

Estimating Full Longwave and Shortwave Radiative Transfer with Neural Networks of Varying Complexity

RYAN LAGERQUIST^{a,b}, DAVID D. TURNER^b, IMME EBERT-UPHOFF^{a,c} AND JEBB Q. STEWART^b

^a Cooperative Institute for Research in the Atmosphere, Colorado State University, Fort Collins, Colorado

^b National Oceanic and Atmospheric Administration/Global Systems Laboratory, Boulder, Colorado

^c Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, Colorado

(Manuscript received 2 February 2023, in final form 22 June 2023, accepted 24 August 2023)

ABSTRACT: Radiative transfer (RT) is a crucial but computationally expensive process in numerical weather/climate prediction. We develop neural networks (NN) to emulate a common RT parameterization called the Rapid Radiative Transfer Model (RRTM), with the goal of creating a faster parameterization for the Global Forecast System (GFS) v16. In previous work we emulated a highly simplified version of the shortwave RRTM only—excluding many predictor variables, driven by Rapid Refresh forecasts interpolated to a consistent height grid, using only 30 sites in the Northern Hemisphere. In this work we emulate the full shortwave and longwave RRTM—with all predictor variables, driven by GFSv16 forecasts on the native pressure–sigma grid, using data from around the globe. We experiment with NNs of widely varying complexity, including the U-net++ and U-net3+ architectures and deeply supervised training, designed to ensure realistic and accurate structure in gridded predictions. We evaluate the optimal shortwave NN and optimal longwave NN in great detail—as a function of geographic location, cloud regime, and other weather types. Both NNs produce extremely reliable heating rates and fluxes. The shortwave NN has an overall RMSE/MAE/bias of 0.14/0.08/−0.002 K day^{−1} for heating rate and 6.3/4.3/−0.1 W m^{−2} for net flux. Analogous numbers for the longwave NN are 0.22/0.12/−0.0006 K day^{−1} and 1.07/0.76/+0.01 W m^{−2}. Both NNs perform well in nearly all situations, and the shortwave (longwave) NN is 7510 (90) times faster than the RRTM. Both will soon be tested online in the GFSv16.

SIGNIFICANCE STATEMENT: Radiative transfer is an important process for weather and climate. Accurate radiative transfer models exist, such as the RRTM, but these models are computationally slow. We develop neural networks (NNs), a type of machine learning model that is often computationally fast after training, to mimic the RRTM. We wish to accelerate the RRTM by orders of magnitude without sacrificing much accuracy. We drive both the NNs and RRTM with data from the GFSv16, an operational weather model, using locations around the globe during all seasons. We show that the NNs are highly accurate and much faster than the RRTM, which suggests that the NNs could be used to solve radiative transfer inside the GFSv16.

KEYWORDS: Longwave radiation; Radiative transfer; Shortwave radiation; Deep learning; Machine learning; Neural networks

1. Introduction

Radiative heating is a main driver of Earth's climate and the only process by which Earth can exchange energy with the rest of the universe; radiative transfer (RT) is the governing theory. In RT studies the electromagnetic spectrum is often separated into the shortwave part (wavelength $\leq 4 \mu\text{m}$), which is mostly emitted by the sun, and the longwave part ($\geq 4 \mu\text{m}$), which is mostly emitted by Earth—both its surface and atmosphere.¹ The global distribution of top-of-atmosphere

(TOA) incoming shortwave radiation is controlled largely by geographic variations in the solar zenith angle and surface albedo, with low (high) zenith angle and albedo at the low (high) latitudes.² This sets up a strong meridional gradient in TOA incoming shortwave radiation, with higher values at lower latitudes. The global distribution of TOA outgoing longwave radiation is somewhat similar, because warmer surfaces (at lower latitudes) emit more longwave radiation than colder surfaces. However, the longwave distribution is more complicated, because longwave radiation interacts more strongly with atmospheric gases. Overall, the low latitudes have a surplus of net radiation (TOA incoming shortwave minus TOA outgoing longwave), while the high latitudes have a deficit. This imbalance maintains the meridional temperature gradient we observe, as well as driving the global atmospheric circulation, including a strong poleward heat flux produced by baroclinic waves (Wallace and Hobbs 2006).

¹ The 4- μm threshold is not an exact constant; sometimes other values are used.

Supplemental information related to this paper is available at the Journals Online website: <https://doi.org/10.1175/JTECH-D-23-0012.s1>.

Corresponding author: Ryan Lagerquist, ralager@colostate.edu

² Clouds (both liquid and ice; Tang et al. 2020) and aerosols (Myhre et al. 2013) also play a major, though highly uncertain, role in Earth's shortwave radiation budget.

