# Exploring pathways to more accurate machine learning emulation of atmospheric radiative transfer

Peter Ukkonen<sup>1</sup>

<sup>1</sup>Danish Meteorological Institute

November 23, 2022

#### Abstract

Machine learning (ML) parameterizations of subgrid physics is a growing research area. A key question is whether traditional ML methods such as feed-forward neural networks (FNNs) are better suited for representing only specific processes. Radiation schemes are an interesting example, because they compute radiative flows through the atmosphere using well-established physical equations. The sequential aspect of the problem implies that FNNs may not be well-suited for it. This study explores whether emulating the entire radiation scheme is more difficult than its components without vertical dependencies. FNNs were trained to replace a shortwave radiation scheme, its gas optics component, and its reflectance-transmittance computations. In addition, a novel recurrent NN (RNN) method was developed to structurally incorporate the vertical dependence and sequential nature of radiation computations. It is found that a bidirectional RNN with an order of magnitude fewer model parameters than FNN is considerably more accurate, while offering a smaller but still significant 4-fold speedup over the original code on CPUs, and a much greater speedup on GPUs. The RNN predicts fluxes with less than  $1\$  error, and heating rates computed from fluxes have a root-mean-square-error of 0.16 K days<sup>{-1}</sup> in offline tests using a year of global data. Finally, FNNs emulating gas optics are very accurate while being several times faster. As with RNNs emulating radiative transfer, the smaller dimensionality may be crucial for developing models that are general enough to be used as parameterizations.

## Exploring pathways to more accurate machine learning emulation of atmospheric radiative transfer

### Peter Ukkonen $^{1,2}$

 $^1\mathrm{Danish}$  Meteorological Institute $^2\mathrm{Niels}$ Bohr Institute, University of Copenhagen

#### Key Points:

1

2

3

4 5

6

7	• Feed-forward and recurrent neural networks (NN) were developed to emulate a
8	shortwave radiation scheme, as well as its components
9	• The recurrent NN has far better accuracy than usual approaches, while offering
10	a significant speedup especially on GPUs

• Using NNs for gas optics is 3 times faster and does not sacrifice accuracy

 $Corresponding \ author: \ Peter \ Ukkonen, \ \texttt{peter.ukkonen@nbi.ku.dk}$ 

#### 12 Abstract

Machine learning (ML) parameterizations of subgrid physics is a growing research area.
 A key question is whether traditional ML methods such as feed-forward neural networks
 (FNNs) are better suited for representing only specific processes. Radiation schemes are
 an interesting example, because they compute radiative flows through the atmosphere
 using well-established physical equations. The sequential aspect of the problem implies
 that FNNs may not be well-suited for it.

This study explores whether emulating the entire radiation scheme is more difficult than its components without vertical dependencies. FNNs were trained to replace a shortwave radiation scheme, its gas optics component, and its reflectance-transmittance computations. In addition, a novel recurrent NN (RNN) method was developed to structurally incorporate the vertical dependence and sequential nature of radiation computations.

It is found that a bidirectional RNN with an order of magnitude fewer model pa-25 rameters than FNN is considerably more accurate, while offering a smaller but still sig-26 nificant 4-fold speedup over the original code on CPUs, and a much greater speedup on 27 GPUs. The RNN predicts fluxes with less than 1% error, and heating rates computed 28 from fluxes have a root-mean-square-error of  $0.16 \text{ K day}^{-1}$  in offline tests using a year 29 of global data. Finally, FNNs emulating gas optics are very accurate while being sev-30 eral times faster. As with RNNs emulating radiative transfer, the smaller dimensional-31 ity may be crucial for developing models that are general enough to be used as param-32 eterizations. 33

#### <sup>34</sup> Plain Language Summary

Numerical weather and climate simulations are being performed at increasingly resolution, making the energy cost of simulations significant. Computing how solar and terrestrial radiation interact with Earth's atmosphere, surface, and clouds is one of the most computationally expensive parts in climate models especially. This has invited efforts to replace these computations with predictions from a neural network, which is approximative but considerably faster than physical radiation computations.

In this paper, different ways of emulating a radiation code with neural networks have been explored. Its main contribution is developing a novel emulation method based on recurrent neural networks, which more closely resemble physical radiative transfer computations. The accuracy is found to be considerably higher than with traditional neural network approaches which use an order of magnitude more model parameters.

#### 46 **1** Introduction

<sup>47</sup> Climate and weather simulations are being performed at increasingly high resolu<sup>48</sup> tions. The implications for energy use are significant: even with an atmospheric model
<sup>49</sup> fully ported to a state-of-the-art GPU supercomputer, kilometer-scale global simulations
<sup>50</sup> consume 596 MWh energy per simulated year (Fuhrer et al., 2018). This is the same as
<sup>51</sup> the yearly electricity consumption of 161 average EU households in 2018 (Odyssee-Mure,
<sup>52</sup> 2021). For the energy costs of earth system simulations not to become untenable, both
<sup>53</sup> hardware and and algorithmic improvements are needed.

An algorithmic development which could improve both the accuracy and computational efficiency of weather and climate simulations is the use of machine learning (ML) methods to represent sub-grid diabatic processes. Recent years has seen an influx of papers on this topic where the typical approach has been training neural networks (NNs) or random forests on coarse-grained data from cloud-resolving or high-resolution simulations, and representing all sub-grid processes with a single model (Rasp et al., 2018;

Brenowitz & Bretherton, 2018; Gentine et al., 2018; Yuval et al., 2021). Results have in 60 many cases been promising: NN-parameterized simulations have shown to reproduce sev-61 eral features of high-resolution simulations not found in coarse-resolution ones (Rasp et 62 al., 2018; Gentine et al., 2018). Issues with instability, model drift or energy conserva-63 tion have been widely reported, but also overcome; for instance by using loss functions 64 or model architectures which incorporate conservation laws (Beucler et al., 2019, 2021) 65 This is in itself an impressive feat considering the challenge of the problem. However, 66 all of these simulations have used highly idealized aquaplanet setups. It has yet to be 67 demonstrated that unified NN parameterizations can improve realistic climate simula-68 tions, which are much more complex and require reliable predictions across different cli-69 mates. 70

One of the most time-consuming components in coarse-resolution simulations (50%)71 in ECHAM) is the radiation scheme. This has led to attempt to replace the entire ra-72 diation scheme with machine learning models, the outputs being column radiative fluxes 73 and/or heating rates (radiation was also included in the subgrid physics emulation in many 74 of the aforementioned studies). Impressive speed-ups (1-2 orders of magnitude) relative 75 to the physical parameterization have been obtained using this approach, but it is un-76 clear if the accuracy and reliability is sufficient for state-of-the-art numerical weather pre-77 diction (NWP) and climate simulations. For instance, surface fluxes deviated by  $10 \text{ Wm}^{-2}$ 78 from the reference simulation in a prognostic evaluation with a climate model (Pal et 79 al., 2019). Recently, Song and Roh (2021) developed NNs to emulate a radiation scheme 80 in a regional NWP setting. In offline tests with independent data, predicted shortwave 81 radiation had a root-mean-square-error of roughly  $0.2 \text{ K day}^{-1}$  in heating rates and 20 82  $Wm^{-2}$  in fluxes. 83

