20.7%).

375     The difference between SNN and SWA was large for small sampling ratios (10% and 30%,

376     respectively), which is thought to be because SWA can better generalize the training results

377     compared to common NN (Izmailov et al., 2018). Because all of the data covering natural

378     variability can be obtained, this benefit of using SWA is expected to exert a strong influence

379     and improve the performance in the real-case simulation.

380       These results suggest that datasets based on a 10% sampling ratio with the smallest

381     redundancy should be used. The activation function is an important hyperparameter that can

382     significantly affect the performance of emulator because it is used not only in the learning

383     process but also in the emulator code (within the WRF model). The SWA results using 16

384     activation functions (Tanh, Arctan, Tanhshrink, Sigmoid, Logsigmoid, SiLU, Softsign,

385     Softplus, Mish, Hardtanh, Hardsigmoid, Hardswish, ReLU, LeakyReLU, ELU, and SELU)

386     are compared with the results obtained by SNN based on Tanh in Fig. 2, together with the

387     RMSEs for 24 h over the 1000-km domain. The mean and standard deviation of RMSEs

388     varied by $2.21\pm0.12$ K day$^{-1}$ for LW heating rate, $0.98\pm0.06$ K day$^{-1}$ for SW heating rate,

389     $12.19\pm1.63$ W m$^{-2}$ for LW flux, $118.93\pm19.58$ W m$^{-2}$ for SW flux, and $0.86\pm0.10$ K for

390     surface temperature. Some activation functions (e.g., Arctan and Hardswish) showed worse

391     performance than SNN. The lowest error among the SWA experiments was observed when

392     Tanh was used. This feature is in line with many emulator studies based on Tanh

393     (Krasnopolsky et al., 2005, 2008, 2010; Belochitski et al., 2011; Roh and Song, 2020;

394     Belochitski and Krasnopolsky, 2021; Chantry et al., 2021; Song and Roh, 2021; Song et al.,

395     2021), and we therefore used Tanh for subsequent experiments.

396    Figure 3 shows the temporal and horizontal evolution for the LW/SW upward fluxes at

397    the top (LWUPT/SWUPT), surface temperature, and precipitation rate at 1-min intervals. The

398    control run, SNN, and SWA results (radt = 20 s) were compared with those of WRF60 (radt =

399    20 m). The SNN, SWA, and WRF60 have the same computational cost with an 84%

400    reduction compared to the control run. The control run shows evolutionary features in two

401    directions (i.e., positive and negative X directions) that are initialized at the center position (0

402    km). The highest SWUPT (an indicator of deep clouds) and the lowest surface temperature

403    areas were observed along the positive X direction. These areas are associated with a

404    squalline precipitating system. This squalline feature was not evident in Roh and Song (2020),

405    probably because of a strong interaction between radiation and microphysics in the small

406    domain (50 km), although this experiment showed the squalline feature in the microphysics

407    scheme only. In the negative X direction, low LWUPT and high SWUPT (an indicator of

408    clouds) and low surface temperature areas are characterized by non-precipitating clouds (e.g.,

409    anvils). The forecast error is more evident in the cloud areas. Interestingly, WRF60 showed

410    discontinuous features for LWUPT and SWUPT, which are direct outputs from the radiation

411    scheme, because the radiation scheme was used 60 times (radt = 20 m) less than the dt of 20 s.

412    This problem was not found in the results of SNN and SWA because radt of 20 s was used, as

413    in the control run. Overall, evolutionary features of the squalline system appear to have been

414    properly simulated in both SNN and SWA.

415        The time series of the RMSEs for the five variables are shown in Fig. 4. The simulation

416    was initialized at midnight and then integrated for 24 h. The zero SW heating rate and flux

417    (i.e., nighttime) were excluded from the analysis. In WRF60, the RMSEs for the LW heating

418    rate and flux tended to increase substantially with integration time because the error due to

419    the use of the infrequent radiation scheme accumulated during integration. The RMSEs of

420    SW heating rate and flux were largest around noon in association with the strong incident SW

421  radiation. The RMSEs of LW heating and flux decreased substantially after sunset when the

422  effects of the SW radiation disappeared. The SNN results show an improved RMSE pattern

423  as a whole compared to WRF60, with improvements evident for all variables before noon.

424  However, the RMSE improvements tended to weaken in the afternoon. This clearly reveals

425  the fundamental problem of radiation emulator, which is associated with accumulated errors

426  during integration (Krasnopolsky et al., 2008; Song et al., 2021). Using SWA alleviated the

427  problem that appeared when using SNN. Before 4 h, SWA showed a larger error than SNN

428  for the LW heating rate, flux, and surface temperature. However, after 4 h, SWA produced

429  significantly lower RMSEs for all variables. The RMSE improvements associated with SWA

430  were evident in relation to the SW radiation during daytime. The largest improvement among

431  the five variables was observed in the SW heating rate, as seen in Fig. 1. Around sunset and

432  afterwards, the RMSE improvements gained by using SWA tended to decrease, indicating

433  that the results are affected by the daily solar cycle; this assumption can be confirmed using

434  the results obtained over multiple days in the subsequent real case simulations (i.e., one

435  week). The total statistics of the ideal simulations are summarized in Table 3. In terms of the

436  total improvement for the five variables compared with WRF60, the performance of the SNN

437  with 60-fold speedup was located between WRF9 with 9-fold speedup (radt = 3 m) and

438  WRF30 with 30-fold speedup (radt = 10 m). In contrast, the SWA results were even better

439  than those of WRF9. Note that WRF9 performed the best among the infrequent uses of

440  radiation scheme with radts of 1 m to 5 m. These results suggest that SWA can produce more

441  accurate and fast results compared with the operational method based on infrequent radiation

442  scheme.

443      Before examining the real case simulation, we further examined the effect of multiple

444  hidden layers (i.e., DNN) on the SWA emulator under the idealized squalline framework.

445  Here, we focus on six speedup conditions of 15, 30, 45, 60, 90, and 120 times for the

446 radiation process, which correspond to 360, 180, 120, 90, 60, and 45 neurons in a single

447 hidden layer. For each speedup condition, we considered DNN structures with two to five

448 hidden layers that have the same numerical complexity as a single hidden layer. For example,

449 in relation to 60-fold speedup, 90, 68-68, 58-58-58, 52-52-52-52, and 47-47-47-47-47

450 neurons were used for one, two, three, four, and five hidden layers, respectively. Figure 5

451 shows that the use of a single hidden layer produced the lowest error among all experiments

452 under the same speedup conditions. Note that dark gray colors predominated in the single

453 hidden layer (Fig. 5) and the use of multiple hidden layers showed 7.41–9.80% degradation

454 compared to the single hidden layer on an average of six speedup cases in terms of the mean

455 RMSE improvement for five variables compared with WRF60. This is thought to be related

456 to the reduction in the number of neurons used for the DNN and provides experimental

457 evidence for the conceptual argument by Belochitski and Krasnopolsky (2021) that the

458 nonlinearity of the DNN can be rapidly increased owing to the complex structure of hidden

459 layers, which can lead to more unstable generalization such as nonlinear extrapolation.

460 Vapnik (2019) also noted that the use of DNN does not always guarantee the best solution for

461 a given problem. However, this result was based on one ideal case from which we cannot

462 draw general conclusion regarding the usefulness of the DNN in developing radiation

463 emulator.

464   As described in the Data and Methods section, the real case simulation was primarily

465 based on KLAPS, which is one of the operational NWP models in the KMA. The training

466 sets were based on the period between 2009 and 2019. The 48 days that were not used for

467 training data were used as the validation sets to optimize the hyperparameters in the SWA.

468 This can be considered as offline testing, whereas the final evaluation for the year 2020

469 connected with WRF modeling was tested online. Unlike the online prognostic test, which is

470 affected by the integration of the numerical model, the accuracy of the offline test should be

471    relatively high because the error does not accumulate. In the offline test, we mainly examined

472    the optimization of the batch size and learning rate in the SWA method. The batch size is an

473    important hyperparameter in determining the fundamental difference between SNN, which is

474    based on sequential training (batch size = 1), and SWA, which is based on batch training

475    (batch size > 1). Reducing the batch size (i.e., the use of mini-batches) and learning rate can

476    lead to better performance in general; however, Smith et al. (2018) insisted that batch size

477    and learning should be proportional to each other. Here, we empirically forced a proportional

478    relationship of $10^{-4}$ between batch size and learning rate based on the relationship observed in

479    the SNN (1 and $10^{-4}$). Because the use of too small batch size (i.e., less parallelization) led to

480    a substantial increase in the training speed, we empirically set the minimum batch size as 100.

481    The batch size was extended to 1000 with 100 intervals and 9000 with 1000 intervals. The

482    corresponding learning rates were 0.001 to 0.9. Figure 6 shows the validation results for the

483    LW/SW heating rates and LW/SW fluxes. Here, 12 months, land/ocean, and clear/cloud

484    results were averaged. The fraction of land over the entire domain was 45.3% and the mean

485    fraction of cloud was assumed to 50%. Regardless of the batch sizes and learning rates used,

486    SWA exhibited superior performance compared to SNN. On average of 10 experiments, the

487    RMSEs of the LW/SW heating rates and LW/SW fluxes were improved by 3.15%, 8.68%,

488    7.92%, and 9.70%, respectively, compared with the RMSEs obtained using SNN (0.4740 K

489    day$^{-1}$, 0.1968 K day$^{-1}$, 3.9140 W m$^{-2}$, and 21.6417 W m$^{-2}$, respectively). Among the 10

490    experiments, the result obtained with a batch size of 500 and a learning rate of 0.05 showed

491    the best performance with RMSE improvements by 3.21%, 10.21%, 8.18%, and 11.59% for

492    the LW/SW heating rates and LW/SW fluxes, respectively. Similar to the ideal simulation,

493    there were relatively large improvements in the SW outputs. These results reveal the

494    characteristics by which SWA strengthens generalization at the expense of training accuracy

495    (Izmailov et al, 2018). The obtained settings (500 and 0.05) were thus used to evaluate the

496    final performance of the online testing results in the real-case simulation.

