Revisiting Machine Learning Approaches for Short- and Longwave Radiation Inference in Weather and Climate Models, Part II: Online Performance

Guillaume Bertoli¹, Salman Mohebi², Firat Ozdemir³, Jonas Jucker⁴, Stefan Rüdisühli³, Fernando Perez-Cruz², and Sebastian Schemm³

¹ETHZ ²Swiss Data Science Center, ETH Zurich ³ETH Zurich ⁴Center for Climate Systems Modeling (C2SM)

March 29, 2024

Abstract

This paper continues the exploration of \gls{ml} parameterization for radiative transfer for the \gls{icon}. Three \gls{ml} models, developed in Part I of this study, are coupled to \gls{icon}. More specifically, a UNet model and a bidirectional \gls{rnn} with \gls{lstm} are compared against a random forest. The \gls{ml} parameterizations are coupled to the \gls{icon} code that includes OpenACC compiler directives to enable \glspl{gpu} support. The coupling is done through Infero, developed by ECMWF, and PyTorch-Fortran. The most accurate model is the bidirectional \gls{rnn} with hysics-informed normalization strategy and heating rate penalty, but the fluxes above 15\,km height are computed with a simplified formula for numerical stability reasons. The presented setup enables stable aquaplanet simulations with \gls{icon} for several weeks at a resolution of about 80\,km and compare well with the physics-based radiative transfer solver ecRad. However, the achieved speed up when using the emulators and the minimum required memory usage relative to the \gls{gpu}-enabled ecRad depend strongly on the \gls{n} architecture. Future studies may explore physics-constraint emulators that predict heating rates inside the atmospheric model and fluxes at the top.











Revisiting Machine Learning Approaches for Shortand Longwave Radiation Inference in Weather and Climate Models, Part II: Online Performance

Guillaume Bertoli¹, Salman Mohebi², Firat Ozdemir², Jonas Jucker⁴, Stefan Rüdisühli¹, Fernando Perez-Cruz^{2,3}, and Sebastian Schemm¹

¹Institute for Atmospheric and Climate Science, ETH Zurich, Zurich, Switzerland ²Swiss Data Science Center, ETH Zurich and EPFL, Zurich, Switzerland ³Computer Science Department, ETH Zurich, Zurich, Switzerland ⁴Center for Climate Systems Modeling, ETH Zurich, Zurich, Switzerland

Key Points:

ues is required.

1

2

3

4

5

6

7 8 9

10

17

11	•	The Fortran-based and OpenACC-enhanced ICON weather and climate model is
12		coupled to a neural network radiation solver developed in Python and both run
13		in tandem on graphic processing units (GPUs). The resulting speed-up critically
14		depends on the architecture of the neural network.
15	•	ICON with the radiation emulator runs stable for several weeks with a negligible
16		difference to ecRad, but tuning at the top of the model at very low pressure val-

 Future research could explore physics-informed radiation emulators that predict heating rates inside the atmosphere and fluxes at the surface and model top as auxiliary fields.

 $Corresponding \ author: \ Guillaume \ Bertoli, \ \tt guillaume.bertoli@env.ethz.ch$

21 Abstract

This paper continues the exploration of machine Learning (ML) parameterization for ra-22 diative transfer for the ICOsahedral Nonhydrostatic weather and climate model (ICON). 23 Three ML models, developed in Part I of this study, are coupled to ICON. More specif-24 ically, a UNet model and a bidirectional recurrent neural network (RNN) with long short-25 term memory (LSTM) are compared against a random forest. The ML parameteriza-26 tions are coupled to the ICON code that includes OpenACC compiler directives to en-27 able GPUs support. The coupling is done through Infero, developed by ECMWF, and 28 PyTorch-Fortran. The most accurate model is the bidirectional RNN with physics-informed 29 normalization strategy and heating rate penalty, but the fluxes above 15 km height are 30 computed with a simplified formula for numerical stability reasons. The presented setup 31 enables stable aquaplanet simulations with ICON for several weeks at a resolution of about 32 80 km and compare well with the physics-based radiative transfer solver ecRad. How-33 ever, the achieved speed up when using the emulators and the minimum required mem-34 ory usage relative to the GPU-enabled ecRad depend strongly on the Neural Network 35 (NN) architecture. Future studies may explore physics-constraint emulators that pre-36

dict heating rates inside the atmospheric model and fluxes at the top.

³⁸ Plain Language Summary

Machine learning (ML) methods could drastically accelerate existing parts of weather 39 and climate models. This research explores machine learning methods to replace the ra-40 diation solver responsible for computing the solar and terrestrial radiative fluxes in the 41 ICON model. Three ML models are trained for this task: a random forest, which com-42 bine decision trees to make accurate prediction, and two neural networks, an increasingly 43 popular deep learning model which learn from data to perform tasks using interconnected 44 mathematical function. Because the ML models and ICON are implemented with dif-45 ferent programming languages, a auxiliary software is used to couple them. Simulations 46 with the ML models show accurate results for two-weeks simulations but the accelera-47 tion depends strongly on the ML method. 48

49 1 Introduction

This paper explores machine Learning (ML) parameterizations of simulated radia-50 tive transfer in the atmosphere. This second part of our two-part study presents a de-51 tailed description of how the offline trained ML models, which are developed using Py-52 Torch (Paszke et al., 2017) and Tensorflow (Abadi et al., 2015), are incorporated into 53 the Fortran code (Kedward et al., 2022) of the ICOsahedral Nonhydrostatic weather and 54 climate model (ICON) (MPI-M et al., 2024; Giorgetta et al., 2022), enabled for graph-55 ics processing units by means of OpenACC (OpenACC, version 3.2, 2022), and the steps 56 required to make the radiation emulator and ICON performant on graphic processing 57 units (GPUs) on Piz Daint, a Cray XC50 machine at the Swiss National Supercomput-58 ing Centre (CSCS). 59

As discussed in detail in the first part, radiative transfer emulation has a long his-60 tory in atmospheric modelling. Previous research includes studies that emulate the en-61 tire radiative transfer code (Chevallier et al., 1998; Krasnopolsky et al., 2005; Ukkonen, 62 2022; Roh & Song, 2020; Pal et al., 2019; Lagerquist et al., 2021), while others focus on 63 replacing individual components of existing radiative schemes (Ukkonen et al., 2020; Veer-64 man et al., 2021). Some studies use radiative flux as a training target (Chevallier et al., 65 1998; Ukkonen, 2022), while others directly predict the resulting heating rates (Krasnopolsky 66 et al., 2005; Roh & Song, 2020; Pal et al., 2019; Lagerquist et al., 2021). The advantages 67 and disadvantages of emulating shortwave and longwave radiative fluxes or directly pre-68 dicting the resulting heating rates are known: Direct prediction of heating rates allows 69 the calculation of radiative flux convergence to be omitted, thus avoiding numerical sta-70

bility problems, especially when the predicted flux profile is not smooth. However, the
prediction of fluxes seems to be a necessity for an earth system model (ESM), since the
radiative flux serves as an input to other components, such as the land model. It is also
a relevant output variable as it serves as an input to impact models (e.g., solar energy
production) and can be compared with measurements, which is important for validation.
For a more detailed literature review, the reader is referred to Part I of this study.