Although these differences seem large compared to parameterization errors for clear-84 sky radiation (Hogan & Matricardi, 2020, Figure 5,7), they are less so relative to the 85 noise caused by Monte Carlo Independent Column Approximation (McICA) which rep-86 resents cloud sub-grid cloud variability stochastically and is used in many climate and 87 weather models (Räisänen et al., 2005). Stable climate simulations incorporating ML have 88 been demonstrated in several studies, with the differences in prognostic tests being sim-89 ilar to the models internal variability (Krasnopolsky et al., 2010). Yet again, a realis-90 tic climate is not sufficient evidence for accuracy. Detailed evaluation of fluxes and heat-91 ing rates across the whole atmosphere using fully independent data is rarely presented. 92 Heating rates are particularly prone to errors in the upper stratosphere: Yuval and OGor-93 man (2020) emulated subgrid tendencies from specific processes and found ML predic-94 tions of radiative heating rates in upper layers to be poor, and had to use the original 95 parameterization above 11.8 km, while Yuval et al. (2021) made the cut-off at 13.8 km. 96

An as alternative to emulating the entire radiation code, one can use ML for pre-97 dicting optical properties while still computing fluxes using a traditional solver. This may 98 be easier from a physical and algorithmic perspective since the former relies on empir-99 ical methods (look-up table interpolation) and has no dependency between adjacent at-100 mospheric layers, while the latter requires solving the radiative transfer equations to com-101 pute radiative flows through an atmospheric column. (Here a parallel can be drawn to 102 the dynamical or "resolved" part of large-scale models, since they both rely on solving 103 104 well-established physical equations to compute flows). Ukkonen et al. (2020) and Veerman et al. (2021) demonstrated high accuracy of NN-predicted optical properties of the gaseous 105 atmosphere, which were 3-4 times faster than the original RRTMGP gas optics scheme. 106 In the former study the NN gas optics were combined with a refactored radiative trans-107 fer solver to speed up clear-sky flux computations by a factor of 2-3, while the fluxes and 108 heating rates were almost identical to the original scheme when evaluated against line-109 by-line radiation computations. 110

While NN methods are powerful algorithms capable of modeling complex relationships, it is not clear that regular feed-forward neural networks (FNNs) are algorithmi-



Figure 1: Feed-forward neural networks, shown to illustrate the potential algorithmic issue with using them to model radiative transfer as is commonly done by stacking the vertical profiles of input variables (such as temperature and pressure) into one FNN input column. Because the variables or nodes are only connected horizontally to nodes in other layers, the vertical dependencies between atmospheric layers (Variable 1, Variable 2...) can only be represented indirectly through the horizontal connections (weights) to shared nodes in one or more hidden layers. Information does not propagate directly in the vertical direction, as it does in radiative transfer equations. Figure adapted from Aldakheel et al. (2021).

cally well-suited for radiative transfer problems which involve computing radiative flows 113 between mediums. In the case of radiation parameterizations used in weather and cli-114 mate models, radiative flows in a column are computed layer by layer, requiring several 115 iterations through a column. Emulating a radiation scheme by stacking vertical profiles 116 of several variables into a single input column of a machine learning model, and predict-117 ing profiles of fluxes and/or heating rates as a single output column, means that infor-118 mation needs to propagate between different inputs or nodes corresponding to adjacent 119 layers (as in the physical equations). This does not occur "directly" in an FNN where 120 nodes are connected horizontally to nodes in other layers, but not vertically to nodes in 121 the same layer (Figure 1). These vertical dependencies can of course be represented via 122 the weights connected to at least one hidden layer, but it is unclear how this can be done 123 accurately with simple neural network architectures. 124

In many ways, the results obtained in previous studies are impressive, as not only 125 does the NN approach skip explicit layer-to-layer computations, but also explicit spec-126 tral computations. Radiation codes have evolved for many decades, and the current state-127 of-the-art is to combine the two-stream approximation to the one-dimensional radiative 128 transfer equation (Meador & Weaver, 1980) with the correlated-k-distribution (CKD) 129 method (e.g., Goody et al., 1989) for the spectral integration. CKD can accurately re-130 solve broadband fluxes (i.e. fluxes integrated over the electromagnetic spectrum, which 131 relate to heating rates) while reducing the number of monochromatic computations by 132 many orders of magnitude compared to line-by-line methods. If it was true that an NN 133 could reduce the problem further by several orders of magnitude, not incorporate any 134 physical laws, and still be accurate and reliable, this would essentially mean that cur-135 rent parameterizations include wasted computations. 136

In this study, we aim to shed some light on the suitability of neural networks to replace radiation computations by addressing the following research questions: 1) Can FNNs closely emulate an entire radiation scheme, i.e. directly predict fluxes
 or heating rates with similar accuracy to existing parameterizations?

2) Is it easier to predict fluxes and heating rates accurately by only emulating com putations without a vertical dependency, such as gas optics or reflectance-transmittance
 computations using FNNs?

3) Do recurrent neural networks (RNNs) which structurally incorporate the vertical dependence of radiation computations better emulate radiative transfer then FNNs?
4) How does the trade-off between efficiency and accuracy vary across the differ-

146 147

ent emulation strategies?

To help answer these questions, we train neural nets to emulate: A. the entire radiation scheme (gas optics and radiative transfer combined), B. gas optics, and C. the reflectance-transmittance computations in the solver. Method A, which maps atmospheric conditions to fluxes or radiative heating rates, has been used in several papers but here a novel method based on RNNs is developed and compared to the standard approach using FNNs.

These emulation strategies are then compared in terms of accuracy and general-154 ization through offline validation with independent profiles, acquired from reanalysis data, 155 that span a wide range of atmospheric conditions. Since the goal is to evaluate how well 156 simple neural networks can emulate complex radiative transfer computations, this pa-157 per restricts itself to shortwave computations accounting for clouds, where the need to 158 consider scattering results in a much harder problem. Generation of training data, model 159 implementation, and verification is carried out using the recently developed RTE+RRTMGP 160 radiation scheme (Pincus et al., 2019). 161

Below, we introduce the data and codes (Section 2), followed by an overview of the different emulation strategies and associated machine learning methodologies (Section 3). The results in terms of accuracy and speed-up are then presented (Section 4) and discussed in the context of previous literature (Section 5). Finally, conclusions are given in Section 6.

#### <sup>167</sup> 2 Data and codes

#### 2.1 Data

168

Global CAMS reanalysis (Inness et al., 2019) data was acquired for 2009-2018, with 169 4 dates and 2 times of day (03 and 15 UTC 1.2, 1.5, 1.8 and 1.11) in order to encom-170 171 pass seasonal and diurnal variability of atmospheric fields. Model level variables consist of temperature, pressure, cloud liquid water and ice mixing ratio, and mixing ratios of 172 five gases that are radiatively active in the shortwave: water vapor, ozone, carbon diox-173 ide, methane and nitrous oxide. The radiation computations also account for oxygen and 174 nitrogen, but these are assumed constant (with mole fractions of 0.209 and 0.781, respec-175 tively) and therefore not included in NN inputs. The gases correspond to all the gas species 176 considered by RRTMGP-SW, with the exception of nitrogen dioxide which was not avail-177 able in CAMS. The single-level variables obtained were surface pressure, 2-metre tem-178 perature, and forecast albedo. We also computed true solar zenith angles for the pur-179 pose of model evaluation, but when generating training data, the solar angle of each col-180 umn was assigned a random value between 0 and 90. The total solar irradiance at top 181 of atmosphere is assumed constant at  $1412 \text{ Wm}^{-2}$ . 182

To avoid over-representation of polar regions in the training data, the CAMS data was interpolated from a longitude-latitude grid to a global 320 km resolution triangular grid as specified for the ICON model (Zängl et al., 2015), while keeping the the original vertical grid of 60 layers (top at 10 Pa). Each year consists of  $5120 \times 8 = 40960$  columns. Data was partitioned into validation (the year 2014), testing (2015, in which 09 and 21 UTC data was additionally included), and training (remaining 8 years in 2009-2018) subsets. It may be noted that the evaluation data does not sample new climate conditions,
but given the high variability and dimensionality of fields associated with column-wise
radiation computations, even one "in-sample" testing year should give some indication
of model generalization.