RT is also crucially important for day-to-day weather prediction, because it causes differential diabatic heating. In numerical weather prediction (NWP), this diabatic heating is a subgrid-scale process and is therefore parameterized by a separate RT model. The most accurate RT models are line-by-line models (Turner et al. 2004; Mlawer and Turner 2016), but these are far too slow for NWP. A popular compromise is the Rapid Radiative Transfer Model (RRTM; Mlawer et al. 1997), a hybrid physical–statistical model that is nearly as accurate as line-by-line models but millions of times faster. The RRTM, like most RT models, performs one-dimensional RT, assuming that RT occurs only in the vertical. Faster variants—such as the RRTM for global climate models (RRTMG; Pincus and Stevens 2013), RRTMG Parallel (RRTMG-P; Pincus et al. 2019), and RRTMG-K (Baek 2017)—are often used in NWP as well. However, the RRTM and its variants are still computationally expensive, accounting for 20%–50% of the total computing of the host NWP model (e.g., Cotronei and Slawig 2020). We have elected to emulate the RRTM³ because, by using more quadrature points, it is more accurate than the RRTMG.

This has motivated a body of work on using neural networks (NN; Part II of Goodfellow et al. 2016), an algorithm from machine learning (ML), to emulate RT models, dating back to Chevallier et al. (1998). ML-based emulation of RT and other subgrid-scale processes almost always uses NNs, so we omit other ML algorithms from this review. The main advantage of NNs is that they can accurately model complex relationships (hence “universal function approximators”; see, e.g., Sonoda and Murata 2017) and are much faster than the RRTM and its variants at inference time, i.e., when applying a trained NN to predict on new data. The main disadvantage is that they are purely statistical models and, without physical constraints, may not generalize well to conditions outside the range of their training data, such as future climates. Also, adding predictor variables to an NN requires complete retraining. An overall disadvantage of replacing parameterizations such as the RRTM is that the host NWP models are very sensitive to changes in parameterizations (Boukabara et al. 2019; Rasp 2020; Muñoz-Esparza et al. 2022). Thus, even if the RT emulator has very small errors in offline testing (outside the NWP model), during online testing (inside the NWP model) these errors may accumulate or cause undesired feedbacks to other components of the NWP model, degrading the quality of the overall weather forecast. However, if successfully integrated into an NWP model, an NN-based RT emulator can decrease computing requirements by orders of magnitude.

The current article expands on work presented in Lagerquist et al. (2021, henceforth L21). Differences between this work and L21 are listed at the end of the introduction. The following review focuses on recent work in RT emulation, especially work published after L21. We divide recent work into four categories:

emulating RT in climate models, emulating RT in weather models, emulating only part of an RT model such as gas optics, and miscellaneous.

In climate modeling, Pal et al. (2019) developed an RT emulator for the superparameterized Energy Exascale Earth System Model (SP-E3SM) and found in online testing that the emulator produces a similar climate to the original RT model. Beucler et al. (2021) used climate-invariant NNs to emulate both RT and other subgrid-scale processes in climate models. They ensured climate invariance by rescaling three predictor variables for the NN—temperature, humidity, and latent heat flux—to forms that are not projected to increase with global warming. Without rescaling, applying the trained NN to future climates involved extrapolating (e.g., applying the NN to temperatures higher than any seen in the training data), which degraded performance. Beucler et al. found that rescaling allows their NN to predict subgrid-scale processes well, including RT, in a climate 8 K warmer than the climate used for training. Belochitski and Krasnopolsky (2021) used an emulator developed in 2011 for the Climate Forecast System (CFS) and integrated it into version 16 of the Global Forecast System (GFSv16). They found that the emulator generalized well between the host models without retraining—i.e., the GFSv16 with the emulator produced a similar climate to the GFSv16 with the original RRTMG parameterization. However, this success was achieved only after changing the number of heights and prognostic variables in the GFSv16 to match the CFS.

In weather modeling, much recent work has been done at the Korea Meteorological Administration (KMA). Roh and Song (2020) became the first to emulate RT at cloud-resolving resolution, developing NNs for a 250-m version of the Weather Research and Forecasting (WRF) Model. However, this work was limited by focusing on a single idealized squall-line simulation. Song and Roh (2021, hereafter SR21) developed a more general RT emulator for use in the Korea Local Analysis and Prediction System (KLAPS), an operational version of the WRF used by the KMA. When tested online in KLAPS, the NN produced similar instantaneous temperature and precipitation fields to the original RRTMG-K parameterization, suggesting that the NN may be suitable for operational use. Kim and Song (2022, hereafter KS22) used automatic hyperparameter tuning⁴ to find the best learning rate and training-batch size for the same KLAPS application, improving the performance of the NN further. Last, researchers at the ECMWF are currently working to integrate NN-based RT emulators into an operational model, namely, the Integrated Forecasting System (Chantry et al. 2022, 2023).