In Part I of this study (Bertoli et al., 2024), the decision was made to employ a sin-77 gle Neural Network (NN) to predict both short- and longwave up- and downward fluxes, 78 using the resulting heating rates as an additional penalty to the loss function during the 79 training. Additionally, a physics-informed normalization strategy proved beneficial for 80 enhanced training and inference accuracy. Longwave training data was normalized us-81 ing the Stefan-Boltzmann law and shortwave data was normalized using the cosine of 82 the solar zenith angle multiplied by the solar constant. This additional physics-informed 83 regularization also benefited the random forest (RF) baseline model. In general, the ac-84 curacy of the RF scaled with its memory usage, but at fixed memory usage the RF was 85 typically outperformed by most NN architectures. The only exception was at the top of 86 the atmosphere (TOA), where only the most advanced NN, a recurrent neural network 87 (RNN), showed higher accuracy than the RF in predicting fluxes. However, a particu-88 lar advantage of the RF, besides its simplicity during (re)training, was that the predicted 89 flux profile is relatively smooth, and hence the derived heating rates did not yield spu-90 rious peaks as observed with some NNs. This feature make the RF an ideal candidate 91 for low-cost, fast and simple radiation flux emulation. 92

The structure of Part II is as follows: Section 2 describes the implementation details for the ML radiative transfer parametrizations. Section 3 presents offline and online results. We conclude the paper with a summary in Section 4.

96 2 Methods

We start with the dataset construction for training the ML models, followed by the architecture choices for the RF and NN emulators, including normalization of inputs and outputs, the loss function for training the NNs, and a modification to the computation of fluxes above 15 km height to ensure ICON's stability at height-levels with very low pressure. Finally, we explain the coupling of the ML emulators with ICON using Infero (Antonino Bonanni, 2022) developed at ECMWF and PyTorch-Fortran (Alexeev, 2022).

103 **2.1 Dataset**

A two-years-long ICON aquaplanet simulation is performed with a physics timestep of 3 min and a horizontal resolution of approximately 80 km (ICON grid R02B05). The radiation parameterization ecRad is called every time-step and at full spatial resolution. The simulation uses 70 vertical levels and thus 71 vertical half levels, ranging from index 70 at the surface to index 0 at the TOA (65 km height). The solar zenith angle is held constant at equinox.

The first year of the simulation is considered as a spinup phase during which no 110 data are stored. Starting at the beginning of the second year, the required inputs and 111 outputs are stored every 3 h and 3 min (183 min output interval). This is slightly differ-112 ent from what is done in Part I (Bertoli et al., 2024), where data are stored every 3 h. 113 The new strategy allows for a more complete coverage of different Sun positions over time 114 in the dataset, which proved beneficial for the training. After 61 simulated days (480 time-115 steps stored every 183 min), the dataset contains all possible angles based on the 3 min 116 time-step interval. All NNs trained with this dataset outperform the old models. 117

The dataset is split into two parts: the initial 270 days for model training and the last 88 days for testing. A 7-day gap is included between these datasets to allow the test set distribution to vary from the training set. Within the initial segment, two sub-sections are established. The first and last 20 days of the training dataset serve as a validation dataset for the ML models during training, aiding in determining when to halt the training process using an early stopping criterion. The remaining 230 days are designated as the actual model training dataset.

In ICON, the dynamical core, the horizontal diffusion, the tracer advection, the fast physics and then the slow physics processes are solved sequentially. The radiative transfer parameterization is part of the slow physics processes, which are solved in parallel. It is therefore essential to store the states of the inputs after they have been updated by the fast physics but before the slow physics update. For this reason, we modified the ICON code to extract the states of the input variables to ecRad in the middle of the sequential time splitting.

132

2.2 Radiation emulation: ML architecture

In Part I (Bertoli et al., 2024), we explored different ML architectures as possible emulators of ecRad. First, an RF emulator composed of 10 trees served as a baseline model and was used to assess the performance of the NNs. While RFs could get very accurate, their memory footprints quickly surpassed 100GB and they thus became prohibitively memory-intensive. For this reason, RF size was constrained by imposing a minimum leaf size equal to 0.01% of the training dataset size.

Second, three NN architectures were explored. A baseline feed-forward multilayer 139 perceptron (MLP) served to evaluate the performance of more complex architectures. 140 This MLP architecture, since it was outperformed by more complex architectures in Part I, 141 is not considered further in this paper for the online tests. A convolutional NN was con-142 structed, more precisely a UNet, which allowed us to reduce the number of parameters 143 significantly. Lastly, a RNN was constructed, more precisely a bidirectional RNN with 144 long short-term memory (LSTM), which was the most accurate model. The RNN closely 145 imitates the ecRad parameterization, which solves the effect of each atmospheric layer 146 sequentially. The accuracy gain of the RNN compared to the UNet might result from 147 the fact that the RNN has access to its own prediction in the layers above (for down-148 ward fluxes) and below (for upward fluxes). However, as discussed in Section 3.2, the 149 RNN emulator is significantly slower than the UNet emulator and requires more mem-150 ory. In Tab. 1, the exact number of layers and number of neurons for each of the two NNs 151 are listed. To improve the accuracy of the emulators, normalization strategies and cus-152 tom loss functions are employed as detailed in Bertoli et al. (2024) and briefly reiterated 153 here. 154

The ML model outputs normalizes short- and longwave up- and downward fluxes. 155 The normalization is chosen such that for each atmospheric column, the model returns 156 values that are approximately in the range [0-1] as in Ukkonen (2022). An exception are 157 columns without incoming solar radiation, where the model returns zero shortwave fluxes 158 at each height by definition. Each atmospheric column's shortwave fluxes are divided by 159 the cosine of the solar zenith angle multiplied by the solar constant (approx. $1400 \,\mathrm{kW \, m^{-2}}$). 160 Atmospheric columns whose cosine of solar zenith angle is smaller than 10^{-4} are not nor-161 malized in the shortwave because these are truncated in ICON for numerical stability 162 reasons at the day-night interface. Following the Stefan-Boltzmann law (Petty, 2006)[Chap-163 ter 6.1.3] for the emission of a black body, the longwave fluxes are divided by the fourth 164 power of the surface temperature multiplied by the Stefan-Boltzmann constant. This nor-165 malization strategy improves the results of all NNs and of the RF. Each input field is 166 also normalized by subtracting its mean and dividing by its standard deviation computed 167 from the training dataset. Note that the input features are normalized across all heights, 168

Table 1: NN architectures.