The amount of training samples depends on the emulation method (Table 1). When 193 training an emulator for the whole radiation scheme, the model inputs are columns of 194 atmospheric variables, resulting in  $40960 \times 8 = 327680$  training samples. For training 195 other models, which take input variables defined at a single spectral and/or atmospheric 196 layer, the potential training data is enormous, especially for the reflectance-transmittance 197 model which operates on individual q-points. In this case, we extract random samples 198 from the data, limiting the number of training samples to roughly 50 million (reflectance-199 transmittance) or 10 million (gas optics). 200

201

#### 2.2 RTE+RRTMGP

RTE+RRTMGP (Pincus et al., 2019) is a recently developed radiation scheme for 202 dynamical models combining two codes: Radiative Transfer for Energetics (RTE), which 203 computes fluxes given a description of boundary conditions, source functions and opti-204 cal properties of the atmosphere, and RRTM for General circulation model application-205 sParallel (RRTMGP), which computes optical properties and source functions of the gaseous atmosphere. The combined package can be used to compute broadband radiative fluxes 207 from input profiles of temperature, pressure and gas concentrations. The gas optics scheme 208 RRTMGP uses a k-distribution based on state-of-the-art spectroscopy, and has 256 q-209 points in the longwave and 224 g-points in the shortwave, which is high compared to many 210 other schemes. RRTMGP continues to evolve and preliminary reduced-resolution k-distributions 211 with roughly half the number of q-points (similar to the predecessor code RRTMG), was 212 available at the time of writing, but in this study the original 224 q-point model is used. 213 When profiling code this should favour the approach of emulating the entire radiation 214 scheme, as this method avoids explicit g-point computations while the runtime of the 215 original code (as well as emulators of components) is proportional to number of g-points. 216 Indeed, reducing the number of *g*-points, for instance by using full-spectrum correlated-217 k-methods, is a promising way to improve the accuracy/speed trade-off in radiation schemes 218 (Hogan, 2010). 219

220

#### 2.3 RTE+RRTMGP-NN

In this work, a refactored version of RTE+RRTMGP developed in tandem with 221 NN emulators for RRTMGP (Ukkonen et al., 2020) is used in order to utilize existing 222 NN code infrastructure and to get a more meaningful measure of the speedup given by 223 emulators. The refactored version (RTE+RRTMGP-NN) has columns as the outermost 224 dimension in both RRTMGP and RTE and therefore avoids expensive array transposes, 225 and also features smaller efficiency optimizations such as an optional inlining of the broad-226 band flux computation inside a column loop for reduced memory use. (This feature was 227 at the time of writing available in RTE+RRTMGP). 228

The NN inference and I/O code in RTE+RRTGMP-NN is based on neural-Fortran 229 (Curcic, 2019) but has been optimized for efficiency by packing (or re-interpreting us-230 ing pointers, when possible) the data into batches, resulting in the core operations - mul-231 tiplying layer weights with inputs - becoming a matrix-matrix multiplication that is del-232 egated to a BLAS library. Other changes include fusing the activation and bias additions, 233 as well as GPU support based on OpenACC directives and the NVIDIA cuBLAS library. 234 The end result is a highly efficient Fortran implementation of feed-forward neural net-235 works that can be used in production code. 236

The data generation workflow consisted of acquiring reanalysis data, pre-processing 237 it into yearly netCDF files that can be read by RTE+RRTMGP (for instance, gas mix-238 ing ratios were converted to mole fractions), and performing shortwave cloudy-sky com-239 putations which account for gases and clouds (but not aerosols) to generate the input-240 output pairs for machine learning. The radiation computations account for scattering, 241 and cloud optical properties were generated with a cloud optics extension in RTE+RRTMGP 242 that is based on Mie calculations. Aerosols are not included; for the purpose of evalu-243 ating the ability of NNs to emulate the physical radiation code, this should not be im-244 portant as the aerosol optical properties are simply added to the optical properties from 245 clouds and gases. The NNs were designed and trained in Tensorflow (https://www.tensorflow 246 .org) using the Keras front-end (https://keras.io), but PyTorch (https://pytorch 247 .org) code was also written to facilitate further research. A Python script was used to 248 convert the Keras models into ASCII files from which neural-Fortran loads the model 249 weights. 250

- **3** Emulation strategies
- 252 253

#### 3.1 FNN-RadScheme - emulation of the full radiation scheme using feedforward neural networks

Emulating the full radiation scheme is the best approach from the perspective of 254 efficiency, since explicit layer-to-layer computations as well as spectral computations can 255 be avoided. Internally, the radiation scheme computes many intermediate variables with a higher dimensionality than the parameterization input and outputs: first RRTMGP 257 computes gas optical properties (optical depth and single-scattering albedo) at each q-258 point and model level. The cloud optical properties (optical depth, single-scattering albedo, 259 and asymmetry parameter) are then generated for each spectral band and model level 260 and added to the gas optical properties. The radiative solver takes the optical proper-261 ties and boundary conditions (incoming solar flux, zenith angle, and surface albedo) and 262 performs radiative transfer computations for each q-point, resulting in upward and down-263 ward fluxes  $F_{\perp}$ ,  $F_{\uparrow}$  (total flux, given by diffuse plus direct flux) and direct shortwave fluxes  $F_{\downarrow,dir}, F_{\uparrow,dir}$  for each g-point and model half-level. Finally, broadband fluxes are obtained 265  $F_{\perp}$ ,  $F_{\uparrow}$  by summing the spectral fluxes together. In the NN approach, the broadband fluxes 266 are predicted directly from profiles of gas and cloud mixing ratios. This is very efficient, 267 but assumes that the spectral and vertical dependencies can be represented by the NN 268 mapping. 269

RTE+RRTMGP was used to generate output downward and upward flux profiles from profiles of gas concentrations, temperature, pressure, cloud ice and water mixing ratios, as well as the scalar variables surface albedo and cosine of the solar zenith angle. The NN outputs in this study consist only of downward and upward fluxes, and is smaller compared to other studies. Direct downward flux is omitted; while this variable would likely be needed in the host model, its computation is more straightforward and it's not needed for heating rates, and therefore less interesting for NN emulation.

Earlier studies (Krasnopolsky et al., 2010; Roh & Song, 2020; Pal et al., 2019) have predicted heating rates (HR) profiles directly as NN output, often omitting prediction of flux profiles completely and instead adding scalar flux variables at the surface and topof-atmosphere as additional NN outputs (Krasnopolsky et al., 2010; Roh & Song, 2020). Here we test predicting fluxes, while HR is given by the vertical divergence of net fluxes at each model layer *i* as in physical radiation codes:

$$\left(\frac{dT}{dt}\right)_{\text{SW radiation}} = -\frac{g}{c_p} \frac{F_{i+1/2, \text{ SW}} - F_{i-1/2, \text{ SW}}}{p_{i+1/2} - p_{i-1/2}},\tag{1}$$

where  $F_{i+1/2, SW}$  is the difference between the downward and upward SW fluxes at the interface between model layers *i* and *i* + 1,  $c_p$  is the specific heat a constant pressure, *g* is the gravitational constant and  $\frac{\partial T}{\partial t}$  is the rate of temperature change.