497        Figure 7 represents the spatial distribution of LWUPT, SWUPT, and skin temperature for

498    a real-case example (typhoon HAISEN, 12LST September 17, 2020). The typhoon is the

499    most extreme weather phenomena that occur over the Korean peninsula. Since it was

500    initialized on 00LST September 1, this case corresponds to a 6.5-day forecast result; thus, the

501    radiation scheme used 28,080 time-steps (with a radt of 20 s). Note that this is a more long-

502    term result compared with the 12-h forecast result for typhoon SANBA in Song and Roh

503    (2021). Despite the 156-h forecast, the SNN and SWA emulator results show similar patterns

504    to the WRF control run, with differences in the detailed patterns. The LWUPT and SWUPT

505    around the typhoon were characterized by low and high values, respectively; mainly over the

506    northern part of the Korean Peninsula. These areas were also connected to cold surface

507    temperatures. During the event, the RMSEs for LWUPT and SWUPT in the SNN (SWA)

508    were improved by 11.11% (10.89%) and 6.08% (6.84%), respectively, compared to WRF60

509    (13.68 W m$^{-2}$ and 138.92 W m$^{-2}$). However, SNN exhibited a 15% higher RMSE for skin

510    temperature. This feature was significantly improved by using SWA, with a 1% decrease in

511    RMSE compared to WRF60, implying that SWA produces more stable results.

512        More generalized evaluations of the total cases are shown in Fig. 8, in which 48 real-case

513    simulations are presented. Each simulation was initialized on the 1st, 8th, 15th, and 22nd of

514    each month in 2020 and then integrated for one week. Thus, 29–31 days in each month were

515    excluded from the analysis. Each RMSE at a given 5-km grid in Fig. 8 represents a statistical

516    result for a one-week forecast over 48 cases in 2020. As shown in Fig. 7, both SNN and SWA

517    tended to improve the forecast accuracy of LW/SW fluxes compared with WRF60, and SWA

518    showed further reduced RMSEs for LW flux, SW flux, and skin temperature than SNN.

519    Relatively large errors of LW flux and skin temperature remain in the mountainous area of

520    North Korea. A more quantitative analysis is presented in Fig. 9. The RMSE time series

521    denotes a statistical result over 226×274 grids (excluding ±4 boundary points) and 48 weeks

522    at 3-h intervals (totaling 166 million data points). In Fig. 9a, the RMSE for the LW flux under

523    WRF60 tended to increase rapidly before 2 day, and then steadily fluctuated with diurnal

524    perturbation observed after 2 day. The improvements in the RMSE of the LW flux for SNN

525    (compared to the WRF60) decreased substantially from 15.5% before 1 day to only 1.4%

526    after 6 days (Fig. 9a). This represents a weakness in the radiation emulator that the

527    accumulation of errors caused by the NN approximation can be rapidly amplified in long-

528    term forecasts. However, because the SWA method is effective in reducing the uncertainty,

529    the RMSE improvements seen in the LW flux were 19.7% before 1 day and 9.0% after 6 day

530    (Fig. 9a). In particular, the RMSE of the LW flux after 6 day was 7.8% lower using SWA

531    than that obtained using SNN. For the SW flux (Fig. 9b), the time series of the RMSEs were

532    relatively similar to those for the LW flux. Looking at the maximum RMSEs of SW flux

533    around noon, the SNN and SWA emulators showed smaller RMSEs until 5 day, whereas the

534    SNN results produced the largest error after 5 day. Thus, we can assume that the rapid

535    increase in the RMSE of the LW flux is also affected by SW radiation. Note that the mean

536    RMSE of SW flux for the SNN decreased by 8.8% after 5 day, whereas that of the SWA

537    improved by 6.3% compared to WRF60. For skin temperature, both emulator results showed

538    degradation after 4 day (Fig. 9c). The maximum RMSEs of skin temperature during both

539    daytime and nighttime were larger than those of WRF60, whereas SWA was better than SNN.

540    Skin temperature is not a direct output of the radiation scheme, and it can interact with other

541    processes in a complex manner. In determining skin temperature, it is thought that the

542    influence of clouds (e.g., the amount and location of clouds) will be greater than that of the

543    radiation process. This can lead to an interpretation of Fig. 9d, which shows the evaluation

544    results with 2-m temperature observations in South Korea. In Fig. 9d, while the RMSEs were

545 distributed over 1.9–2.7 K, the difference obtained from the various experiments was

546 relatively small. The final RMSEs are listed in Table 4. The RMSEs were 2.2438 K for

547 WRF60, 2.2466 K for SNN, and 2.2563 K for SWA, and their difference was much smaller

548 than the observation error (0.1 K). Similar results were also found in the evaluation of

549 precipitation compared with the gauge-radar merged observations in South Korea (Fig. 10),

550 with RMSEs of 12.1987–12.3120 mm (Table 4). The standard deviation of the RMSEs was

551 only 0.4% of the mean RMSE obtained for precipitation. As noted by Song and Roh (2021),

552 because the control run also had errors as compared with observation, the error induced by

553 the use of a radiation emulator can be insignificant in terms of observation. Instead, the

554 uncertainty associated with clouds can play a more important role in determining surface

555 temperature. Even so, these results imply that the radiation emulators in this study produce

556 accurate one-week forecasts at the NWP level, in addition to a significant 60-fold speedup. In

557 this context, the use of SWA guarantees robust results in terms of speed, accuracy, and

558 stability. The RMSEs for both emulators were between those of WRF30 and WRF60 (Table

559 4).

560　　　When multiple hidden layers and a small number of neurons (i.e., keeping the same 60-

561 fold speedup) were considered, the RMSEs for the one-week forecast changed (Table 4).

562 Here, 90, 68, 58, 52, and 48 neurons were used in 1–5 hidden layers (1 h to 5 h), respectively.

563 Among the five SWA experiments using the different numbers of hidden layers, the use of

564 two hidden layers showed the lowest RMSEs for LW/SW fluxes and skin temperature,

565 exhibiting 0.4–1.3% lower RMSEs compared with the use of one hidden layer. As a result,

566 the RMSEs of LW/SW fluxes and skin temperature were improved by 12.6%, 8.0%, and 4.4%

567 compared with those of WRF60. The use of four and five hidden layers resulted in a worse

568 performance than the results obtained with one hidden layer. This implies that there is an

569 optimal number of hidden layers for a given problem. Gentine et al. (2018) and Pal et al.

570 (2019) also used three and eight hidden layers as the optimal numbers of hidden layers,

571 respectively, when developing their emulators. In a similar context, the use of an optimizer

572 for tuning hyperparameters (e.g., Hertel et al., 2020), including the number of neurons and

573 hidden layers, may improve the accuracy of the training data, but it does not always

574 guarantee the generalized performance using independent test data (e.g., the overfitting

575 problem). However, the RMSEs for 2-m temperature and precipitation among the

576 experiments using different hidden layers changed within 1%, implying that the operational

577 use of the developed emulator is possible as it is now.

578 **4. Summary and Conclusions**

579     This study examined the performance of a radiation emulator based on SNN and SWA

580 training methods under idealized squalline and real case (over the Korean peninsula)

581 frameworks. Both frameworks used the WRF model with 5-km horizontal resolution, 39

582 vertical layers, a model/radiation time step of 20 s, and the RRTMG-K radiation scheme.

583 Ideal and real case simulations were integrated for 24 h and 168 h, respectively. Input

584 variables of 157–187 (ideal) and 158–190 (real), and 42 output variables were prepared, and

585 90 neurons with a single hidden layer were used in the NN training. The variables were

586 further separated into four categories (LW/SW and clear/cloud) in the ideal simulation and 96

587 categories (LW/SW, clear/cloud, land/ocean, and 12 months) in the real case simulation. The

588 weight and bias coefficients obtained from the NN training were implemented in the WRF

589 model by replacing the RRTMG-K code. The resultant radiation process was speed up 60

590 times with a total reduction in the computation time of 84–87%. In the ideal simulation,

591 sensitivity experiments were conducted examining the sampling ratio, activation functions,

592 and number of hidden layers. Regardless of the sampling ratios, SWA improved the RMSEs

593 by 10% as compared to SNN. At a sampling ratio of 10%, the performance increased even

594 further to 13.2%. Compared to the infrequent use of radiation scheme by 60 times, SNN

595    improved RMSEs by 5.8–14.1% for five forecast variables, and SWA further increased these

596    improvements by 18.2–26.9%. Among the 16 activation functions, the use of Tanh showed

597    the best performance. However, even if multiple hidden layers were considered, the

598    performance was not superior to that of the single hidden layer in the ideal simulation. The

599    final performance of the SWA was better than operational methods based on infrequent

600    radiation scheme by 3 to 60 times, suggesting improvements in both accuracy and speed for

601    SWA emulator. The ideal framework served as the testbed for various sensitivity experiments

602    before the real case simulation, which requires significant computational effort.