Model	Architecture Overview
UNet	The 2D features are broadcasted and concatenated with the 3D features (along height axis with size 70). A UNet with convolutional units [128, 256, 512, 1024] and pooling layers [1, 2, 5, 7] are then applied. Finally, a 1D convolutional layer with 4 layers along the height axis and a final dense layer maps to the outputs to size 71x4.
RNN	The 2D features are broadcasted and concatenated with the 3D input (along height axis with size 70). The features are then concatenated with a constant vector (all ones) such that the height is equal to the output height (i.e., 71). The feature vectors at each height level are then passed through an MLP with units [128, 256]. A Bidirectional LSTM layers with units [128, 256, 512] is applied. Finally, a dense layer is applied which maps the 71x512 hidden units to the 71x4 output features.

which means that for each field (temperature, cloud cover, etc.), only a single mean and standard deviation are computed. This approach is adopted because, at higher atmospheric levels, certain fields that should be zero end up numerically near, but not equal to zero. Normalizing across all heights prevents the unnatural scaling that would occur if they were normalized independently at each height.

For the custom loss function, the radiative heating rates for each atmospheric layer are computed by an ICON routine, using a finite-difference approximation of the vertical flux derivative. Obtaining accurate heating rates from the predicted fluxes is essential for updating the thermodynamic equation in ICON. Therefore, the mean squared error of the heating rates calculated from the emulated flux is added to the loss function of the NNs. The loss function thus consists of two parts: the mean squared error of the flux prediction and the mean squared error of the heating rates.

A challenge of the ML emulation of ecRad found during online inference is obtain-181 ing accurate heating rates at atmospheric levels above 35 km height. At those heights, 182 the air mass between two model levels is extremely small and small errors in the height 183 profile of the fluxes can result in large heating rate errors which can then break ICON 184 simulations. The proposed loss function already ensures improved heating rate estima-185 tions, but it turns out not to be enough to ensure the stability of the ICON model at 186 its top. For this reason, the computation of the fluxes above 15 km height is modified¹. 187 At those heights, a simplified formula is used to compute the heating rates. The fluxes 188 above the chosen level H = 30 (15 km height) are constructed by multiplying the fluxes 189 at level H by a set of constants $\beta_{H-1}, \ldots, \beta_0$, which are optimized based on the train-190 ing set: 191

$$f_{l,k} = \beta_k f_{l,H},\tag{1}$$

where $f_{l,k}$ is the chosen flux for the atmospheric column l at height level k and β_k are

different for short- and longwave up- and downward fluxes. See Appendix A for more details.



Figure 1: Schematic illustrating the integration of the ML emulator into ICON with Infero. Infero requires both the input and output data to be copied back and forth between GPU and CPU memory (steps 1, 3, 5 and 7), which substantially slows down the ICON simulation. First, the input data from ICON are copied from GPU to CPU memory (step 1). Infero obtains a CPU pointer to these ICON inputs (step 2). The "Tensorflow for C" backend then copies the inputs to the GPU memory (step 3), calls the NNs (or RF) for the fluxes predictions on GPUs (step 4) and finally copies back the outputs to the CPU memory (step 5). Finally, a CPU pointer to the Infero outputs is created (step 6) and the output data are copied back to GPU memory (step 7).

197

2.3 Implementation of the ML emulators into ICON

The implementation of the ML models with Python into the Fortran code of ICON 198 using Infero (Antonino Bonanni, 2022) is discussed here. Infero works with models built 199 with the Tensorflow Python library (Abadi et al., 2015) or with models in the ONNX for-200 mat (Bai et al., 2019). The RF is trained with the Scikit-learn (Pedregosa et al., 2011) 201 Python library and then ported to the ONNX format. The NNs are implemented with Ten-202 sorflow. In Figure 1, we show a schematic of how Infero integrates into ICON running 203 on GPUs. Infero requires the inputs to be in CPU (referred to as host) memory and af-204 ter an ML prediction on the GPUs (referred to as device), it returns the outputs there. 205 This is a limitation for a weather or climate model running on GPUs as it leads to back-206 and-forth copies of data between the separate memory spaces. Note however that for most 207 models, which are running solely on the central processing unit (CPU), this is not an is-208 sue and Infero then becomes a suitable coupler. Furthermore, Infero is developed by ECMWF, 209 which are increasing their efforts in using ML methods to improve current weather mod-210 els (ECMWF, 2023). Infero may thus see improvements in the future, allowing for an 211 optimized coupling on the GPUs. Hence, we still report next how Infero calls the NNs 212 when coupled to ICON, although an alternative coupler allowing direct access of the data 213 on the GPUs could be preferred for GPU-enabled weather and climate models. 214

¹ Note that the method we explain here does not appear in Bertoli et al. (2024) since the problem became apparent only after the emulator was integrated into ICON.



Figure 2: Schematic illustrating the integration of the ML emulator into ICON with PyTorch-Fortran. In contrast to Infero, no copies between CPU and GPU memory are required. First, a pointer to ICON inputs is created (step 1). Then, the NNs is called (step 2). Finally, a pointer to PyTorch-Fortran outputs is created so that ICON can handle the predicted fluxes (step 3).

In this paper, we experiment with the ICON version running fully on GPUs. The 215 data are first copied from the GPU to the CPU (step 1 in Figure 1). Infero can then ac-216 cess the input data (step 2) and compute the outputs with the NNs or RF (steps 3–5). 217 Internally, Infero uses a "Tensorflow for C" backend, which will copy the input data to 218 GPU memory (step 3), run the ML model on the GPUs (step 4) and then copy the out-219 put data back to CPU memory (step 5). Finally, the output data are accessed by ICON 220 (step 6) and copied back to GPU memory (step 7). The input and output data are thus 221 copied twice each. This drastically slows down the computation of the radiative fluxes 222 and diminishes the possible speed-up from using an ML model instead of ecRad. 223

An alternative to Infero is PyTorch-Fortran (Alexeev, 2022). Figure 2 shows how 224 the ML models are integrated into ICON with PyTorch-Fortran. PyTorch-Fortran re-225 quires the models to be built with the Pytorch library instead of Tensorflow. The main 226 advantage of Pytorch-Fortran is that both ML emulator and ICON can run purely on 227 GPUs. PyTorch-Fortran first obtains pointers to the input data, then runs the ML model 228 and return pointers to the output data to ICON. Therefore, it completely avoids copy-229 ing data between CPU and GPU memory. This is a major advantage compared to In-230 fero. A limitation is, however, that PyTorch-Fortran does not work with models built 231 in the ONNX format. This renders it incompatible with the Scikit-learn library which con-232 tains a variety of ML models. 233

234 3 Results

235

3.1 Offline performance

Figure 3 shows the mean absolute error (MAE) of the different ML models. The 236 RF is the most accurate above 15 km height for the heating rates and the downward fluxes. 237 It is however the least accurate model below 15 km height, where accuracy is the most 238 important. The RNN outperforms the UNet at all heights for both the fluxes and heat-239 ing rates. The effect of the simplified formula (1) used to compute the fluxes above level 240 30 is shown in the heating rates profile. This simplified equation improves the shortwave 241 heating rates accuracy above 35 km height and the longwave heating rates accuracy above 242 25 km for both NNs. On pressure coordinates, we observe that the increase in the heat-243 ing rates error from levels 30–25 to level 0 at the TOA, is concentrated in a small pres-244 sure interval from around level 100 hPa up to the model top. This compares well with 245 results of the literature shown in pressure coordinates (Chevallier et al., 1998; Ukkonen, 246 2022; Liu et al., 2020) or height coordinates with logarithmic scale (Lagerquist et al., 2021). 247



Figure 3: MAE over the test set for the RF, the UNet and the RNN models. Above level 30, the simplified Equation (1) is used to compute the fluxes for the UNet and the RNN. The MAEs of the RNN and UNet without the simplified Equation (1) is shown in black. In the last column, the heating rates MAE is shown in pressure coordinates.