Computing heating rates from fluxes ensures physical consistency and energy con-286 servation Yuval et al. (2021). On the other hand, it can result in large errors in HR be-287 cause NN-predicted fluxes tend to be noisy and HR are very sensitive to the vertical gra-288 dient in fluxes, especially in the stratosphere where pressure is low. The problem can be 289 alleviated by taking special care in the NN design. Firstly, normalizing the fluxes by the 290 downward direct flux at the top layer of each column (incoming flux multiplied with the 291 cosine of the solar zenith angle) is found to reduce errors in fluxes. Effectively this phys-292 ically re-scales the output values to a range between 0 and 1. Although in some cases 293 the flux at a lower layer can exceed the incoming flux (Jiang et al., 2005), the training 294 data only had a handful of values above 1. Therefore the flux scaling is combined with 295 a sigmoid activation in the output layer to constrain outputs within the 0-1 range, which 296 was found to reduce errors. Second, a custom loss function can be used to explicitly min-297 imize the error in both flux and heating rates: 298

$$loss = \alpha (y - y_{pred})^2 + (1 - \alpha)(HR - HR_{pred})^2,$$

where y is the target value (scaled flux),  $y_{pred}$  is the NN output, HR is the heat-299 ing rate computed using equation (1), and  $\alpha$  is a manually tuned coefficient controlling 300 how much heating rates are weighted relative to fluxes. In practice, the benefit from us-301 ing a hybrid loss function was limited by the heating rates being very noisy when not 302 predicted directly, and the sensitivity of computed HR to flux errors in the upper atmo-303 sphere. This issue with noisy heating rates when predicting fluxes, manifesting in large 304 swings in the training losses (not shown), seems to be specific to FNNs as it was not seen 305 with RNNs (section 3.2). 306

Figure 2 compares the flux and heating rate errors for models using different pre-307 dictands and scaling methods. Included in the comparison is a model which predicts heat-308 ing rates profiles directly in addition to fluxes at the boundaries, as in Krasnopolsky et 309 al. (2010); Roh and Song (2020). Heating rate errors are much smaller using this method. 310 However, adding the full flux profiles as output in addition to heating rate profiles (182 311 outputs in total) led to very poor predictions at the surface in quick tests with models 312 with up to 256 neurons in two layers (not shown). To avoid the issue with physically in-313 consistent heating rates and fluxes at the boundaries, and to allow an equal footing with 314 other emulation strategies (sections 3.2 - 3.4) which can all produce flux profiles, method 315 c) in Figure 2 is used in the final evaluation despite the larger heating rate errors. This 316 may be a questionable choice, but for operational implementation the conservation of 317 energy is important, and is only guaranteed when predicting fluxes and computing heat-318 ing rates from those (the incoming solar radiation will then be equal to the energy heat-319 ing the atmosphere and the surface). 320

The FNNs have 128 neurons in three hidden layers. Two hidden layers produced only slightly worse results, but a model with a single hidden layer and 128-192 neurons had substantially larger errors.

324 325

#### 3.2 RNN-RadScheme - emulation of the full radiation scheme using bidirectional recurrent neural networks

While the feed-forward neural network can predict heating rate profiles and scalar fluxes reasonably well, on paper it still appears ill-suited for predicting radiative flows due to the lack of inter-node connections in a NN layer. The FNN approach also has the drawback of being tied to vertical resolution of the training data, as the number of inputs and outputs are fixed. A type of NN which can avoid this problem is found in the



Figure 2: Impact of scaling, loss function and predictand on the vertical profiles of mean absolute error in downwelling flux (left column), upwelling flux (middle column) and heating rate (right column) for the validation data from 2014 with randomly sampled solar zenith angles. The outputs of the different feed-forward NN models are unnormalized fluxes (a), fluxes scaled by the incoming flux (b-c), and heating rate profiles plus three flux scalar variables (d). Adding heating rate to the loss function is helpful when predicting scaled fluxes (c); with a regular loss function (b) the heating rate errors reach up to 20 K day<sup>-1</sup> at the top of atmosphere (the x-axis has been cropped at 10 K day<sup>-1</sup>). All fluxes are total (direct + diffuse) shortwave fluxes. Overall mean absolute error (MAE) is annotated, with the number in parenthesis indicating the MAE value as a percentage of the column and layer mean of the variables, which only have positive values for physically computed SW radiation. When testing each method, three separate NNs were trained, and the results with smallest heating rate errors were saved and compared.

recurrent neural network (RNN), in which connections form a directed graph. RNNs are usually applied to problems associated with temporal sequences. A RNN layer takes the input at a given sequence, updates its internal state, and then processes the next point in the sequence. The sequential nature is not present in an FNN where the output of one layer forms the input of a separate NN layer with different weights. The internal state allows the RNN to have memory so that prior inputs, i.e. from earlier in time when dealing with a temporal problem, can influence the current prediction.

This idea can be exploited for radiative transfer by letting the sequence be repre-338 sented by vertical levels. However, a basic RNN is not appropriate because the radia-339 tive fluxes at a given level depend not only on conditions at the levels above but also on 340 the levels below. Fortunately, information can propagate from future states in a bidi-341 rectional RNN (BiRNN). A BiRNN comprises of two RNNs of opposite directions con-342 nected to the same output, meaning that one RNN begins from the beginning of the se-343 quence and moves in the positive direction, while the other begins from the end of the 344 sequence and moves in the negative direction. A single BiRNN layer approach, as illus-345 trated in Figure 3 was tested. In this method the input for a given atmospheric layer is 346 used to predict the scaled downward flux at the bottom of this layer (the next half-level) 347 as well as the upward flux at the top of the layer (the previous half-level). Two output 348 variables then remain; the downward flux at the top and upward flux above the surface. 349 The first of these is actually an input and used here for scaling the fluxes, while the lat-350 ter can be physically computed from the downward flux above the surface times the sur-351 face albedo. 352

The above approach is elegant, but requires the albedo to be a broadband quan-353 tity. This happens to be true for the data used here, but may not be a valid assumption 354 generally. Furthermore, the inconsistency in how upward fluxes are computed led to larger 355 heating rate errors at the surface for a BiRNN model which otherwise performed well 356 (not shown). To remedy these issues, the model structure can be refined to output the 357 full flux profile at layer interfaces (nlay + 1), despite the inputs being defined at lay-358 ers (nlay). One way of achieving this is by concatenating layer-wise RNN outputs with 359 the output from a dense NN layer, which takes as input the albedo(s) and/or other sur-360 face quantities. This more complex approach is illustrated in Figure 4. A third RNN layer, 361 where the information propagates downward, has also been added; this was found to work 362 better than just two RNNs (one BiRNN). The structure in Figure 4 was inspired by the 363 physical equations in the radiative transfer solver, and resembles them quite closely. Three 364 vertical iterations are used there, too: one to compute direct downwelling flux starting 365 from top-of-atmosphere, one starting from the surface and computing the albedos at each 366 level using the adding-doubling method (Hansen, 1971), and a final downward pass from 367 the top-of-atmosphere to compute upward and downward fluxes. 368

While three vertical iterations within the NN model reduce the potential for speedup, 369 on the other hand the number of hidden neurons needed for accurate results is very small. 370 Here we evaluate a model using only 16 neurons in each of the three RNN layers. For 371 the RNN layers, we use gated recurrent units (GRU), which are more complex than sim-372 ple RNN lavers. A GRU laver consists of an "update gate" and a "reset gate". Here the 373 former decides if the cell state should be updated with the past (accumulated) state or 374 375 not, while the reset gate allows the network to forget past information. It is not clear how these mechanisms specifically benefit radiative transfer, but they have been found 376 to alleviate problems with vanishing gradients by allowing information to be passed with-377 out going through a nonlinear activity, thus helping preserve information from earlier 378 states. For radiation such information could relate to optical properties, or reflectances 379 and transmittances, as computed in prior states. In practice, GRU layers gave substan-380 tially better results than simple RNN layers. 381