603        In the real case simulation, the training sets were prepared for the period 2009 to 2019. To

604    optimize batch size and learning rate, independent validation sets were prepared. After 10

605    sensitivity experiments based on the SWA, the optimal batch size and learning rate were

606    determined to be 500 and 0.05, respectively. This contributed to the mean RMSE

607    improvement averaging 8.30% for the four variables (LW/SW heating rates and fluxes)

608    compared to the SNN that was based on sequential training with one batch size. In a case

609    study, both emulators properly simulated the 156-h forecast patterns of typhoon HAISEN

610    (12LST September 17, 2020). However, SWA showed better performance for predicting skin

611    temperature with a 14% reduction in the RMSE compared to SNN. The final evaluation was

612    performed for 2020. Here, 48 cases were initialized from 1, 8, 15, and 22 days of each month,

613    which were then integrated over one week. Compared to WRF60, SNN showed 8.8% and 4.7%

614    RMSE improvements for LW and SW fluxes; however, these improvements deceased

615    significantly after a 5-day forecast, resulting the RMSE of skin temperature was increased by

616    1.8%. By contrast, the use of the SWA solved this problem, and the resultant RMSE

617    improvements were 12.3%, 7.2%, and 3.2% for LW flux, SW flux, and skin temperature,

618    respectively, compared to WRF60. These RMSEs were further improved by the use of two

619    hidden layers, to 12.6%, 8.0%, and 4.4%. This is in contrast to the ideal experiment, which

620  showed the best performance under the use of a single hidden layer. Therefore, we can

621  conclude that the use of multiple hidden layers can be helpful for optimizing forecast

622  accuracy, but it does not always guarantee better performance owing to the constraint of

623  computational cost (i.e., a smaller number of neurons should be used in the DNN). When

624  compared with surface temperature and precipitation observations, the maximum RMSE

625  difference between experiments (control run, infrequent methods of radiation scheme, and

626  emulators) was less than 1%, confirming the robustness of the developed emulators.

627  The radiation emulators in this study will replace the radiation scheme of the KMA

628  operational short-range weather forecasting model over the Korean peninsula. The one-year

629  evaluation suggests that the use of this scheme can contribute to maintaining accuracy while

630  significantly improving the computational speed of the NWP model. Operational

631  implementation should be more technically optimized through the combination of the

632  radiation emulator and its infrequent use (Song and Roh, 2021), and the use of compound

633  parameterization (Song et al., 2021). In this study, the advantages of SWA with better

634  generalization are emphasized. The strengths of SWA for long-term integration can be

635  beneficial for developing a radiation emulator that can be used for seasonal prediction or

636  multi-model climate simulations that require high computational costs (e.g., O'Neill et al.,

637  2016). Furthermore, it can be also applied to improve the NN emulation studies for other

638  physical parameterizations (Brenowitz and Bretherton, 2018; Gentine et al., 2018; Rasp et al.,

639  2018; Wang et al., 2019; Chantry et al., 2021; Mooers et al., 2021). Various sensitivity

640  experiments on important hyperparameters (activation functions, hidden layers, batch sizes,

641  and learning rates) are worthwhile. These efforts will provide guidance for future

642  development toward the total replacement of numerical weather–climate forecasting models

643  using machine learning emulators.

644  **Acknowledgements**

650 **Data Availability Statement**

651 The datasets and all sources codes are available at https://doi.org/10.5281/zenodo.5638436.

652 **References**

653 Bae, S. Y., Hong, S.-Y., & Tao, W.-K. (2019). Development of a single-moment cloud
654     microphysics scheme with prognostic hail for the Weather Research and Forecasting
655     (WRF) model. *Asia-Pacific Journal of Atmospheric Sciences,* 55, 233–245.
656     https://doi.org/10.1007/s13143-018-0066-3.

657 Baek, S. (2017). A revised radiation package of G-packed McICA and two-stream
658     approximation: Performance evaluation in a global weather forecasting model. *Journal*
659     *of Advances in Modeling Earth Systems,* 9, 1628–1640.
660     https://doi.org/10.1002/2017MS000994.

661 Belochitski, A., Binev, P., DeVore, R., Fox-Rabinovitz, M., Krasnopolsky, V., & Lamby, P.
662     (2011). Tree approximation of the long wave radiation parameterization in the NCAR
663     CAM global climate model. *Journal of Computational and Applied Mathematics*, 236,
664     447–460. https://doi.org/10.1016/j.cam.2011.07.013.

665 Belochitski, A., & Krasnopolsky, V. (2021). Robustness of neural network emulations of
666     radiative transfer parameterizations in a state-of-the-art General Circulation Model.
667     *Geoscientific Model Development, Discussions,* https://doi.org/10.5194/gmd-2021-114
668     (in press).

669 Bottou, L. (2012). Stochastic Gradient Descent tricks. *Neural Networks: Tricks of the Trade.*
670     *Lecture Notes in Computer Science*, 7700. Springer, Berlin, Heidelberg.
671     https://doi.org/10.1007/978-3-642-35289-8_25.

672 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural network
673     unified physics parameterization. *Geophysical Research Letters,* 45, 6289–6298.
674     https://doi.org/10.1029/2018GL078510.

675 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Machine
676     learning emulation of gravity wave drag in numerical weather forecasting. *Journal of*
677     *Advances in Modeling Earth Systems*, *13*, e2021MS002477.
678     https://doi.org/10.1029/2021MS002477.

Chevallier, F., Chéruy. F., Scott, N. A., & Chédin, A. (1998). A neural network approach for a fast and accurate computation of a longwave radiative budget. *Journal of Applied Meteorology,* 37, 1385–1397. https://doi.org/10.1175/1520-0450(1998)037.

Chevallier, F., Morcrette, J.-J., Chéruy, F., & Scott, N. A. (2000). Use of a neural-network-based long-wave radiative-transfer scheme in the ECMWF atmospheric model. *Quaterly Journal of the Royal Meteorological Society,* 126, 761–776. https://doi.org/10.1002/qj.49712656318.

Clough, S. A., Iacono, M. J., & Moncet, J.-L. (1992). Line-by-line calculation of atmospheric fluxes and cooling rates: Application to water vapor. *Journal of Geophysical Research,* 97, 15761–15785. https://doi.org/10.1029/92JD01419.

Clough, S. A., Shephard, M. W., Mlawer, E. J., Delamere, J. S., Iacono, M. J., Cady-Pereira, K., Boukabara, S., & Brown, P. D. (2005). Atmospheric radiative transfer modeling: a summary of the AER codes. *Journal of Quantitative Spectroscopy and Radiative Transfer,* 91, 233–244. https://doi.org/10.1016/j.jqsrt.2004.05.058.

Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could machine learning break the convection parameterization deadlock? *Geophysical Research Letters,* 45, 5742–5751. https://doi.org/10.1029/2018GL078202.

Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., Chiara, G., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., Thépaut, J.-N. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society,* 146, 1999– 2049. https://doi.org/10.1002/qj.3803.

Hertel, L., Collado, J., Sadowski, P., Ott, J., & Baldi, P. (2020). Sherpa: Robust hyperparameter optimization for machine learning. *SoftwareX,* 12, 100591, https://doi.org/10.1016/j.softx.2020.100591.

Hogan, R. J., & Bozzo, A. (2015). Mitigating errors in surface temperature forecasts using approximate radiation updates. *Journal of Advances in Modeling Earth Systems*, 7, 836–853. https://doi.org/10.1002/2015MS000455.

Iacono, M. J., Delamere, J. S., Mlawer, E. J., Shephard, M. W., Clough, S. A., & Collins, W. D. (2008). Radiative forcing by long-lived greenhouse gases: Calculations with the AER radiative transfer models. *Journal of Geophysical Research,* 113, D13103. https://doi.org/10.1029/2008JD009944.

Izmailov, P. Podoprikhin, D., Garipov, T., Vetrov, D., and Wilson, A. G. (2018). Averaging weights leads to wider optima and better generalization. *Conference on Uncertainty in Artificial Intelligence (UAI) 2018,* https://arxiv.org/abs/1803.05407.

Jiménez, P. A., Dudhia, J., González-Rouco, J. F., Navarro, J., Montávez, J. P., García-Bustamante, E. (2012). A revised scheme for the WRF surface layer formulation. *Monthly Weather Review*, 140, 898–918. https://doi.org/10.1175/MWR-D-11-00056.1.

Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). New approach to calculation of atmospheric model physics: Accurate and fast neural network emulation of longwave radiation in a climate model. *Monthly Weather Review*, 133, 1370–1383. https://doi.org/10.1175/MWR2923.1.

Krasnopolsky, V. M., Fox-Rabinovitz, M. S., Tolman, H. L., & Belochitski, A. A. (2008). Neural network approach for robust and fast calculation of physical processes in numerical environmental models: Compound parameterization with a quality control of larger errors. *Neural Networks,* 21, 535–543. https://doi.org/10.1016/j.neunet.2007.12.019.

Krasnopolsky, V. M., Fox-Rabinovitz, M. S., Hou, Y. T., Lord, S. J., & Belochitski, A. A. (2010). Accurate and fast neural network emulations of model radiation for the NCEP coupled Climate Forecast System: Climate simulations and seasonal predictions. *Monthly Weather Review*, 138, 1822–1842. https://doi.org/10.1175/2009MWR3149.1

Krasnopolsky, V. M. (2014). NCEP neural network training and validation system: Brief description of NN background and training software. Environment Modeling Center, NCEP/NWS, NOAA. https://doi.org/10.7289/v5qr4v2z.

Kwon, Y. C., & Hong, S. (2017). A mass-flux cumulus parameterization scheme across gray-zone resolutions. *Monthly Weather Review*, 145, 583–598. https://doi.org/10.1175/MWR-D-16-0034.1.

Li, M., Zhang, T., Chen, Y., and Smola, A. J. (2014). Efficient mini-batch training for stochastic optimization. *Proceedings of the 20th ACM SIGKDD international conference on knowledge discovery and data mining,* 661–670, https://doi.org/10.1145/2623330.2623612.

Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring deep learning architectures for longwave radiative transfer. *Geoscientific Model Development*, 13, 4399–4412. https://doi.org/10.5194/gmd-13-4399-2020.