248

3.2 Online performance

ICON simulations are performed over 3 weeks with the different ML emulators. Each simulation is restarted from the end of the 23rd month of the two-year simulation that produced the training and testing datasets. Recall that the ML models are trained with months 13 to 21. The experiments related to Figures 4 and 5 are performed on the Piz Daint HPC with NVIDIA Tesla P100 16GB GPUs (NVIDIA, 2016), while the runtime presented in Table 2 are performed on the Balfrin CSCS machine with NVIDIA A100 Tensor Core GPU (NVIDIA, 2022).

In Figure 4, the global mean temperatures at the TOA, at 1 km height and at the surface from ICON simulations with different radiation emulators are compared to the reference simulation with ecRad. At the TOA, the RNN model provides by far the highest accuracy. For the UNet model, the global temperature is dropping fast to below 135 K at the end of the three-weeks-long simulation. We extended the simulation with the UNet to two months and although the global temperature continues to drop at the TOA, this



Figure 4: Global mean temperature in Kelvin at the TOA, at 1 km height and at the surface for three-weeks-long ICON simulations with four different radiation parameterizations: RF, RNN, UNet and ecRad. The horizontal dotted line corresponds to the ecRad global mean temperature at time 0. Each vertical axis is centered on the ecRad global mean temperature at time 0 and has a 2 K range of values. At the TOA, the RNNs is the only ML model that provides an accurate global mean temperature with respect to ecRad.

does not seem to affect ICON's stability. Furthermore, no perturbations of the lower lev-262 els due to the decrease of temperature at the top is observed. The temperature of the 263 RF simulation is significantly more accurate than the UNet simulation although the temperature is increasing over time at the TOA. It is however the least accurate model at 265 1 km height and at the surface during the first ten days where accuracy matters the most. 266 Extensive online simulation horizons with RNN are still competitive with the ecRad ref-267 erence. Throughout a time horizon of two months, the temperature difference between RNN and ecRad simulations never exceed 2 K at TOA and 0.3 K at surface and 1 km height 269 levels. The superior stability of the RNN makes it a good candidate for multi-years ICON 270 climate simulations. 271

Figure 5 presents meridional vertical cross sections of heating rates along the prime meridian 2 d, 4 d, and 6 d into the simulation using ecRad and the RNN model. Both simulations are started from the same restart file after a 23 months of simulation with ecRad. This is possible because our implementation allows us to switch between ecRad and the RNN model during runtime. After two days, both parameterizations exhibit nearly iden-



Figure 5: Meridional vertical cross-sections along the prime meridian of net heating rates in pressure coordinates as computed by ICON with (left) ecRad and (b) the RNN radiation emulator. The rows show (top to bottom) instantaneous fields $(a, b) \ 2d$, $(c, d) \ 4d$ and $(e, f) \ 6d$ into the simulation, as well as (g, h) the mean over the first 10 d. The bottom row shows a 10-day average. Both simulations have been started from the same restart file saved after 23 months of simulation with ecRad.

277 tical heating rate profiles. By day 4, differences emerge near the equator. By day 6, they have become more pronounced above 600 hPa, which is expected as small differences be-278 tween the simulations start to grow over time, similar as in ensemble runs. While the 279 mean heating rates align in both simulations, a small discrepancy arises at 5° S where 280 positive heating rates between 400 hPa and 100 hPa are underestimated on average in 281 the RNN-based simulation. This is evident in the instantaneous heating rate displayed 282 for 4 days and 6 days into the simulations in Figure 5. At this stage it is unclear whether 283 this indicates a systematic emulation bias, deviating model trajectories due to growing 284 perturbations (akin to different ensemble members) or a combination of both. 285

In Tab. 2, the run times of the radiation and the whole physics are shown for 400 286 ICON time-steps, corresponding to 20 hours of simulated time. The RNN model is twice 287 as slow as ecRad and requires more than 8 GPUs (Nvidia A100 tensor core GPU with 288 80 GB of memory) to run. In comparison, the slim UNet model is as fast as ecRad with 289 12 GPUs, twice as fast with 8 GPUs and three times as fast with 4 GPUs. It can even 290 fit on 2 GPUs in contrast to ecRad. As an alternative to the RNN model, a smaller model 291 is trained, with similar MAE than the RNN model described in Table 1, with fewer LSTM 292 layers and a single set of weights for the first MLPs of size [128, 256]. This tuning re-293 duces the number of trainable parameters from 15 million to 5 million and the size of 294

Table 2: Comparison of run time in seconds for 400 ICON time-steps using ecRad and the RNN and the UNet emulators as a function of the number of GPU nodes (Nvidia A100 tensor core GPU with 80 GB of memory). Out-of-memory (OOM) issues are indicated. Shown in brackets are runtimes of the all physics parameterization combined.

No. GPUs	2	4	8	12
UNet	92 [116]	50 [65]	35 [49]	46 [57]
RNN	OOM	OOM	OOM	108 [117
RNN optimized	OOM	143 [156]	77 [78]	48 [59]
ecRad	OOM	144 [184]	73 [99]	51 [73]

the RNN from 64 MB to 22 MB. The smaller model fits on 4 GPUs and the total run 295 time of all physics parameterizations is now below that including ecRad, and the radi-296 ation emulation itself is now comparably fast between the two, with variations between 297 the number of GPUs used. Note that, for the ICON simulation corresponding to the run-298 time given in Table 2, the radiation parameterization is called every time-step and is solved at full resolution. The total runtime of all physics is thus dominated by the radiation 300 runtime. The runtime ratio between ecRad and the machine learning models for the ra-301 diative process aligns with a similar ratio for the sum of all physical phenomena as seen 302 in Table 2. The ratio for ecRad and the full physics may differ slightly due to due dif-303 ferent workload on the Balfrin system used for these experiments. This experiment shows 304 that ML-based radiation emulators running on GPUs are not per se faster than a highly 305 optimized physics-based GPU-enabled solver like ecRad, which is written in Fortran with 306 OpenACC compiler directives. 307