Some groups have used NNs to emulate only the gas-optics step of an RT model. Gas optics maps the physical/chemical state of the atmosphere to a profile of spectral optical depths, and the solver—the second and last step of an RT model—maps the optical depths to heating rates and fluxes

³ Specifically, version 2.7.1 of the shortwave model, covering the 0.2–12.2- μm band, and version 3.3 of the longwave model, covering the 3.07–1000- μm band.

⁴ A hyperparameter is an NN parameter that, unlike the weights and biases, cannot be adjusted during training. A hyperparameter must be tuned by trial and error, i.e., training many NNs with different values.

(Veerman et al. 2020). Specifically, gas optics converts temperature, pressure, and chemical concentrations into quantities that directly determine how much radiation is emitted, absorbed, and scattered in different directions (Veerman et al. 2020). Gas optics is an empirical algorithm in many RT models, relying on observations stored in large lookup tables, whereas the RT solver is a physical algorithm, relying on well-known equations. Because large lookup tables are computationally slow, gas optics is ripe for acceleration by NNs; because gas optics is already empirical, acceleration by NNs does not remove physical knowledge from the RT model. Ukkonen et al. (2020) emulated the gas-optics scheme in the RRTMGP and found that at most locations on Earth, the emulator introduces an RMSE of $<0.5 \text{ W m}^{-2}$ in fluxes and $<0.1 \text{ K day}^{-1}$ in heating rates for both the shortwave and longwave. Veerman et al. (2020) also emulated gas optics in the RRTMGP, obtaining similar results. Stegmann et al. (2022) emulated gas absorption in the Community Radiative Transfer Model, which is used in the observation operator for satellite data assimilation. Last, Ukkonen (2022) tested the use of NNs for three different emulation tasks: only the gas-optics scheme, only the reflectance-transmission calculations in the RT solver, and the full RT model. They found that replacing only the gas-optics scheme leads to the most accurate emulation, followed by replacing the full RT model. However, this study is limited by focusing only on shortwave RT for cloudy profiles. Geiss et al. (2022) emulated the aerosol-optics scheme of an RT model, using NNs with novel architectures, and found that connections between nonadjacent NN layers—which are uncommon in the literature—yielded the best performance.

NNs have additionally been used to simulate three-dimensional RT (Meyer et al. 2022; Yang et al. 2022) and hyperspectral RT (Le et al. 2020). Also, one study (Liu et al. 2020) has explored the effect of NN architectural complexity on RT accuracy. They compared fully connected and convolutional NNs,⁵ finding that convolutional NNs achieve slightly better performance but not enough to justify the added computational cost. However, they focused only on longwave RT in clear-sky conditions, and their errors were quite large (e.g., heating-rate errors often $\gg 1 \text{ K day}^{-1}$ near the surface). L21 explored U-net (Ronneberger et al. 2015) and U-net++ models (Zhou et al. 2020), convolutional NNs designed for image-to-image translation. In offline evaluation, they found that U-net++ models outperform fully connected NNs in general and outperform traditional U-nets for profiles with multilayer cloud, where RT is the most complex. See their supplemental section Cd for this architectural comparison.

In this work we use NNs—specifically, the U-net++ and U-net3+ architectures—to emulate the full RRTM. “Full” means that we emulate everything: both gas optics and the RT solver, for both the shortwave and longwave, including all

predictor variables. This contrasts with L21, where we emulated a simplified shortwave RRTM without aerosols, trace gases, or information on the particle size distribution (PSD) of hydrometeors. Our eventual aim is to integrate the NN-based emulators into the GFSv16, a global model with hybrid pressure–sigma coordinates in the vertical. Thus, we train the NNs with GFSv16 data from locations around the globe on the native pressure–sigma grid—in contrast to L21, we trained with data from 30 sites in the Northern Hemisphere on a standard height grid.