Figure 3: A RNN-based approach to predicting radiative fluxes. Input variables defined at N model layers  $(X_0, X_1...X_N)$  form the sequential input to the Bidirectional RNN (BiRNN), while the output consist of two scalar values: the (scaled) upward flux at the upper layer boundary (the N+1 layer boundaries are referred to as *levels*), and downward flux at the lower boundary. Note that the figure shows the *unrolled* network structure; there is actually just one BiRNN layer which forms a directed graph to itself by saving a hidden state h; or two hidden states  $h_{\downarrow}, h_{\uparrow}$  in the case of the BiRNN which internally comprises of a forward and backward RNN (not shown). The auxiliary scalar inputs, albedo  $\alpha$  and cosine of the solar zenith angle  $\mu$ , are incorporated through a dense layer (DNN) which predicts the initial states of the BiRNN  $h_{0,\downarrow}, h_{N,\uparrow}$ . The diagram depicts input variables in gray, output variables in light yellow, and NN layers in light blue. The upward flux near the surface (dark yellow) is not an NN output but computed explicitly from the albedo and the downward flux at the surface.



Figure 4: A RNN-based approach to predicting fluxes at layer interfaces (N+1) from layer-wise inputs (N), consisting of three RNNs to mimic two-stream radiative transfer equations with scattering. The first RNN  $(RNN_{\downarrow})$  has a forward (downward) direction, and when it reaches the end of the sequence, i.e. the last vertical layer at N, its hidden state  $h_{\downarrow,N}$  is concatenated ("concat") with the surface albedo and fed to a dense layer, whose output is then concatenated with the RNN sequence. The dense layer essentially replaces the RNN at the boundary, where layer-wise inputs are missing. Hereafter the sequence has a length of (N+1), and is connected to a backward/upward RNN  $(RNN_{\uparrow})$ . Then, the first two sequences are concatenated (as is usually done in a bidirectional RNN) and connected to a third and final RNN  $(RNN2_{\downarrow})$ . Finally, the sequential output from this RNN is connected to a dense layer which predicts two values, the upwelling and downwelling fluxes scaled by incoming flux.

#### 3.3 FNN-RRTMGP - emulation of gas optics only

Successful NN emulators for RRTMGP gas optics have been developed in earlier work: in Veerman et al. (2021), average flux errors were within 0.5 Wm<sup>-2</sup> of RRTMGP, while in Ukkonen et al. (2020) root-mean-square-errors (RMSE) in heating rate with respect to line-by-line results were virtually identical with RRTMGP. Here an identical NN methodology as in Ukkonen et al. (2020) is used, which involves predicting absorption and Rayleigh cross-sections with two separate NNs.

The main advantage of using neural networks for gas optics is efficiency: whereas 389 the original kernel computes optical properties separately for each band and each mi-390 nor gas species (the absorption due to two major gases in a band is computed separately 391 and parameterized to account for overlap in the absorption spectra), the NN can take 392 all gases as one input vector and predict the optical properties for all q-points as one out-393 put vector. Consequently, minor greenhouse gases can be included with almost no ad-394 ditional cost. NNs are also suitable for predicting optical properties from a physical per-395 spective, since the original kernel relies on empirical look-up-tables and incorporates no 396 physical laws explicitly. Further benefits are generalization to arbitrary vertical grids by 397 predicting layer-wise optical properties normalized by number of molecules (cross-sections), 398 and that a NN can treat gas overlap implicitly. In theory, a novel NN gas optic model 399 could be trained directly on data generated with line-by-line radiation codes to avoid er-400 rors associated with gas overlap assumptions, but the data generation would be a sig-401 nificant computational challenge. 402

403 404

382

#### 3.4 FNN-RefTrans - emulation of reflectance-transmittance computations

Training NNs to emulate the radiative transfer solver was considered for this work, but because RTE and other solvers perform computations per *g*-point, an emulator which respects the underlying physics and similarly operates on *g*-points is unlikely to be more efficient (broadband fluxes could be predicted directly, but the inputs are still defined per *g*-point).

An alternative is focusing on computations of reflectance and transmittance in the shortwave solver. While the efficiency drawback of explicit *g*-point computations remain, this may be more promising for FNNs since the problem has a simpler nonlinear inputoutput mapping which does not include vertical dependencies. The reflectance-transmittance computations (kernel sw\_two\_stream) are furthermore the slowest part of RTE and exhibit a high sensitivity to numerical precision.

Simple neural networks are able to predict direct and diffuse reflectance and trans-416 mittances with high accuracy (Figure 5). However, when implementing the NNs into the 417 radiation code it was discovered that even very small inaccuracies overall (with R-squared 418 > 0.999 for each variable) can translate into significant RMSE and maximum errors in 419 net fluxes; typically tens and hundreds of  $Wm^{-2}$  respectively. A possible explanation 420 is a larger sensitivity for errors at specific values of reflectance and transmittance, spe-421 cific q-points (which contribute to broadband flux more strongly than others), or spe-422 cific atmospheric levels, or just a high sensitivity in the dependence of flux on reflectance 423 and transmittance in general. For instance, predicting intermediate values of transmit-424 tance accurately may be more important than values near zero, since the latter case is 425 likely to be associated with radiation being fully extinguished (reflected or absorbed). 426 The distribution of the predictands is highly skewed with such intermediate values be-127 ing rare, and as a result are also associated with larger errors when employing a regu-428 lar loss function. 429

To combat this problem, one could devise custom loss functions, data transformation, or synthetic data generation to create more samples for the important but under-



Figure 5: Comparison of the predicted (y-axis) and true (x-axis) reflectance and transmittance values using the validation data set and final REFTRANS model, which has 12 neurons in two hidden layers. The colors on the scatter plot correspond to the occurrence on a log-scale.

represented parts of the distribution. A simple data transformation which reduced errors in radiative flux was to take the square root of the output prior to training, which makes the distribution more Gaussian (albeit still highly non-Gaussian). We then tested custom loss functions which give smaller weights to intermediate values of all four outputs, and/or a smaller weight to diffuse transmittance, but found no clear improvements in predicted fluxes. Figure 5 shows the validation performance of the final model.

438

#### 3.5 Summary of model architectures and methodologies

The architecture and pre-processing used for the different NN emulators are de-439 scribed in Table 1. The model hyperparameters (number of hidden neurons, hidden lay-440 ers and activation functions) as well as suitable pre-processing methods - which were of-441 ten found to be more important than NN hyperparameters - were tuned by hand. The 442 objective of this laborious tuning process was to find a reasonable trade-off between ac-443 curacy and model complexity, which determines the computational cost. This restricted 444 the reflectance-transmittance emulator to a very simple NN model, as it turned out to 445 be difficult to surpass the efficiency of the original computations. For the FNN emulat-446 ing the entire radiation scheme, efficiency was less of a consideration, as the inference 117 code using this method was very fast regardless. 448

The hyperparameters of the gas optics emulator were taken from Ukkonen et al. (2020). All models were trained using the Adam optimizer (Kingma & Ba, 2015) and the early-stopping method, which stops training when the validation error has not improved for a certain number of epochs (here 28). The batch size was set to 1024.