ML models size could be reduced further with, for example, automatic mixed pre-308 cision (Carilli, 2024), where some operations are done with half precision instead of full 309 precision, and dynamic quantization (Dynamic Quantization, 2024), which reduces the 310 resolution of the model's weight. ICON and PyTorch are written in two different lan-311 guages (Fortran with OpenACC and C++ with CUDA) which access the same GPU mem-312 ory space when coupled together. It is yet unclear how both languages interact regard-313 ing data access and further profiling would be required to optimize how PyTorch should 314 share the data access with the rest of the ICON code. This is however beyond the scope 315 of this study. By comparison, ICON and ecRad are both written in Fortran with Ope-316 nACC directives. Note also that in (Ukkonen & Hogan, 2024), the authors restructured 317 the ecRad code and improved its runtime performance by a factor of up to 12. The code's 318 restructuring is designed for a CPU usage but the improved parallelism could benefit a 319 GPU implementation of ecRad. Depending on its performance on GPUs, this optimized 320 version of ecRad could outperform the ML models presented here. 321

322 4 Summary

In this two-part study, an ML emulation of the ecRad radiative transfer param-323 eterization is built for the GPU-enabled ICON weather and climate model. In Part I, 324 through a series of offline tests the most accurate ML model has been identified as a bidi-325 rectional recurrent NNs with long-short memory layers and additional physics-informed 326 normalization of input and output features, as well as an additional heating rate related 327 loss term in the objective function (Bertoli et al., 2024). In this work, Part II, a signif-328 icant technical advancement is made by integrating the ML radiative transfer param-329 eterization into ICON. 330

Since the air in the upper atmospheric layers is substantially less dense than lower 331 layers, a small error in the flux profile results in a large heating rate error in the upper 332 levels, which during long integration can cause spurious temperature trends near the model 333 top. Model tops are known to often require additional tuning not seen during offline train-334 ing. For example, Brenowitz and Bretherton (2019) removed the model top from train-335 ing to stabilise the online performance of their ML-based convection scheme. To mit-336 igate the aforementioned problem with high heating rates, which is easily overlooked in 337 the offline testing, a simplified formula is used to calculate the fluxes in the damping layer 338 of ICON near the model top, which reduces the heating rate error significantly and keeps 339 the model free from any temperature drifts and very close to the original ecRad simu-340 lation for several simulated weeks. However, future research is needed to further improve 341 the emulation at the top of the model, in particular for shortwave radiation, and to in-342 crease the reliability of the results at scales beyond weather forecasting. A potential way 343 forward is to train an emulator that infers the TOA and surface level shortwave and long-344 wave fluxes plus the heating rates on all levels within the atmosphere. 345

To seamlessly connect the Fortran code with the NNs implemented with Python, 346 the Infero coupler from ECMWF was explored initially (Section 2.3). Infero requires that 347 the ML model inputs and outputs are provided in the CPU memory. However, the ver-348 sion of ICON examined in this paper operates entirely on GPUs. Consequently, using 349 Infero leads to needless data transfers between CPU and GPU memory, causing notable 350 delays in the ML parameterization compared to ecRad. Note however that the next gen-351 eration of hardware, like the Grace Hopper chips (NVIDIA, 2023) chosen for the next 352 CSCS supercomputer Alps (CSCS et al., 2021), reduces the overhead of CPU-GPU copies 353 and can even expose a shared CPU-GPU memory space to the user. This could make 354 Infero more competitive, even in its current form. There appear to be no obstacles pre-355 venting the adaptation of Infero for complete GPU utilization, thereby eliminating the 356 need for data transfers between hardware components. As such, the limitation of Infero 357 exposed in this paper could potentially be nonexistent in future versions of this software. 358 In this paper, to avoid CPU-GPU copies, the PyTorch-Fortran library (Alexeev, 2022) 359 is adopted as an alternative solution. This approach enables direct processing of ML in-360 puts and outputs on the GPU. The integrated system of GPU-enabled ICON, by means 361 of OpenACC, and ML emulator implemented with PyTorch is deployed on the Piz Daint 362 and Balfrin systems at CSCS, leveraging GPU capabilities for enhanced performance. 363

To the best of our knowledge, this is the first time that a full-fledged weather and 364 climate model in combination with an ML-based parameterization developed in PyTorch 365 has been run completely on GPUs. The performance gain compared with simulations 366 using the original ecRad radiation solver depends critically on the complexity of the NNs 367 architecture, and not all tested NNs are per se faster than the traditional physics-based 368 code. Performance gains reported in past studies may stem from the fact that the em-369 ulated parameterisation was originally run on CPUs while the ML emulator was run on 370 GPUs. Also in terms of memory consumption, we find that the memory footprint of ecRad 371 is smaller compared to the RNN, albeit larger compared to the UNet architecture. We 372 cannot rule out the possibility that a more sophisticated tuning of the NNs architectures 373 would result in a higher speed-up, but this holds also true for the original radiation solver 374 (Ukkonen & Hogan, 2024). 375

³⁷⁶ Open Research Section

The data were generated using the ICON climate model described in Prill et al. (2023). The software is available at https://www.icon-model.org/. The codes to reproduce the results of this paper will be made available in https://gitlab.renkulab .io/deepcloud/rfe. Data to reproduce results of this work will be hosted at ETH Research Collection https://www.research-collection.ethz.ch/ (with a DOI) together with the ICON runscript used to generate the full dataset. ETH Zurich's Research-Collection adheres to the FAIR principles and data is stored for at least 10 years.

384 Acknowledgments

This work was funded through a grant by the Swiss Data Science Center (SDSC grant C20-03). This research was supported by computation resources provided by the EXCLAIM project funded by ETH Zurich (CSCS project number d121). The Center for Climate Systems Modeling (C2SM) at ETH Zurich is acknowledged for providing technical and scientific support. Sebastian Schemm and Stefan Rüdisühli are supported by the European Research Council, H2020 European Research Council (grant no. 848698).