Manners, J., Thelen, J.-C., Petch, J., Hill, P., & Edwards, J. M. (2009). Two fast radiative transfer methods to improve the temporal sampling of clouds numerical weather prediction and climate models. *Quaterly Journal of the Royal Meteorological Society,* 135, 457–468. https://doi.org/10.1002/qj.385.

Mandt, S., Hoffman, M. D., & Blei, D. M. (2017). *Journal of Machine Learning Research,* 18, 1–35. https://www.jmlr.org/papers/volume18/17-214/17-214.pdf.

Meyer, D., Hogan, R. J, Dueben, P. D, & Mason, S. L. (2021). Machine learning emulation of 3D cloud radiative effects. *Journal of Advances in Modeling Earth Systems,* https://arxiv.org/abs/2103.11919. (in revision)

Mooers, G., Pritchard, M., Beucler, T., Ott, J., Yacalis, G., Baldi, P., & Gentine, P. (2021). Assessing the potential of deep learning for emulating cloud superparameterization in climate models with real-geography boundary conditions. *Journal of Advances in Modeling Earth Systems*, 13, e2020MS002385. https://doi.org/10.1029/2020MS002385

O'Neill, B. C., Tebaldi, C., van Vuuren, D. P., Eyring, V., Friedlingstein, P., Hurtt, G., Knutti, R., Kriegler, E., Lamarque, J.-F., Lowe, J., Meehl, G. A., Moss, R., Riahi, K., &

Sanderson, B. M. (2016). The Scenario Model Intercomparison Project (ScenarioMIP) for CMIP6. *Geoscientific Model Development,* 9, 3461–3482. https://doi.org/10.5194/gmd-9-3461-2016.

Pal, A., Mahajan, S., & Norman, M. R. (2019), Using deep neural networks as cost-effective surrogate models for Super-Parameterized E3SM radiative transfer. *Geophysical Research Letters,* 46, 6069–6079. https://doi.org/10.1029/2018GL081646.

Pauluis, O., & Emanuel, K. (2004). Numerical instability resulting from infrequent calculation of radiative heating, *Monthly Weather Review,* 132, 673–686. https://doi.org/10.1175/1520-0493(2004)132.

Pincus, R., & Stevens, B. (2013). Paths to accuracy for radiation parameterizations in atmospheric models. *Journal of Advances in Modeling Earth Systems,* 5, 255–233. https://doi.org/10.1002/jame.20027.

Pincus, R., Mlawer, E. J., & Delamere, J. S. (2019). Balancing accuracy, efficiency, and flexibility in radiation calculations for dynamical models. *Journal of Advances in Modeling Earth Systems,* 11, 3087–3089. https://doi.org/10.1029/2019MS001621.

Roh, S., & Song, H.-J. (2020). Evaluation of neural network emulations for radiation parameterization in cloud resolving model. *Geophysical Research Letters,* 47, e2020GL089444. https://doi.org/10.1029/2020GL089444.

Shin, H. H., & Hong, S. (2015). Representation of the subgrid-scale turbulent transport in convective boundary layers at gray-zone resolutions. *Monthly Weather Review*, 143, 250–271. https://doi.org/10.1175/MWR-D-14-00116.1.

Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., Wang, W., Powers, J. G., Duda, M. G., Barker, D. M., & Huang, X.-Y. (2019). A description of the Advanced Research WRF model version 4. *NCAR Technical Notes.* https://doi.org/10.5065/1DFH-6P97.

Smith, S. L., Kindermans, P.-J., Ying, C., & Ye, Q. V. (2018). Don't decay the learning rate, increase the batch size. *6th International Conference on Learning Representations (ICLR 2018),* https://arxiv.org/abs/1711.00489.

Song, H.-J., & Roh, S. (2021). Improved weather forecasting using neural network emulation for radiation parameterization. *Journal of Advances in Modeling Earth Systems,* 13, e2021MS002609, https://doi.org/10.1029/2021MS002609.

Song, H.-J., Roh, S., & Park, H. (2021). Compound parameterization to improve the accuracy of radiation emulator in a numerical weather prediction model. *Geophysical Research Letters,* 48, e2021GL095043, https://doi.org/10.1029/2021GL095043.

Tewari, M., Chen, F., Wang, W., Dudhia, J., LeMone, M., Mitchell, K., Ek, M., Gayno, G., Weigel, J., & Cuenca, R. (2004). Implementation and verification of the unified Noah land surface model in the WRF model. *20th Conference on Weather Analysis and Forecasting/16th Conference on Numerical Weather Prediction,* American Meteorological Society, Seattle, WA, 11 – 15 Jan.

800    Ukkonen, P., Pincus, R., Hogan, R. J., Nielsen, K. P., & Kaas, E. (2020). Accelerating
801            radiation computations for dynamical models with targeted machine learning and code
802            optimization. *Journal of Advances in Modeling Earth Systems,* 12, e2020MS002226.
803            https://doi.org/10.1029/2020MS002226.

804    Vapnik, V. N. (2019). Complete statistical theory of learning. *Automation and Remote*
805            *Control,* 80, 1949–1975. https://doi.org/10.1134/S000511791911002X.

806    Veerman M. A., Pincus, R., Stoffer, R., van Leeuwen, C. M., Podareanu, D., & van
807            Heerwaarden, C. C. (2021). Predicting atmospheric optical properties for radiative
808            transfer computations using neural networks. *Philosophical Transactions of the Royal*
809            *Society A,* 379, 20200095. https://doi.org/10.1098/rsta.2020.0095.

810    Wang, J., Balaprakash, P., & Kotamarthi, R. (2019). Fast domain-aware neural network
811            emulation of a planetary boundary layer parameterization in a numerical weather
812            forecast model. *Geoscientific Model Development,* 12, 4261–4274.
813            https://doi.org/10.5194/gmd-12-4261–2019.

814 **Table 1.** List of inputs and outputs for longwave (LW) and shortwave (SW) emulators. The
815 numbers of inputs decreased by 157 and 158 for ideal and real cases under clear conditions,
816 respectively, because cloud fractions were not used.

| Inputs (ideal case) | # |
|---|---|
| Pressure | 1–39 |
| Temperature | 40–78 |
| Water Vapor | 79–117 |
| Ozone | 118–156 |
| Cloud Fraction | 157–186 |
| Skin Temperature (LW) | 187 |
| Solar Constant × Cosine Zenith Angle (SW) | 187 |
| Inputs (real case) | # |
| Pressure | 1–39 |
| Temperature | 40–78 |
| Water Vapor | 79–117 |
| Ozone | 118–156 |
| Cloud Fraction | 157–188 |
| Skin Temperature (LW) | 189 |
| Surface Emissivity (LW) | 190 |
| Solar Constant × Cosine Zenith Angle (SW) | 189 |
| Surface albedo (SW) | 190 |
| Outputs | # |
| Heating Rate (LW, SW) | 1–39 |
| Upward Flux at the Top (LW, SW) | 40 |
| Upward Flux at the Bottom (LW, SW) | 41 |
| Downward Flux at the Bottom (LW, SW) | 42 |

817

818

819     **Table 2.** Definitions of the activation functions used. All empirical coefficients were based
820     on the default settings in pytorch.