Model	FNN-RadScheme	RNN-RadScheme	FNN-RRTMGP	FNN-RefTrans
Emulated compo- nent	Radiation scheme with gas and cloud optics	Radiation scheme with gas and cloud optics	Gas optics	Reflectance- transmittance computations
Input	scalars $\alpha$ and $\mu_0$ + vertical profiles of gas mole fractions, T, p, cloud ice and cloud water	same as for FNN- RadScheme but one layer at a time	gas mole fractions, T, p	$ au, ssa, g, mu, T_{noscat}$
Input size	$2+9 \cdot nlay = 542$	2+9=11	5	7
Output	vertical profiles of broadband fluxes	vertical profiles of broadband fluxes	$ \begin{vmatrix} absorption/Rayleigh \\ cross-sections as a \\ vector of g-points \end{vmatrix} $	$\begin{vmatrix} R_{dif}, T_{dif}, R_{dir}, \\ T_{dir} \end{vmatrix}$
Output size	$2 \cdot nlev = 122$	2	ngpt = 224	4
Hidden layers	Dense, Dense, Dense	RNN, Dense, RNN, RNN, Dense	Dense, Dense	Dense, Dense
Activation func- tions in hidden layers and output layer	relu, relu, relu, sigmoid	tanh, linear, tanh, tanh, sigmoid	softsign, softsign, linear	softsign, softsign, hard sigmoid
Neurons in each hidden layer	128	16	16	12
Total parameters	118,266	5,698	4,208	280
Required itera- tions	ncol	$\left \begin{array}{c} ncol \ (\times \ 3 \ nlay \\ internallly) \end{array}\right $	$\left \begin{array}{c} ncol \times nlay \\ (\times \ 2 \ \text{NN models}) \end{array}\right $	$\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ $
Flexible with re- gards to vertical grid	No	Yes	Yes	Yes
Input scaling	$x_i = \frac{x_i}{\max(x_i)}$	$x_i = \frac{x_i}{max(x_i)}$	$\begin{vmatrix} x = log(x) \text{ for p;} \\ x = x^{\frac{1}{4}} \text{ for H}_2O \\ \text{and } O_3; \\ x_i \\ \frac{x_i - min(x_i)}{max(x_i) - min(x_i)} \end{vmatrix}$	$x = x^{\frac{1}{8}} \text{ for } \tau;$ $x = \frac{x}{max(x)} \text{ for }$ $\tau \text{ (other features already in 0-1 range)}$
Output scaling	$\left  \begin{array}{c} y_i = \frac{y_i}{F_{\downarrow,0}} \end{array} \right $	$y_i = \frac{y_i}{F_{\downarrow,0}}$	$\left  \begin{array}{c} y=y^{\frac{1}{8}};\\ y_i=\frac{y_i-\bar{y}_i}{\sigma_y} \end{array} \right.$	$y = y^{\frac{1}{4}}$

Table 1: Description of the different models. Abbreviations:  $\alpha$  = surface albedo,  $\mu_0$  = cosine of solar zenith angle, T = temperature, p = pressure, *ssa* = single-scattering albedo, g = asymmetry parameter, T = transmittance, R = reflectance, H<sub>2</sub>O = mixing ratio of water vapor, O<sub>3</sub> = mixing ratio of ozone.



Figure 6: Vertical profiles of the error in shortwave downwelling flux (left column), upwelling flux (middle column) and heating rates (right column) for the test data (2015) using different emulation methods: replacing the radiation scheme with (a) a feedforward NN or (b) a bidirectional recurrent NN, (c) replacing only the radiative solver's reflectance-transmittance computations with a FNN, or (d) replacing the gas optics computations with a FNN. The solid and dotted lines show the mean error and mean absolute error, respectively, while the shaded area indicates the 5th and 95th percentile of mean absolute error. For FNN-Radscheme (a) the mean heating rate errors at TOA (0.01 Pa) reach around 3.5 K day<sup>-1</sup> (the x-axis has been cropped). In the annotated statistics, the number in parenthesis gives the error as a percentage of the column and layer mean of the variable.

#### 453 4 Results

454

#### 4.1 Accuracy

The models are evaluated by comparing the final output of the radiation code, fluxes and heating rates, to a reference result computed in double precision using RTE+RRTMGP. (Comparison to a single precision result would be very similar, as the NN errors are larger than those from using reduced precision.)

The errors in flux and heating rate using different emulators is shown in Figure 6. In this offline evaluation based on one year of independent global data, all emulation methods produce fluxes with R-squared values very close to 1 and mean absolute errors around <sup>462</sup> 1% or less. In the case of gas optics emulation (FNN-RRTMGP), there is practically no
<sup>463</sup> error in fluxes. The emulation of the whole scheme (FNN-RadScheme) gives a similar
<sup>464</sup> accuracy in flux compared to emulating only reflectance-transmittance computations (FNN-RefTrans), which is a poor result for the latter method, as it is far more expensive.

Heating rates computed from these fluxes show much larger differences across em-466 ulators. FNN-RadScheme has the largest heating rate errors with a mean absolute er-467 ror (MAE) of 0.50 K day<sup>-1</sup>, or 25.5% when expressed as a percentage of the mean HR 468 in the dataset. The radiation scheme emulator based on recurrent NNs (RNN-RadScheme) 469 produces far more accurate heating rates despite not predicting them directly, with a MAE 470 of 0.07 and RMSE of 0.16 K day<sup>-1</sup>. FNN-RefTrans reproduces heating rates well rel-471 ative to fluxes; with errors well below 0.5 K day<sup>-1</sup> throughout most of the atmosphere 472 despite the flux errors being comparable to FNN-RadScheme. The most accurate heat-473 ing rates are seen with FNN-RRTMGP with a MAE of only 0.02 K day<sup>-1</sup>. 474

For simulating climate, the upwelling flux at the top-of-atmosphere (TOA) is an important quantity. All emulators have small errors in TOA upwelling flux (Figure 7): less than 1 Wm<sup>-2</sup> for all models but FNN-RefTrans. Likewise, the downwelling flux at surface is predicted within roughly 1% by all emulators (Figure 8).

4.2 Speed-up

479

Speedup of the radiation codes were measured on a modern workstation with both 480 reference and NN computations performed in single precision. A fair comparison is en-481 sured by implementing all NN models, with the exception of the RNN, in the RTE+RRTMGP-482 483 NN Fortran code and including the overhead from pre- and post-processing. Principal timings were done using the AMD Ryzen 7 5800H processor and GNU compiler version 181 11 (compiler options -march=native -O3). The matrix-matrix computations in RTE+RRTGMP-485 NN were accelerated using AMD BLIS (https://developer.amd.com/amd-aocl/) ver-486 sion 3.0.6. The Fortran code uses blocking of the columns for better cache performance; 487 for each emulator, an optimal block size was used. All timings represent the best result 488 from three trials. 489

The computation of cloudy-sky fluxes for the 81920 test columns took roughly 18.5 490 seconds using the reference code a single CPU core. By comparison, the FNN-RadScheme 491 computed fluxes in just 0.35 seconds, a 52-fold speed-up. This is similar to what has been 492 reported in other studies (e.g., Song & Roh, 2021). Replacing only the gas optics com-493 ponent with a FNN reduces the the runtime of the gas optics by a factor of 3, but the total runtime by only 25%. This reflects that the solver is the most expensive part of 495 SW radiation computations in optimized RTE+RRTMGP (Ukkonen et al., 2020). Fi-496 nally, the reflectance-transmittance emulator is not faster than the original code, but 40-497 45% slower. This is despite the FNN being a very simple model with only 280 param-498 eters. The slowness of the method can be attributed to it operating on individual spec-499 tral points as does the original code, but not being tailored as the physical equations, 500 resulting in redundant computations. 501

Finally, the RNN and FNN models predicting fluxes was evaluated within Python 502 using the ONNX Runtime Library (ORT) version 1.9.0, first using a CPU (single core). 503 This was necessary because the neural-Fortran library does not support RNNs. These 504 timings do not include pre- and post-processing, but those accounted for less than 10%505 of the runtime of FNN-RadScheme in Fortran. The inference with the RNN emulator 506 took roughly 4.1 seconds using ORT, representing a speed-up of 4.5X over the reference 507 code in Fortran. This is a significant speed-up, but much smaller than obtained with the 508 FNN model. To compare the FNN and RNNs on a single platform, the ONNX timings 509 were also done for FNN-RadScheme, which in this instance took 0.21 seconds. It can be 510 concluded that the recurrent NN approach is roughly 20 times slower than an FNN-based 511 approach on CPUs. The performance on a GPU (RTX 3060 Mobile) was then briefly 512



Figure 7: Global upwelling shortwave flux at top-of-atmosphere for the testing year 2015 as computed with RTE+RRTMGP (a) and the grid box mean differences in this quantity using different emulators (b - e). Bulk error statistics with respect to individual columns are displayed on the right hand side.