391 **References**

410

411

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. 392 Tensorflow: Large-scale machine learning on heterogeneous systems. (2015).393 Retrieved from https://www.tensorflow.org/ (Software available from 394 tensorflow.org) 395 Alexeev, D. (2022). pytorch-fortran v0.4. https://github.com/alexeedm/pytorch 396 -fortran. GitHub. 397 Antonino Bonanni, T. Q., James Hawkes. (2022). Infero 0.1.0. https://github 398
- Antonino Bonanni, T. Q., James Hawkes. (2022). Infero 0.1.0. https://github .com/ecmwf/infero. GitHub.
- Bai, J., Lu, F., Zhang, K., et al. (2019). ONNX: Open neural network exchange.
 https://github.com/onnx/onnx. GitHub.
- Bertoli, G., Ozdemir, F., Perez-Cruz, F., & Schemm, S. (2024). Revisiting machine
 learning approaches for short- and longwave radiation inference in weather and
 climate models, part I: Offline performance. Journal of Advances in Modeling
 Earth Systems.
- Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by coarse-graining. Journal of Advances in Modeling Earth Systems, 11(8), 2728-2744. Retrieved from https://doi.org/ 10.1029/2019ms001711 doi: 10.1029/2019ms001711
 - Carilli, M. (2024). Automatic mixed precision. https://pytorch.org/tutorials/ recipes/recipes/amp_recipe.html. PyTorch.
- 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(11), 1385–1397. Retrieved from
 https://doi.org/10.1175/1520-0450(1998)037<1385:annafa>2.0.co;2
 doi: 10.1175/1520-0450(1998)037<1385:annafa>2.0.co;2
- CSCS, NVIDIA, & Entreprise, H. P. (2021). Swiss national supercomput ing centre, hewlett packard enterprise and nvidia announce world's most
 powerful ai-capable supercomputer. https://www.cscs.ch/science/
- computer-science-hpc/2021/cscs-hewlett-packard-enterprise-and
 -nvidia-announce-worlds-most-powerful-ai-capable-supercomputer.
- 422Dynamic quantization.(2024).https://pytorch.org/tutorials/recipes/423recipes/dynamic_quantization.html.PyTorch.
- ECMWF. (2023). Annual report 2022. https://www.ecmwf.int/sites/default/ files/elibrary/2023/81359-annual-report-2022.pdf.
- Giorgetta, M. A., Sawyer, W., Lapillonne, X., Adamidis, P., Alexeev, D., Clément,
 V., ... Stevens, B. (2022). The ICON-A model for direct QBO simulations on
 GPUs (version icon-cscs:baf28a514). Geoscientific Model Development, 15(18),
 6985–7016. Retrieved from https://gmd.copernicus.org/articles/15/
 6985/2022/ doi: 10.5194/gmd-15-6985-2022
- 431 Kedward, L. J., Aradi, B., Certik, O., Curcic, M., Ehlert, S., Engel, P., ... Van-
- denplas, J. (2022). The state of fortran. Computing in Science I& En-