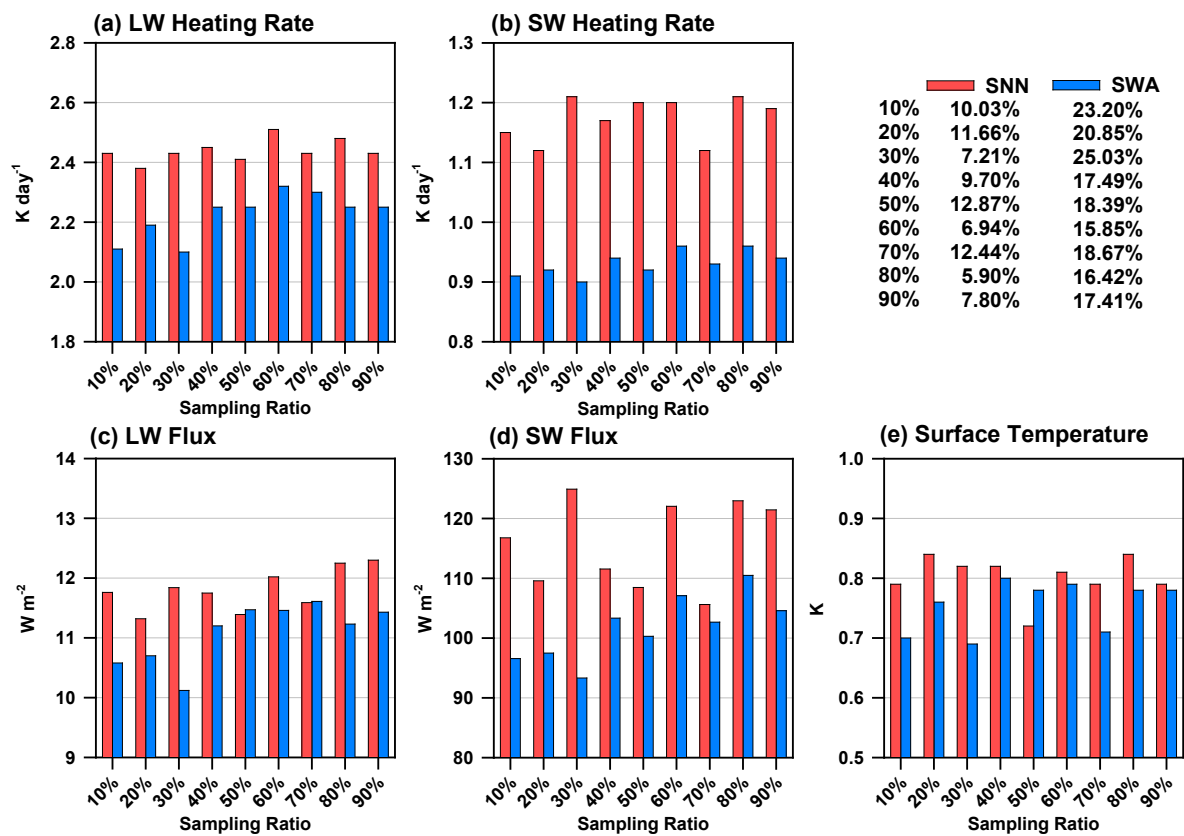| # | Functions | Equations | Ranges |
|---|---|---|---|
| 1 | Tanh | $(\exp(x) - \exp(-x)) \div (\exp(x) + \exp(-x))$ | $-1, 1$ |
| 2 | Arctan | $\tan^{-1}(x)$ | $-\pi/2, \pi/2$ |
| 3 | Tanhshrink | $x - \tanh(x)$ | $-\infty, \infty$ |
| 4 | Sigmoid | $1 \div (1+\exp(-x))$ | $0, 1$ |
| 5 | Logsigmoid | $\log(1 \div (1+\exp(-x)))$ | $-\infty, 0$ |
| 6 | SiLU | $x \div (1+\exp(-x))$ | $0, \infty$ |
| 7 | Softsign | $x \div (1+|x|)$ | $-1, 1$ |
| 8 | Softplus | $\log(1+\exp(x))$ | $0, \infty$ |
| 9 | Mish | $x \times \tanh(\mathrm{softplus}(x))$ | $0, \infty$ |
| 10 | Hardtanh | $[-1, x \leq -1], [x, -1 < x < 1], [1, x \geq 1]$ | $-1, 1$ |
| 11 | Hardsigmoid | $[0, x \leq -3], [x \div 6 + 1 \div 2, -3 < x < 3], [1, x \geq 3]$ | $0, 1$ |
| 12 | Hardswish | $[0, x \leq -3], [x \times (x+3) \div 6, -3 < x < 3], [x, x \geq 3]$ | $0, \infty$ |
| 13 | ReLU | $\max(0,x)$ | $0, \infty$ |
| 14 | LeakyReLU | $\max(0,x) + 0.01 \times \min(0,x)$ | $-\infty, \infty$ |
| 15 | ELU | $[x, x > 0], [\exp(x) -1, x \leq 0]$ | $-1, \infty$ |
| 16 | SELU | $\alpha \times (\max(0,x)+\min(0, \beta \times (\exp(x) -1)))$ <br> $\alpha = 1.0507009873554804934193349852946$ <br> $\beta = 1.6732632423543772848170429916717$ | $-\alpha \times \beta, \infty$ |

821

822

823 **Table 3.** Statistical results of the idealized squalline simulation under the infrequent use of
824 radiation scheme by 9, 30, and 60 times (WRF9, WRF30, and WRF60), and the SNN/SWA
825 emulation results compared to the control run. Total improvement is the relative reduction of
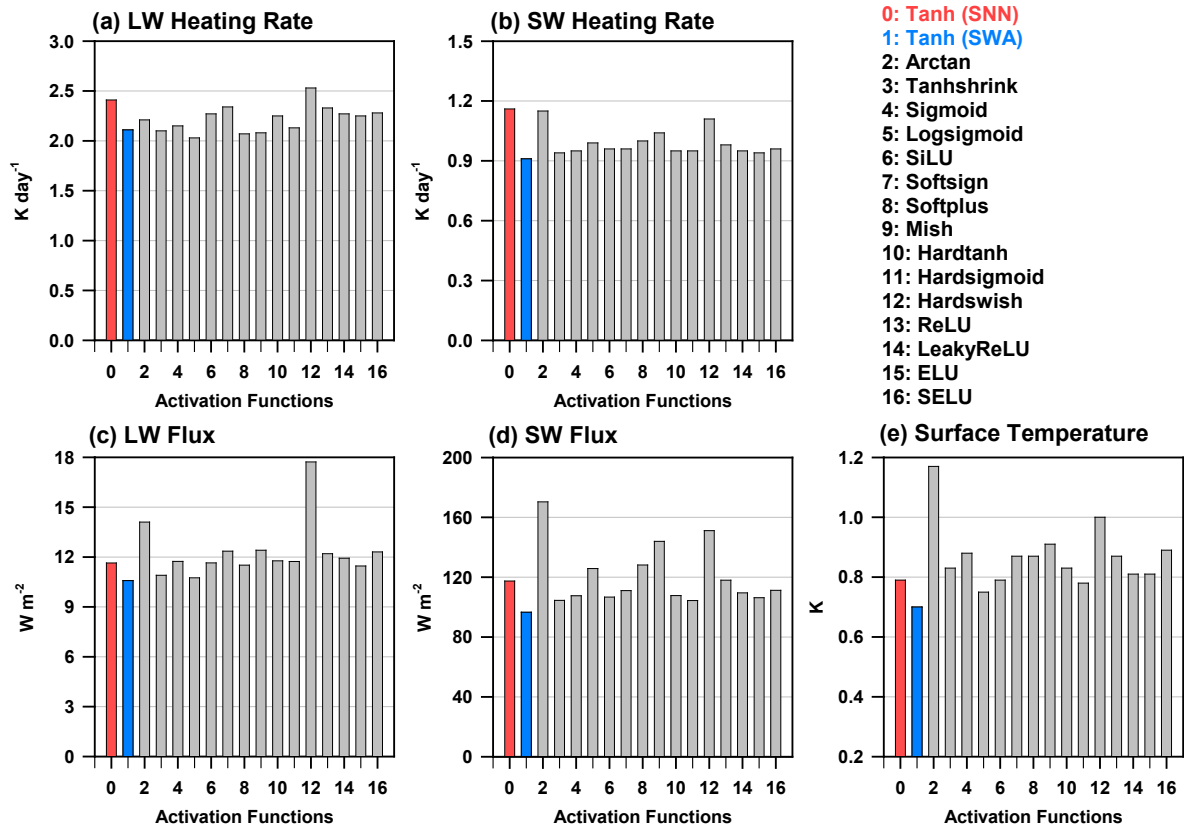826 RMSE (%) in WRF60 for five variables (LW/SW hearing rates, LW/SW flux, and surface
827 temperature).

| Experiments | WRF9 | WRF30 | WRF60 | SNN | SWA |
|---|---|---|---|---|---|
| Radiation time step (radt) | 3 m | 10 m | 20 m | 20 s | 20 s |
| Speedup of radiation | 9 | 30 | 60 | 59.7 | 60.1 |
| Reduced total time | 75.56% | 82.17% | 83.58% | 83.61% | 83.69% |
| LW heating rate [K day$^{-1}$] | 2.40 | 2.57 | 2.58 | 2.43 | 2.11 |
| SW hearing rate [K day$^{-1}$] | 1.16 | 1.20 | 1.24 | 1.15 | 0.91 |
| LW flux [W m$^{-2}$] | 11.12 | 12.28 | 13.29 | 11.76 | 10.58 |
| SW flux [W m$^{-2}$] | 102.08 | 113.43 | 132.15 | 116.78 | 96.56 |
| Surface temperature [K] | 0.72 | 0.77 | 0.92 | 0.79 | 0.70 |
| Total improvement (%) | 14.74 | 8.21 | - | 10.03 | 23.20 |

828

829

830 **Table 4.** Root mean square error (RMSE) results of fluxes and skin temperature ($T_s$) in the
831 real case simulation under the infrequent use of radiation scheme by 15, 30, and 60 times
832 (WRF15, WRF30, and WRF60), the SNN, and the SWA with one to five hidden layers (1 h to
833 5 h), compared to the control run. The results of 2-m temperature ($T_{2m}$) and 3-h accumulated
834 precipitation were produced through comparison with surface observations in South Korea.
835 Note that the RMSE of the control run for 2-m temperature and precipitation observations
836 were 2.2581 K and 12.3526 mm, respectively.