Figure 8: As in Figure 7, but for downwelling flux at surface.

evaluted using ORT. The RNN inference time is reduced to 0.34 seconds on the GPU,
while the FNN inference took a mere 0.046 seconds. The performance gap between the
FNN and RNN-based approaches for radiative transfer is therefore reduced considerably
when using GPUs, here to roughly 7.4X.

#### <sup>517</sup> 5 Benefits of targeted and physics-informed machine learning

All the emulators evaluated here produce very reasonable fluxes, but the large sensitivity of heating rates and noise in the fluxes predicted by feed-forward NNs results in relatively large heating rate errors. Some other studies have sidestepped this issue by predicting heating rates directly, implying a lack of energy conservation which may or may not be an issue in practice but is nonetheless undesirable in an operational setting.

The heating rate errors with flux-predicting FNN seems to be caused by the fact 523 that the NN outputs at different atmospheric levels are not structurally correlated with 524 outputs at other levels. The recurrent NN, which does have incorporate the vertical de-525 pendence, produces far more accurate heating rates. The RMSE of 0.16 K day<sup>-1</sup>, eval-526 uated across the whole atmosphere with the uppermost layer at 10 Pa, is smaller than 527 the errors reported in other studies. For instance, shortwave heating rates had an offline 528 RMSE of 0.5 K day<sup>-1</sup> in Roh and Song (2020) and 0.17 K day<sup>-1</sup> in Song and Roh (2021). 529 In both of these studies, the vertical grid only reached 50 hPa and heating rates were 530 predicted directly with an FNN. With this in mind these initial results with an RNN are 531 very promising, and the errors are in fact similar in magnitude to parameterization er-532 rors associated with the correlated-k distribution method (Hogan & Matricardi, 2020, 533 Figure 7). The drawback of the RNN approach is that its sequential nature, which lets 534 it emulate a radiation parameterization more closely, also makes it less efficient than FNNs. 535 However, a speed-up of more than 4 times is still significant, and when testing with a 536 GPU a much larger speed-up was seen (a factor of 54 relative to a single CPU core). 537

Smooth flux profiles, associated with small heating rate errors, are also seen with FNN-RRTMGP and FNN-RefTrans, demonstrating the advantage of retaining the radiative transfer equations. While the FNN-RefTrans model is considerably slower than the original code, and therefore found to be an unsuccessful emulation target, the gas optics emulation produces extremely accurate results while speeding up the original lookup-table by several factors.

Regarding the choice of output, while it may seem attractive to predict heating rates 544 directly in addition to fluxes at boundaries, it should also be noted that it could lead to 545 larger errors in fluxes: the RMSE in SW flux was around 15  $\mathrm{Wm}^{-2}$  in offline evaluation in both Roh and Song (2020) and Song and Roh (2021). By comparison, the MAE in 547 SW upwelling flux at TOA and downwelling flux at surface were around  $1 \text{ Wm}^{-2}$  or less 548 for both FNN-RadScheme and RNN-RadScheme. It is unclear, however, why tests with 549 a heating rate predicting FNN had relatively small errors in the boundary fluxes in this 550 study (Figure 2). Our experience is that hyperparameters (number of hidden layers, ac-551 tivation functions used in the output layer) and other technical details in how ML mod-552 els are developed can have a substantial impact on the results. Unfortunately, these are 553 not always well documented. A great example is pre-processing of both inputs and out-554 puts, which can have a major impact. Besides such more overlooked aspects, the quan-555 tity of training data can obviously be an important factor. In this study, the number of 556 training profiles was initially an order of magnitude smaller, and model errors significantly 557 worse. 558

How the NN emulators would perform in a prognostic evaluation when embedded in a large-scale model is a critical question. Such experiments were considered to be out of the scope of the present work. Nonetheless, comparison with studies which have performed both offline and online errors should give some indication. These include the stud-

ies mentioned above with larger offline errors, where prognostic evaluation based on a 563 squall-line simulation (Roh & Song, 2020) and a regional NWP simulation (Song & Roh, 564 2021) did not show a significant degradation for precipitation and temperature forecasts, 565 and forecasts were improved relative to infrequent calls of the original scheme at the same 566 computational cost. In another study, NN output consisting of both flux and heating rate 567 profiles had mean errors of a few percentage points in an offline setting (Figure 1 in Pal 568 et al., 2019). In year-long climate simulations, the NN parameterization resulted in time-569 and area-averaged SW surface downwelling fluxes that differed substantially from the 570 reference simulation, but the differences were comparable to the internal variability of 571 the model. 572

#### 573 6 Conclusions

Emulating a sub-component of a physics scheme reduces the potential to speedup, but can greatly improve accuracy and generalization. For operational implementation, the fact that the dimensionality is much smaller is important, because it allows sampling the input space more thoroughly. Accelerating computations of reflectance and transmittance using NNs was not successful, but the gas optics component is relatively straightforward to emulate at high accuracy, and the FNNs are much faster than the look-uptable method of the original code.

It was also found that transforming inputs and outputs prior to training can have a substantial impact on the accuracy in both the physical output variable as well as derived variables which are not directly predicted. Scaling shortwave fluxes by the incoming TOA flux reduces flux errors substantially, but at the expense of heating rate errors when using a feed-forward NN. A loss function which computes the heating rate error alleviated the issue, but predicting heating rates directly (as opposed to fluxes) may be necessary to produce accurate heating rates with a feed-forward NN.

Finally, this study has contributed to more accurate emulation of radiation com-588 putations by developing a recurrent NN method that can predict fluxes at layer inter-589 faces from inputs defined at levels and the surface. The author is not aware of previous 590 work using recurrent NNs to compute radiative fluxes in a vertical column. This method 591 is in principle flexible with regards to the vertical grid. A model of roughly 5,600 param-592 eters which consists of three RNN layers, propagating information in both directions of 593 the vertical column (mimicking radiative transfer computations), is able to predict fluxes and heating rates far better than a FNN with more than 100,000 parameters. Fewer pa-595 rameters, in turn, makes it much easier to build general models which can replace pa-596 rameterizations in real applications. While the speedup offered by the RNN is smaller 597 than with FNNs, it still offered a 4-fold speedup on a CPU and a 54-fold speedup on GPU 598 relative to running the original scheme a single CPU core. Future work should investi-599 gate the RNN approach further by implementation in a large-scale model. 600