 MCSE.2022.3159862 doi: 10.1109/mcse.2022.3159862 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). proach to calculation of atmospheric model physics: Accurate and f network emulation of longwave radiation in a climate model. Mont Review, 133(5), 1370–1383. Retrieved from https://doi.org mwr.2923.1 doi: 10.1175/mwr.2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. 1 ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mod ment, 13(9), 4399–4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.co ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.m content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nvd -tesla-p100-datasheet.pdf. Author. NVIDIA. (2023). NVIDIA A100 Tensor Core GPU. https://www.m content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvdia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:/ .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM raa transfer. Geophysical Research Letters, 46(11), 6069-6079. Ref https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z.,	005). New ap- and fast neural Monthly Weather i.org/10.1175/ V. (2021). Us- model. Journal s://doi.org/10 pring deep learn- Model Develop- 10.5194/gmd-13 ul nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor.
 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). proach to calculation of atmospheric model physics: Accurate and f network emulation of longwave radiation in a climate model. Mont Review, 133(5), 1370–1383. Retrieved from https://doi.or. mwr2923.1 doi: 10.1175/mwr2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mo ment, 13(9), 4399–4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.cc ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.m content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nvidia-z -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.m content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-z -datasheet-nvidia-us-218804-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	005). New ap- and fast neural <i>Monthly Weather</i> i.org/10.1175/ V. (2021). Us- model. <i>Journal</i> s://doi.org/10 oring deep learn- c <i>Model Develop</i> - 10.5194/gmd-13 <i>ul nonhydrostatic</i> lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 proach to calculation of atmospheric model physics: Accurate and f network emulation of longwave radiation in a climate model. Mont Review, 133(5), 1370–1383. Retrieved from https://doi.or mwr2923.1 doi: 10.1175/mwr2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Occanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mod ment, 13(9), 4399–4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.cc ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nvidia-c -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-c -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon Ala, Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., <!--</td--><td>and fast neural Monthly Weather i.org/10.1175/ V. (2021). Us- model. Journal s://doi.org/10 oring deep learn- e Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks</td>	and fast neural Monthly Weather i.org/10.1175/ V. (2021). Us- model. Journal s://doi.org/10 oring deep learn- e Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 network emulation of longwave radiation in a climate model. Mont Review, 133(5), 1370-1383. Retrieved from https://doi.or. mwr2923.1 doi: 10.1175/mwr2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mon ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.cog NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.nr content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.nr content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia- datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author. OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	Monthly Weather i.org/10.1175/ V. (2021). Us- model. Journal s://doi.org/10 pring deep learn- c Model Develop- 10.5194/gmd-13 al nonhydrostatic hel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 Review, 133(5), 1370–1383. Retrieved from https://doi.or/mwr2923.1 doi: 10.1175/mwr2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (1) ing deep learning to emulate and accelerate a radiative-transfer modor of Atmospheric and Oceanic Technology. Retrieved from https://.1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Moment, 13(9), 4399–4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.cc ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.nt content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.nt content/dam/en-zz/Solutions/Data-Center/sla-p100/pdf/nvida-ac-datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM radatasheet. Geophysical Research Letters, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081646 	<pre>i.org/10.1175/ V. (2021). Us- model. Journal s://doi.org/10 oring deep learn- e Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks of rodiative</pre>
 mwr2923.1 doi: 10.1175/mwr2923.1 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mon ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nor weather and climate model. Retrieved from https://icon-model.co ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.nt content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.nt content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authot OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad- transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	V. (2021). Us- rmodel. Journal s://doi.org/10 pring deep learn- <i>Model Develop</i> - 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor.
 Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mod ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.org NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.m content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nvidia- datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ref Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	V. (2021). Us- model. Journal s://doi.org/10 pring deep learn- <i>Model Develop</i> - 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor.
 ing deep learning to emulate and accelerate a radiative-transfer mod of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mod ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.cog ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.nr content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.nr content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ref https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	<pre>model. Journal s://doi.org/10 oring deep learn- c Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks of radiative</pre>
 of Atmospheric and Oceanic Technology. Retrieved from https:// .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mon ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral nov weather and climate model. Retrieved from https://icon-model.cc ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.m content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.m content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	s://doi.org/10 pring deep learn- c Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 .1175/jtech-d-21-0007.1 doi: 10.1175/jtech-d-21-0007.1 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mon ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.cc ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.m content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv retsla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.m content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	pring deep learn- c Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 Liu, Y., Caballero, R., & Monteiro, J. M. (2020). RadNet 1.0: exploring ing architectures for longwave radiative transfer. Geoscientific Mon ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.co ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.nt content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.nt content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	pring deep learn- c Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ df/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 ing architectures for longwave radiative transfer. Geoscientific Mod ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.c ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM radiation for super-parameterized E3SM radiation	c Model Develop- 10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 ment, 13(9), 4399-4412. Retrieved from https://doi.org/10.5 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.cg NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	10.5194/gmd-13 al nonhydrostatic lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 -4399-2020 doi: 10.5194/gmd-13-4399-2020 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.c ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho <i>OpenACC, version 3.2.</i> (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	al nonhydrostatic del.org/ (Soft- ww.nvidia.com/ df/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 MPI-M, DWD, KIT, DKRZ, & C2SM. (2024). ICON: Icosahedral no weather and climate model. Retrieved from https://icon-model.c ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho <i>OpenACC, version 3.2.</i> (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	al nonhydrostatic del.org/ (Soft- ww.nvidia.com/ df/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 weather and climate model. Retrieved from https://icon-model.c ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon <i>OpenACC</i>, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	<pre>lel.org/ (Soft- ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources .uthor. neural networks wf radiative</pre>
 ware available at https://icon-model.org) NVIDIA. (2016). NVIDIA Tesla P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho <i>OpenACC</i>, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ref https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ww.nvidia.com/ lf/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 NVIDIA. (2016). NVIDIA Testa P100 GPU Accelerator. https://www.n content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author <i>OpenACC, version 3.2.</i> (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ww.nvidia.com/ hf/nvidia ww.nvidia.com/ hia-a100 ps://resources uthor. neural networks
 content/dam/en-zz/Solutions/Data-Center/tesla-p100/pdf/nv -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon <i>OpenACC, version 3.2.</i> (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	if/nvidia ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 -tesla-p100-datasheet.pdf. Author. NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho <i>OpenACC</i>, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ww.nvidia.com/ lia-a100 ps://resources .uthor. neural networks
 NVIDIA. (2022). NVIDIA A100 Tensor Core GPU. https://www.n content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a -datasheet-nvidia-us-2188504-web.pdf. Author. NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur as cost-effective surrogate models for super-parameterized E3SM rad transfer. Geophysical Research Letters, 46(11), 6069-6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ww.nvidia.com/ lia-a100 ps://resources uthor. neural networks
 ⁴⁵⁵ Content/dam/en-22/Solutions/Data-Center/aloo/pdf/hVidia-a ⁴⁵⁶ -datasheet-nvidia-us-2188504-web.pdf. Author. ⁴⁵⁷ NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// ⁴⁵⁸ .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Author ⁴⁵⁹ OpenACC, version 3.2. (2022). https://www.openacc.org/. ⁴⁶⁰ Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur ⁴⁶¹ as cost-effective surrogate models for super-parameterized E3SM rad ⁴⁶² transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret ⁴⁶³ https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 ⁴⁶⁴ Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ps://resources .uthor. neural networks
 ⁴⁵⁶ NVIDIA. (2023). NVIDIA GH200 Grace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Authon ⁴⁵⁷ OpenACC, version 3.2. (2022). https://www.openacc.org/. ⁴⁶⁰ Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur ⁴⁶¹ as cost-effective surrogate models for super-parameterized E3SM rad ⁴⁶² transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret ⁴⁶³ https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 ⁴⁶⁴ Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	ps://resources .uthor. neural networks
 ⁴⁵⁷ NVIDIA. (2023). NVIDIA GI200 Glace Hopper Superchip. https:// .nvidia.com/en-us-grace-cpu/grace-hopper-superchip. Autho ⁴⁵⁸ OpenACC, version 3.2. (2022). https://www.openacc.org/. ⁴⁶⁰ Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neur ⁴⁶¹ as cost-effective surrogate models for super-parameterized E3SM rad ⁴⁶² transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret ⁴⁶³ https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 ⁴⁶⁴ Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	neural networks
 OpenACC, version 3.2. (2022). https://www.openacc.org/. Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM radius transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	neural networks
 ⁴⁵⁹ OpenACC, version 3.2. (2022). https://www.openacc.org/. ⁴⁶⁰ Pal, A., Mahajan, S., & Norman, M. R. (2019). Using deep neural as cost-effective surrogate models for super-parameterized E3SM rate transfer. Geophysical Research Letters, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 ⁴⁶⁴ Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	neural networks
 as cost-effective surrogate models for super-parameterized E3SM rad transfer. <i>Geophysical Research Letters</i>, 46(11), 6069–6079. Ret https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	Incurat incovorks
 462 transfer. <i>Geophysical Research Letters</i>, 46(11), 6069–6079. Ret 463 https://doi.org/10.1029/2018g1081646 doi: 10.1029/2018g1081 464 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., 	vi tamalive
463 https://doi.org/10.1029/2018gl081646 doi: 10.1029/2018gl081 464 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z.,	Retrieved from
Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z.,	g1081646
	Lerer. A.
465 (2017). Automatic differentiation in pytorch. In Nips-w.	, ,
Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Gris	Grisel, O.,
467 Duchesnay, E. (2011). Scikit-learn: Machine learning in Python.	non. Journal of
468 Machine Learning Research, 12, 2825–2830.	•
469 Petty, G. W. (2006). A first course in atmospheric radiation. Madison,	lison, Wisconsin:
470 Sundog Publishing.	
471 Prill, F., Reinert, D., Rieger, D., & Zängl, G. (2023). ICON tu	N tutorial 2023:
472 Working with the ICON model. Deutscher Wetterdienst. Ret	Retrieved from
473 https://www.dwd.de/EN/ourservices/nwv_icon_tutorial/pdf_vo	J£
474 icon_tutorial2023_en.pdf?blob=publicationFile&v=3 de	ar_vorume/
475 DWD_PUB/NWV/ICON_TUTORIAL2023	doi: 10.5676/
⁴⁷⁶ Roh, S., & Song, HJ. (2020). Evaluation of neural network emulation	doi: 10.5676/
	doi: 10.5676/
477 tion parameterization in cloud resolving model. <i>Geophysical Research</i>	doi: 10.5676/ lations for radia- Research Letters,
 tion parameterization in cloud resolving model. <i>Geophysical Resea</i> 477 477 477 47(21). Retrieved from https://doi.org/10.1029/2020gl08944 1020 1020 1020 1020 1044 	doi: 10.5676/ lations for radia- Research Letters, 89444 doi: 10
 tion parameterization in cloud resolving model. Geophysical Resea 477 tion parameterization in cloud resolving model. Geophysical Resea 478 479 47(21). Retrieved from https://doi.org/10.1029/2020gl089444 479 .1029/2020gl089444 479 Line parameterization in cloud resolving model. Geophysical Resea 	doi: 10.5676/ lations for radia- Research Letters, 89444 doi: 10
 tion parameterization in cloud resolving model. Geophysical Resea 477 tion parameterization in cloud resolving model. Geophysical Resea 478 47(21). Retrieved from https://doi.org/10.1029/2020gl089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear tion of atmospheric redicting transfer 	doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10
 tion parameterization in cloud resolving model. Geophysical Resea 477 tion parameterization in cloud resolving model. Geophysical Resea 478 47(21). Retrieved from https://doi.org/10.1029/2020gl089444 479 .1029/2020gl089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear 481 tion of atmospheric radiative transfer. Journal of Advances in Mode 481 Sautema 1/(4) Patriaved from https://doi.org/10.1029/2020/2020 	doi: 10.5676/ doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10 e learning emula- <i>Modeling Earth</i>
 tion parameterization in cloud resolving model. Geophysical Reseat 477 tion parameterization in cloud resolving model. Geophysical Reseat 478 47(21). Retrieved from https://doi.org/10.1029/2020gl089444 479 .1029/2020gl089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear 481 tion of atmospheric radiative transfer. Journal of Advances in Mod 482 Systems, 14 (4). Retrieved from https://doi.org/10.1029/202 	di_volume/ doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10 e learning emula- b <i>Modeling Earth</i> 9/2021ms002875
 tion parameterization in cloud resolving model. Geophysical Reseat 477 tion parameterization in cloud resolving model. Geophysical Reseat 478 479 470(21). Retrieved from https://doi.org/10.1029/2020g1089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear 481 tion of atmospheric radiative transfer. Journal of Advances in Mod 482 Systems, 14 (4). Retrieved from https://doi.org/10.1029/200 483 doi: 10.1029/2021ms002875 484 Ukkonen, P. & Hogan, R. L. (2024). Twelve times factor wet accurate 	di_volume/ doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10 e learning emula- <i>Modeling Earth</i> 9/2021ms002875
 tion parameterization in cloud resolving model. Geophysical Reseat 477 tion parameterization in cloud resolving model. Geophysical Reseat 478 479 .1029/2020gl089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear 481 tion of atmospheric radiative transfer. Journal of Advances in Mod 482 Systems, 14 (4). Retrieved from https://doi.org/10.1029/202 483 doi: 10.1029/2021ms002875 484 Ukkonen, P., & Hogan, R. J. (2024). Twelve times faster yet accurate 	doi: 10.5676/ doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10 e learning emula- b <i>Modeling Earth</i> 9/2021ms002875 accurate: A new ral optimiza-
 tion parameterization in cloud resolving model. Geophysical Reseat 477 tion parameterization in cloud resolving model. Geophysical Reseat 478 47(21). Retrieved from https://doi.org/10.1029/2020gl089444 479 .1029/2020gl089444 480 Ukkonen, P. (2022). Exploring pathways to more accurate machine lear 481 tion of atmospheric radiative transfer. Journal of Advances in Model 482 Systems, 14 (4). Retrieved from https://doi.org/10.1029/202 483 doi: 10.1029/2021ms002875 484 Ukkonen, P., & Hogan, R. J. (2024). Twelve times faster yet accurate 485 state-of-the-art in radiation schemes via performance and spectral of 486 tion. Journal of Advances in Modeling Earth Systems, 16(1) Ret 	doi: 10.5676/ doi: 10.5676/ lations for radia- <i>Research Letters</i> , 89444 doi: 10 e learning emula- <i>Modeling Earth</i> 9/2021ms002875 accurate: A new ral optimiza- Retrieved from