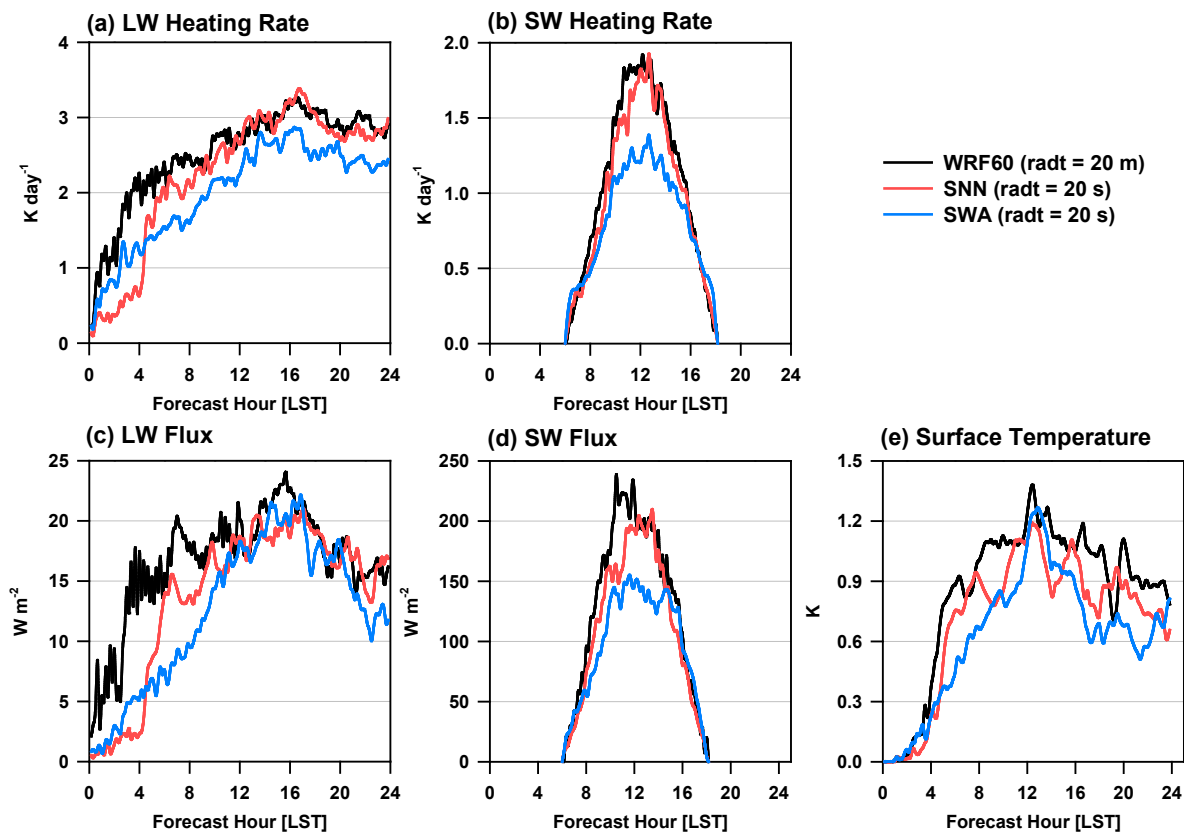| Experiments | LW flux [W m$^{-2}$] | SW flux [W m$^{-2}$] | $T_s$ [K] | $T_{2m}$ [K] | Precipitation [mm] |
|---|---|---|---|---|---|
| WRF15 | 7.8756 | 53.9819 | 0.5371 | 2.2590 | 12.2649 |
| WRF30 | 8.6558 | 57.6258 | 0.5753 | 2.2532 | 12.1987 |
| WRF60 | 10.1513 | 64.8639 | 0.6602 | 2.2438 | 12.2897 |
| SNN | 9.2629 | 61.8149 | 0.6721 | 2.2466 | 12.3120 |
| SWA (1h) | 8.9027 | 60.2215 | 0.6389 | 2.2563 | 12.2551 |
| SWA (2h) | 8.8680 | 59.6838 | 0.6309 | 2.2487 | 12.2944 |
| SWA (3h) | 8.9614 | 59.9000 | 0.6390 | 2.2470 | 12.3060 |
| SWA (4h) | 9.2006 | 60.9223 | 0.6563 | 2.2424 | 12.2800 |
| SWA (5h) | 9.4009 | 62.1192 | 0.6559 | 2.2593 | 12.2230 |

837

838

839

840

**Figure 1.** Sensitivity experiments with the ratio of training sets. The SNN and SWA results are represented by the ratio of training sets to full sets. Statistical values denote the RMSE using 5-km and 20-s intervals over the entire domain and period compared with the control run (radt = 20 s). Compared to the WRF60, the mean reduced RMSEs for five variables and nine ratios are presented in the upper right corner.
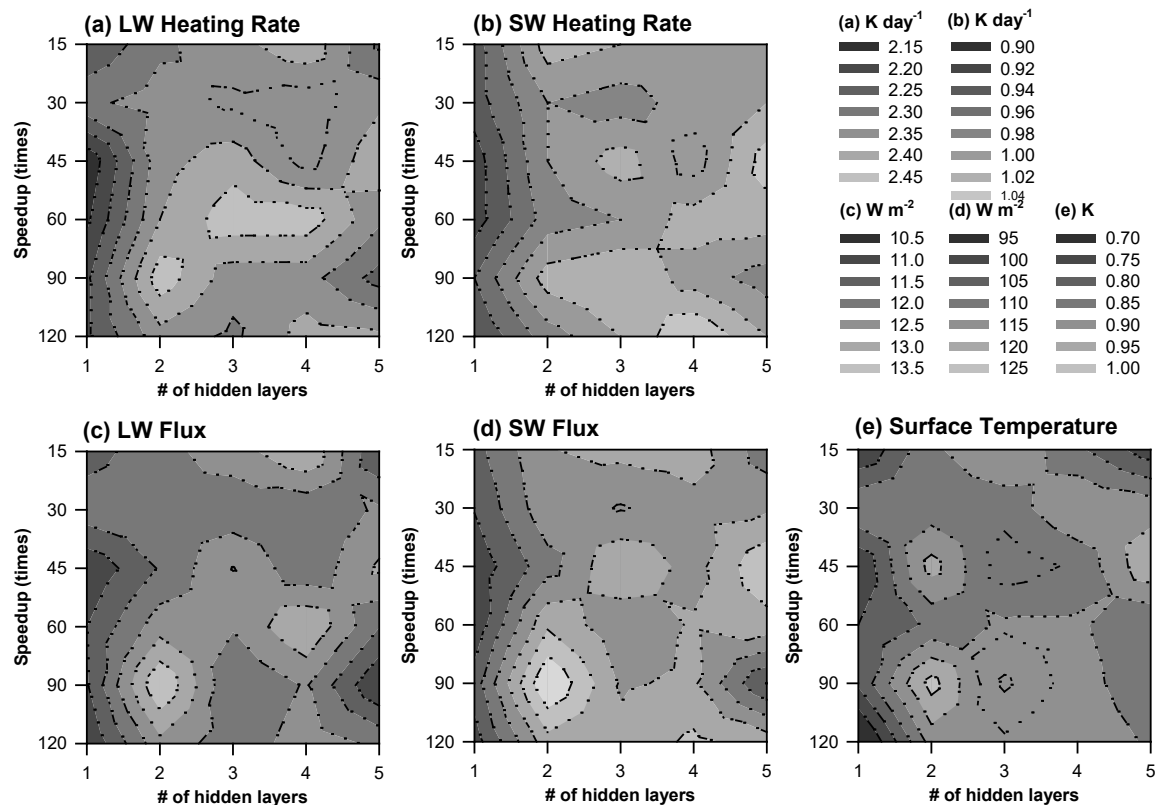
**Figure 2.** Sensitivity experiments with activation functions for (a) LW heating rate, (b) SW heating rate, (c) LW flux, (d) SW flux, and (e) surface temperature. Vertical bars denote the RMSE with 5-km and 20-s intervals over the entire domain and a 24-h period compared with the control run (radt = 20 s). The SNN is displayed as the red bar and the best experiment among the SWA experiments is highlighted as the blue bar.

859
860 **Figure 3.** Evolutionary features for idealized squalline simulation. The control run, WRF60
861 (radt = 20 m), SNN, and SWA results are displayed for LW and SW upward fluxes at the top
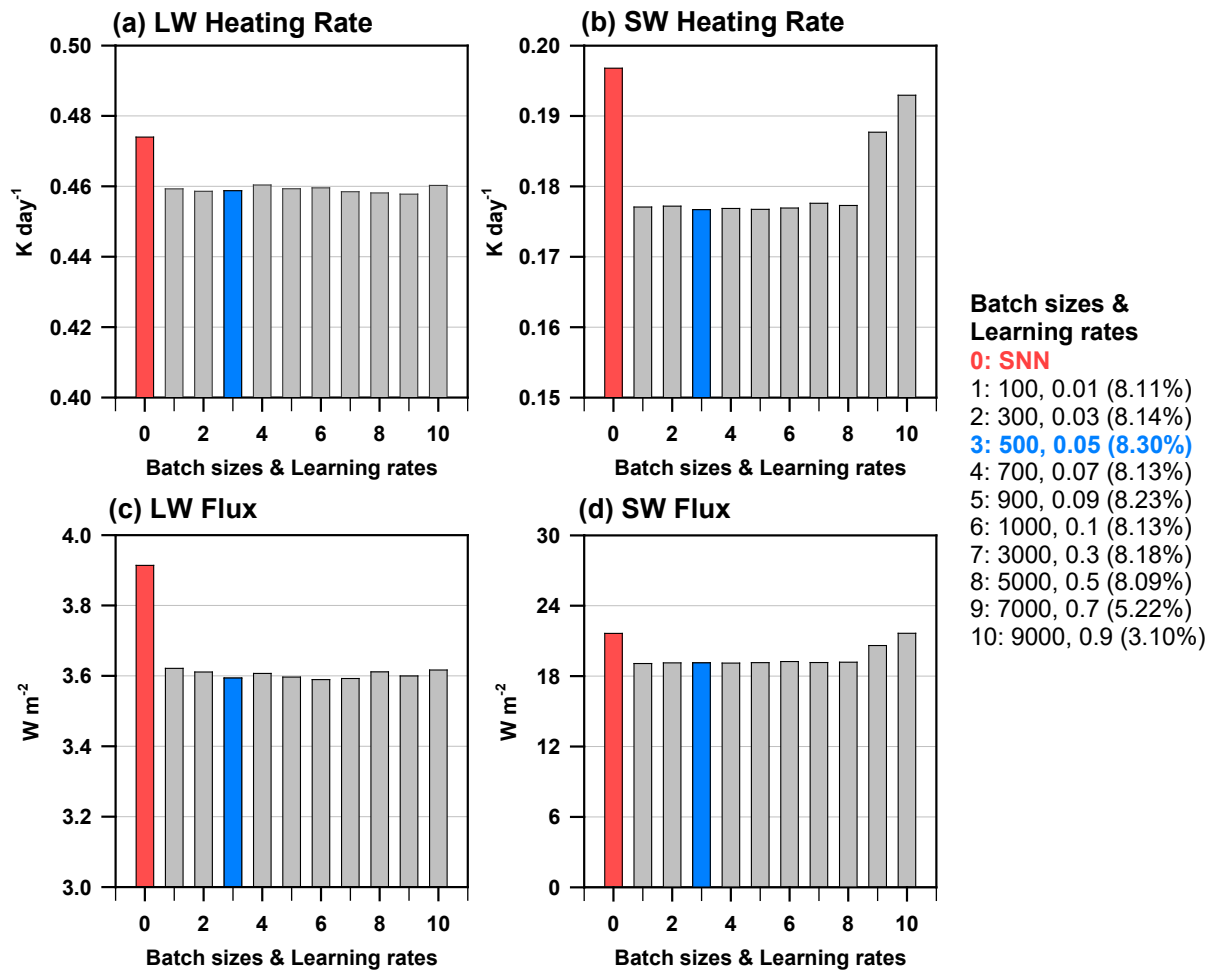862 (LWUPT and SWUPT), surface temperature, and precipitation rate.