#### 601 Acknowledgments

The author would like to thank Kristian Pagh Nielsen for helping him understand ra-602 diation and how radiation schemes interact with other parts of large-scale models, as well 603 as for reviewing the manuscript. The code used in this work is available on Github on 604 a dedicated branch of the RTE+RRTMGP-NN package (https://github.com/peterukk/ 605 rte-rrtmgp-nn/tree/nn\_dev/examples/emulator-training); which includes Python 606 and Fortran code for data retrieval, pre-processing, data generation, model training, and 607 model evaluation. The machine learning input-output data can be found on Zenodo (https:// 608 doi.org/10.5281/zenodo.5564314). The author acknowledges the financial support 609 from the European Commission for the project: "Energy-efficient SCalable Algorithms 610 for weather and climate Prediction at Exascale" (ESCAPE-2) Horizon 2020 project, which 611 has enabled us to do this work. 612

#### References 613

632

636

637

638

645

646

647

648

649

650

651

652

653

- Aldakheel, F., Satari, R., & Wriggers, P. (2021, July). Feed-forward neural net-614 works for failure mechanics problems. , 11(14), 6483. Retrieved from https:// 615 doi.org/10.3390/app11146483 doi: 10.3390/app11146483 616
- Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforc-617 ing analytic constraints in neural networks emulating physical systems. Physi-618 cal Review Letters, 126(9), 098302. 619
- Beucler, T., Rasp, S., Pritchard, M., & Gentine, P. (2019). Achieving conservation 620 of energy in neural network emulators for climate modeling. arXiv preprint 621 arXiv:1906.06622. 622
- Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural net-623 work unified physics parameterization. Geophysical Research Letters, 45(12), 624 6289–6298. doi: 10.1029/2018gl078510 625
- Curcic, M. (2019). A parallel fortran framework for neural networks and deep learning. In Acm sigplan fortran forum (Vol. 38, pp. 4–21). 627
- Fuhrer, O., Chadha, T., Hoefler, T., Kwasniewski, G., Lapillonne, X., Leutwyler, D., 628 ... Vogt, H. (2018, May). Near-global climate simulation at 1 km resolution: 629 establishing a performance baseline on 4888 GPUs with COSMO 5.0., 11(4), 630 1665 - 1681.Retrieved from https://doi.org/10.5194/gmd-11-1665-2018 631 doi: 10.5194/gmd-11-1665-2018
- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018).Could 633 machine learning break the convection parameterization deadlock? Geophysical 634 Research Letters, 45(11), 5742–5751. doi: 10.1029/2018gl078202 635
  - Goody, R., West, R., Chen, L., & Crisp, D. (1989). The correlated-k method for radiation calculations in nonhomogeneous atmospheres. Journal of Quantitative Spectroscopy and Radiative Transfer, 42(6), 539–550.
- Hansen, J. E. (1971).Multiple scattering of polarized light in planetary atmo-639 spheres. part i. the doubling method. Journal of Atmospheric Sciences, 28(1), 640 120 - 125.641
- Hogan, R. J. (2010).The full-spectrum correlated-k method for longwave atmo-642 spheric radiative transfer using an effective planck function. Journal of the at-643 mospheric sciences, 67(6), 2086–2100. 644
  - Hogan, R. J., & Matricardi, M. (2020). Evaluating and improving the treatment of gases in radiation schemes: the correlated k-distribution model intercomparison project (ckdmip). Geoscientific Model Development, 13(12), 6501–6521.
  - Inness, A., Ades, M., Agusti-Panareda, A., Barré, J., Benedictow, A., Blechschmidt, A.-M., ... others (2019).The cams reanalysis of atmospheric composition. Atmospheric Chemistry and Physics, 19(6), 3515–3556.
  - Jiang, S., Stamnes, K., Li, W., & Hamre, B. (2005).Enhanced solar irradiance across the atmosphere-sea ice interface: a quantitative numerical study. Applied optics, 44(13), 2613–2625.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In 654 Y. Bengio & Y. LeCun (Eds.), 3rd international conference on learning rep-655 resentations, ICLR 2015, san diego, ca, usa, may 7-9, 2015, conference track 656 proceedings. Retrieved from http://arxiv.org/abs/1412.6980 657
- Krasnopolsky, V., Fox-Rabinovitz, M., Hou, Y., Lord, S., & Belochitski, A. (2010).658 Accurate and fast neural network emulations of model radiation for the ncep 659 coupled climate forecast system: climate simulations and seasonal predictions. 660 Monthly Weather Review, 138(5), 1822–1842. 661
- Meador, W., & Weaver, W. (1980). Two-stream approximations to radiative transfer 662 in planetary atmospheres: A unified description of existing methods and a new 663 improvement. Journal of Atmospheric Sciences, 37(3), 630–643. 664
- Odyssee-Mure. (2021).Sectoral profile households - electricity consumption 665 per dwelling. Retrieved 2021-02-01, from https://www.odyssee-mure 666

.eu/publications/efficiency-by-sector/households/electricity 667 -consumption-dwelling.html 668 Pal, A., Mahajan, S., & Norman, M. R. (2019).Using deep neural networks as 669 cost-effective surrogate models for super-parameterized e3sm radiative transfer. 670 Geophysical Research Letters, 46(11), 6069–6079. 671 Pincus, R., Mlawer, E., & Delamere, J. (2019). Balancing accuracy, efficiency, and 672 flexibility in radiation calculations for dynamical models. Earth and Space 673 Retrieved from https://www.essoar.org/doi/abs/ Science Open Archive. 674 10.1002/essoar.10500964.2 doi: 10.1002/essoar.10500964.2 675 Räisänen, P., Barker, H. W., & Cole, J. (2005). The monte carlo independent col-676 umn approximations conditional random noise: Impact on simulated climate. 677 Journal of climate, 18(22), 4715–4730. 678 Rasp, S., Pritchard, M. S., & Gentine, P. (2018, September). Deep learn-679 ing to represent subgrid processes in climate models. , 115(39), 9684-680 Retrieved from https://doi.org/10.1073/pnas.1810286115 9689. doi: 681 10.1073/pnas.1810286115682 Roh, S., & Song, H.-J. (2020). Evaluation of neural network emulations for radia-683 tion parameterization in cloud resolving model. Geophysical Research Letters, 684 47(21), e2020GL089444. 685 Song, H.-J., & Roh, S. (2021). Improved weather forecasting using neural network 686 emulation for radiation parameterization. Journal of Advances in Modeling 687 Earth Systems, e2021MS002609. doi: 10.1002/essoar.10506992.1 688 Ukkonen, P., Pincus, R., Hogan, R. J., Nielsen, K. P., & Kaas, E. (2020).Accel-689 erating radiation computations for dynamical models with targeted machine 690 learning and code optimization. Journal of Advances in Modeling Earth Sys-691 tems, 12(12), e2020MS002226. doi: 10.1029/2020ms002226 Veerman, M. A., Pincus, R., Stoffer, R., Van Leeuwen, C. M., Podareanu, D., & 693 Van Heerwaarden, C. C. (2021). Predicting atmospheric optical properties for 694 radiative transfer computations using neural networks. *Philosophical Transac*-695 tions of the Royal Society A, 379(2194), 20200095. 696 Yuval, J., O'Gorman, P. A., & Hill, C. N. (2021). Use of neural networks for stable, 697 accurate and physically consistent parameterization of subgrid atmospheric 698 processes with good performance at reduced precision. Geophysical Research 699 Letters, 48(6), e2020GL091363. 700 Yuval, J., & OGorman, P. A. (2020). Stable machine-learning parameterization of 701 subgrid processes for climate modeling at a range of resolutions. Nature com-702 munications, 11(1), 1–10. 703 Zängl, G., Reinert, D., Rípodas, P., & Baldauf, M. (2015).The icon (icosahedral 704 non-hydrostatic) modelling framework of dwd and mpi-m: Description of the 705 non-hydrostatic dynamical core. Quarterly Journal of the Royal Meteorological 706 Society, 141(687), 563–579. 707