- Ukkonen, P., Pincus, R., Hogan, R. J., Nielsen, K. P., & Kaas, E. (2020).Accel-488 erating radiation computations for dynamical models with targeted machine 489 learning and code optimization. Journal of Advances in Modeling Earth Sys-490 tems, 12(12). Retrieved from https://doi.org/10.1029/2020ms002226 doi: 491 10.1029/2020 ms 002226492 Veerman, M. A., Pincus, R., Stoffer, R., van Leeuwen, C. M., Podareanu, D., & 493 (2021).van Heerwaarden, C. C. Predicting atmospheric optical properties 494
- for radiative transfer computations using neural networks. *Philosophical*
- 496 Transactions of the Royal Society A: Mathematical, Physical and Engineering
- 497 Sciences, 379(2194), 20200095. Retrieved from https://doi.org/10.1098/
- 498 rsta.2020.0095 doi: 10.1098/rsta.2020.0095

Appendix A Modification of the fluxes computation at the upper levels

Four different sets of constants β_k corresponding to the shortwave down, shortwave up, longwave down and longwave up fluxes are constructed. The constants β_k , $k = 0, \ldots, H-1$, are given by

$$\beta_k = \frac{1}{N} \sum_{l=1}^{N} f_{k,l} / f_{H,l},$$
(A1)

where the mean is taken over all atmospheric columns in the training set. The top at-501 mospheric layers in ICON are Rayleigh damping layers, whose purpose is to attenuate 502 oscillations reaching the top boundary. ICON is hence not designed to be accurate at 503 such heights. In particular the ML emulator does not need to emulate perfectly ecRad 504 in the damping layers. The goal of this strategy is to thus sacrifice flux accuracy at the 505 top height levels for an increase in the stability of ICON. This strategy is not used on 506 the random forest, which already provides accurate heating rates at the top levels. This 507 is to be expected because the random forest prediction is average of flux profiles encoun-508 tered in the training set. 509

For the upper levels, where the fluxes are approximated with Formula 1, the heating rates $\partial_t T_k^{rad}$ at level k for k in $0, \ldots, 29$, are given by the following formula:

$$\partial_t T_k^{rad} = C_k (f_{k-1} - f_k)$$

= $C_k f_{30} (\beta_{k-1} - \beta_k)$
= $C_k f_{30} \left(\frac{1}{N} \sum_{l=1}^N \frac{f_{k-1,l}}{f_{30,l}} - \frac{1}{N} \sum_{l=1}^N \frac{f_{k,l}}{f_{30,l}} \right)$
= $\frac{1}{N} \sum_{l=1}^N \frac{f_{30}}{f_{30,l}} C_k (f_{k-1,l} - f_{k,l})$
= $\frac{1}{N} \sum_{l=1}^N \frac{f_{30}}{f_{30,l}} \frac{C_k}{C_{l,k}} \partial_t T_{k,l}^{rad},$ (A2)

where C_k represent the effect of the air mass and humidity and is given by

$$C_k = \frac{1}{m_k(c_d + (c_v - c_d)q_k)},$$

and where m_k an q_k are the air mass and specific humidity at height level k and c_v , c_d 510 are the specific heat of water vapor and dry air at constant volume, assumed constant 511 in ICON. Equation A2 shows that the heating rates obtained from the predicted fluxes 512 are then the approximate weighted mean of heating rates observed during training. It is only an approximate weighted mean because $\frac{1}{N}\sum_{l=1}^{N}\frac{f_{30}}{f_{30,l}}\frac{m_{k,l}(c_d+(c_v-c_d)q_{k,l})}{m_k(c_d+(c_v-c_d)q_k)}$ is not equal to 1 in general. Furthermore, no vertical derivative of the predicted fluxes appear in Equa-513 514 515 tion A2. Additionally, the inverse of the air mass $1/m_k$ is multiplied by the air mass $m_{k,l}$ 516 of observed atmospheric columns during training. This reduces the sensitivity of the heat-517 ing rates to the fluxes prediction. 518

Infero schematic.



PyTorch-Fortran schematic.



MAE vs height.



Shortwave





Mean temperature.

a) Top of atmosphere



heating rates prediction.



Latitude