863
864
865

**Figure 4.** Times series of RMSEs for (a) LW heating rate, (b) SW heating rate, (c) LW flux, (d) SW flux, and (e) surface temperature. The mean reductions in the RMSE for five variables compared to WRF60 are given in the upper right corner.
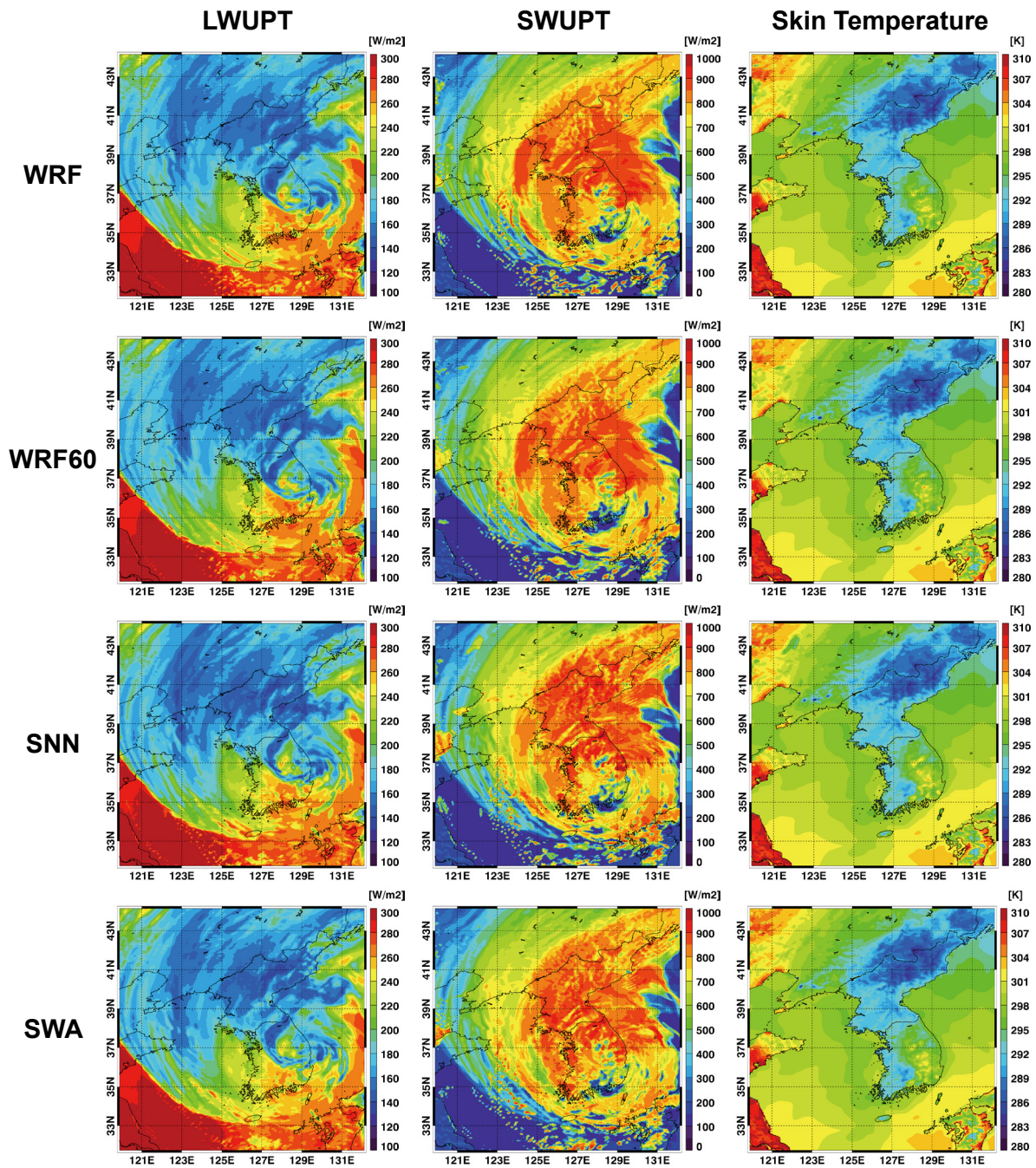
**Figure 5.** Sensitivity experiments with hidden layers and speedups for (a) LW heating rate, (b) SW heating rate, (c) LW flux, (d) SW flux, and (e) surface temperature. The speedups of 15, 30, 45, 60, 90, and 120 times correspond to the use of 360, 180, 120, 90, 60, and 45 neurons for the case of single hidden layer. For the case of multiple hidden layers, the reduced neurons were used to maintain the same numerical complexity and resulting speedup. The values inside each figure denote the RMSE with 5-km and 20-s intervals over the entire domain and a 24-h period compared with the control run.
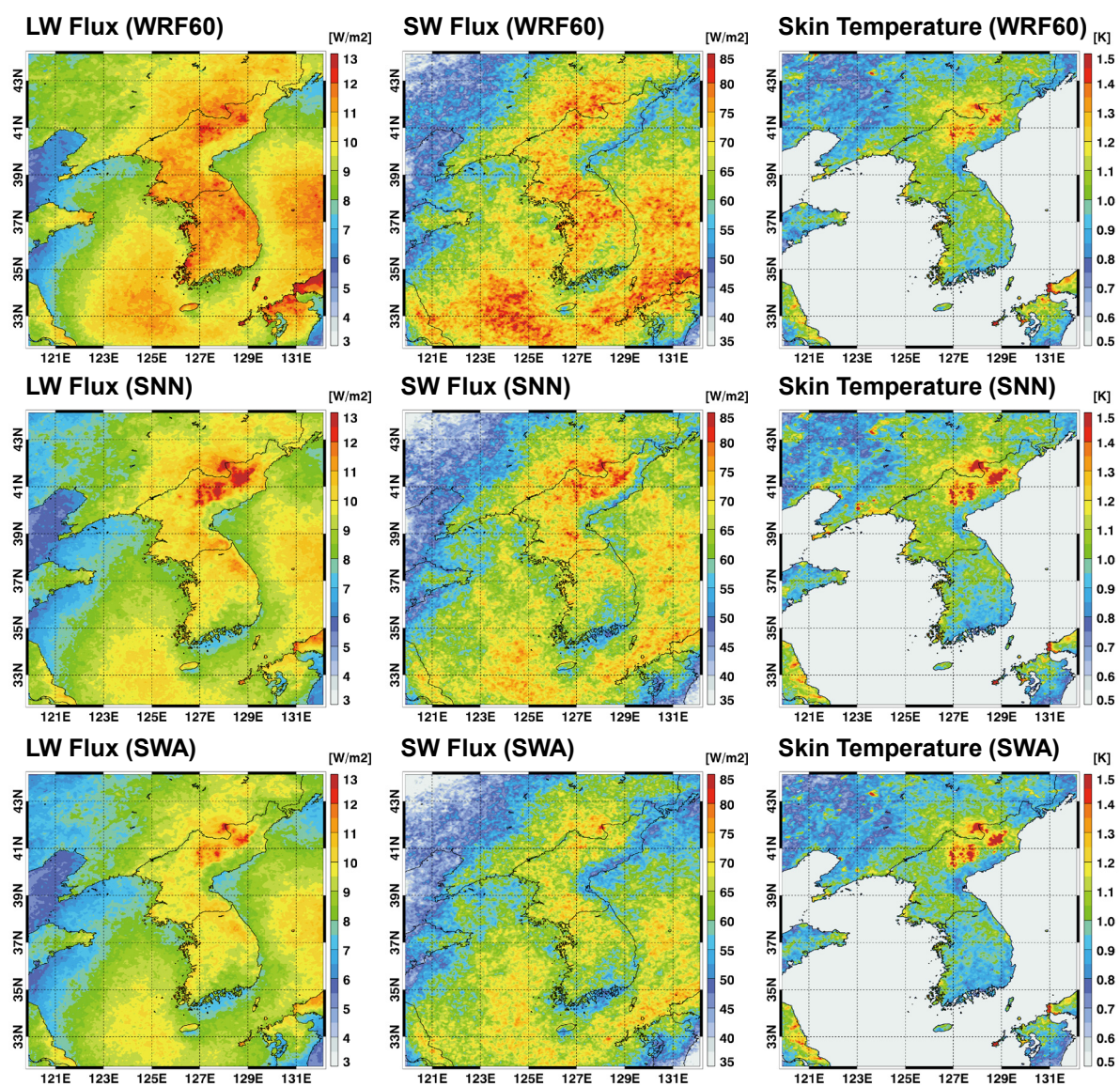
**(a) LW Heating Rate**

**(b) SW Heating Rate**

**(c) LW Flux**

**(d) SW Flux**

**Batch sizes & Learning rates**
**0: SNN**
1: 100, 0.01 (8.11%)
2: 300, 0.03 (8.14%)
**3: 500, 0.05 (8.30%)**
4: 700, 0.07 (8.13%)
5: 900, 0.09 (8.23%)
6: 1000, 0.1 (8.13%)
7: 3000, 0.3 (8.18%)
8: 5000, 0.5 (8.09%)
9: 7000, 0.7 (5.22%)
10: 9000, 0.9 (3.10%)

882
883 **Figure 6.** Sensitivity experiments with batch sizes and learning rates based on the SWA. The
884 RMSE values of (a) LW heating rate, (b) SW heating rate, (c) LW flux, and (d) SW flux for
885 validation sets are given in each figure. The percentages in the right corner denote the mean
886 RMSE improvements for four variables compared with SNN. This is an offline validation
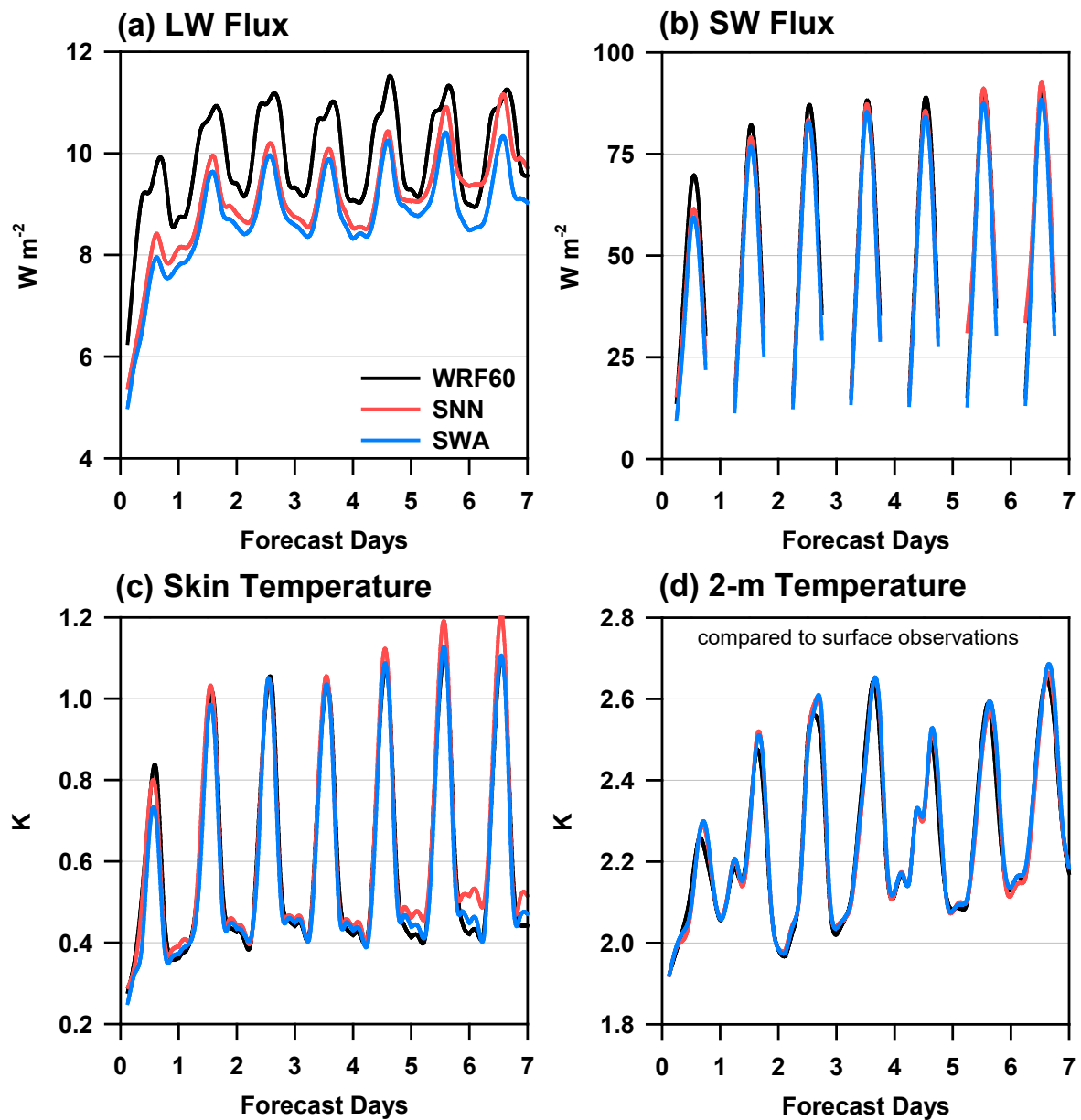887 which is not linked to the WRF simulation.

888

889

**Figure 7.** Example for Typhoon HAISEN (12LST September 7, 2020). Because the initial conditions started at 00LST 1 September 2020, it is 156-h forecast result. The control run, WRF60 (radt = 20 m), SNN, and SWA results are displayed for LW and SW upward fluxes at the top (LWUPT and SWUPT), and surface temperature.
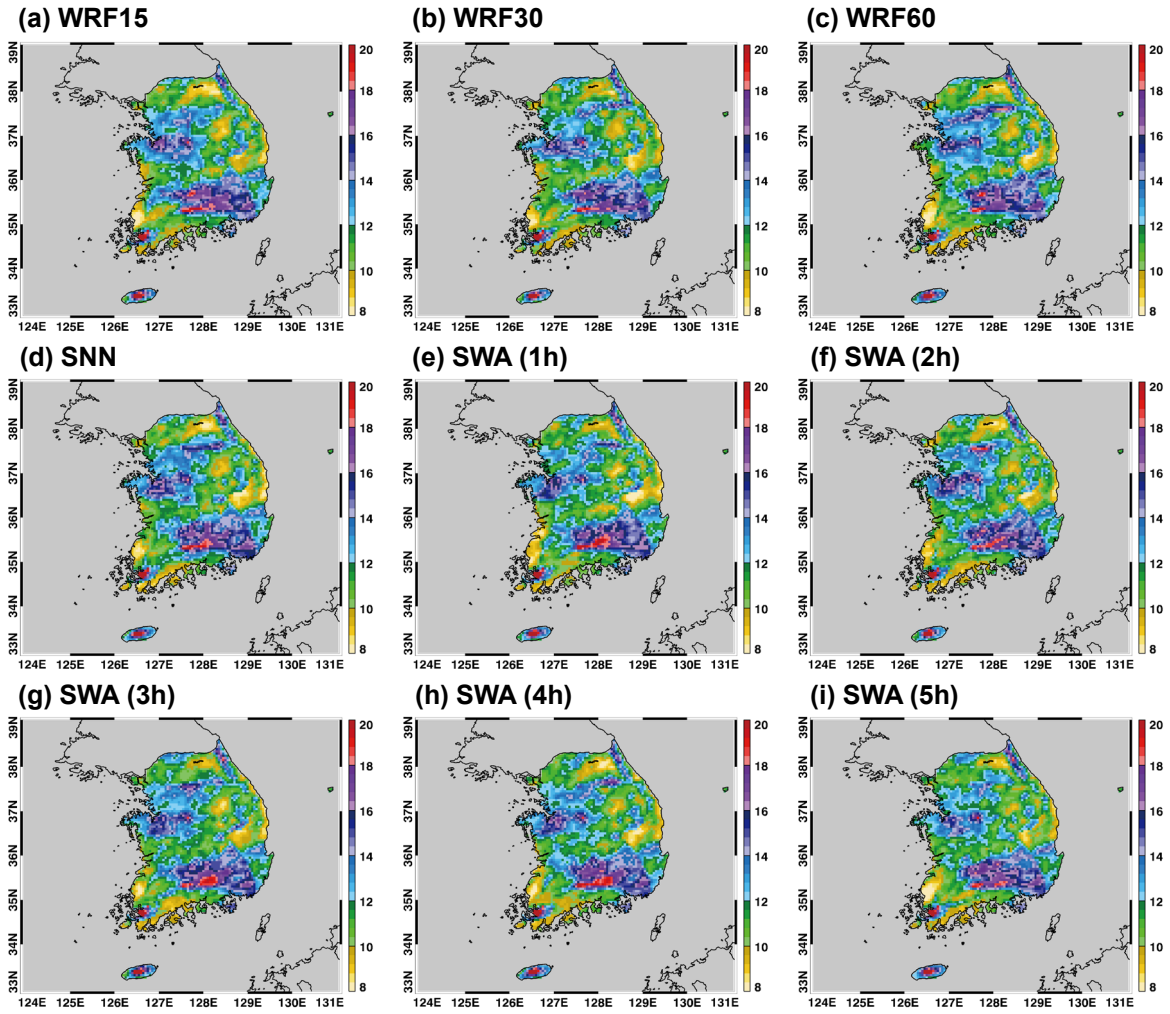
**Figure 8.** RMSE distributions of LW flux, SW flux, and skin temperature ($T_s$) for the WRF60 (radt = 20 m), SNN, and SWA compared with the control run. Each RMSE at a given 5-km grid represents a statistical result for one-week forecasts over 48 simulations of 2020.

**Figure 9.** Times series of RMSEs for (a) LW flux, (b) SW flux, (c) skin temperature, and (d) 2-m air temperature compared with surface observations in South Korea. The RMSE represents a statistical result over the entire domain or points (for 2-m temperature) and one-year period. The WRF60 (radt = 20 m), SNN, and SWA results are compared.

911
912 **Figure 10.** RMSE distributions of 3-h accumulated precipitation (mm) compared with the
913 observations in South Korea. The results of infrequent radiation scheme (WRF15, WRF30,
914 and WRF60), SNN, and SWA (one to five hidden layers; 1 h to 5 h) are compared.