2. Data

This section discusses predictor (input) variables and target (output) variables. The RRTM and the NNs we use to emulate the RRTM have the same target variables and mostly the same predictor variables; the NNs have two extra predictor variables, as discussed in section 2a. Most predictor variables come from the GFSv16, but some are synthetic, because they are difficult to observe and not available in the GFSv16 output files. Because the NNs are built to emulate the RRTM, target variables produced by the RRTM are considered ground truth—“labels” in ML terminology.

a. GFSv16-based predictors

The GFSv16 is a global, nonhydrostatic, operational model with 0.25° horizontal spacing and 127 vertical levels in hybrid pressure–sigma coordinates, extending to the mesopause at $\sim 80 \text{ km}$ above mean sea level.⁶ We have obtained 0000 UTC model runs from the National Environmental Security Computing Center’s (NESCC) High-Performance Storage System (HPSS). The HPSS archive contains most days from 1 September 2018 to 23 December 2020 and forecast lead times of {0, 6, 12, 18, 24, 30, 36} h. We extract 6-, 12-, 18-, 24-, 30-, and 36-h forecast profiles (columns) from locations around the globe. Specifically, for each model solution (i.e., each combination of initialization time and valid time), we randomly select 4000 grid points from the global grid. We extract all predictor variables used by the RRTM that are in GFSv16 output files, listed in Table 1. We also extract two extra variables—the height thickness and pressure thickness of each layer—for use by the NNs but not the RRTM. For the work in L21, where all profiles have the same physical height grid (i.e., the k th pixel always corresponds to the same height in meters), the thickness variables were not necessary. But for the current work, where all profiles have a different physical height grid due to the hybrid coordinates, we found that the thickness variables improve RT estimation by the NNs. These variables are important because they tell the NNs how much “stuff” is in each layer—i.e., how much air there is to heat and how many other molecules there are to interact with radiation, which cannot be determined from molecular concentrations alone.

⁵ Fully connected (or “dense”) NNs treat the predictor variables as independent scalars, while convolutional NNs treat the predictors as images. Thus, convolutional NNs can leverage spatial structures in gridded data, while fully connected NNs cannot. While convolutional NNs are typically applied to 2D or 3D images, they can be applied just as easily to 1D “images”—such as the vertical profiles in this study—and leverage spatial structures therein.

⁶ See 2021 update here: https://www.emc.ncep.noaa.gov/emc/pages/numerical_forecast_systems/gfs/documentation.php.

TABLE 1. Description of GFSv16-based predictor variables. “Vector?” asks whether the variable is a profile or a scalar, and “AGL” = above ground level. Downward LWP at height z is LWC integrated from the top of the profile down to z , and upward LWP at height z is LWC integrated from the bottom of the profile up to z . The definitions of downward IWP, upward IWP, downward WVP, and upward WVP are analogous.

Variable	Units	Predictor for shortwave RT?	Predictor for longwave RT?	Vector?
Solar zenith angle	°	✓		
Surface albedo	—	✓		
Surface temperature	K		✓	
Surface emissivity	—		✓	
Temperature	K	✓	✓	✓
Pressure	Pa	✓	✓	✓
Specific humidity	kg kg ⁻¹	✓	✓	✓
Relative humidity	—	✓	✓	✓
Liquid water content (LWC)	kg m ⁻³	✓	✓	✓
Ice water content (LWC)	kg m ⁻³	✓	✓	✓
Downward liquid water path (LWP)	kg m ⁻²	✓	✓	✓
Downward ice water path (IWP)	kg m ⁻²	✓	✓	✓
Downward water vapor path (WVP)	kg m ⁻²	✓	✓	✓
Upward LWP	kg m ⁻²	✓	✓	✓
Upward IWP	kg m ⁻²	✓	✓	✓
Upward WVP	kg m ⁻²	✓	✓	✓
O ₃ mixing ratio	kg kg ⁻¹	✓	✓	✓
Height	m AGL	✓	✓	✓
Height thickness	m	✓	✓	✓
Pressure thickness	Pa	✓	✓	✓

b. Synthetic predictors

Predictors not in GFSv16 output files are listed in Table 2. We create synthetic data for these predictors, which fall into three categories: particle sizes, trace gases, and aerosols.

1) PARTICLE SIZES

The two relevant variables are ice effective radius ($r_{\text{eff}}^{\text{ice}}$) and liquid effective radius ($r_{\text{eff}}^{\text{liq}}$), both summaries of the PSD (Mitchell et al. 2011). To create a synthetic profile of $r_{\text{eff}}^{\text{ice}}$, we apply the following equation from Mishra et al. (2014), their Fig. 6b) independently to each height in the profile:

$$r_{\text{eff}}^{\text{ice}} = 86.73 \mu\text{m} + (1.07 \mu\text{m } ^\circ\text{C}^{-1})T, \quad (1)$$

where T is the temperature (°C) and each height has a different temperature (Fig. 1a). After Eq. (1), we apply two types of noise to the profile: bulk noise, which shifts the whole

profile to higher or lower values, and structure noise, which changes the structure of the profile (Fig. 1b). For bulk noise, we multiply the whole $r_{\text{eff}}^{\text{ice}}$ profile by $1 + \epsilon_b$, where ϵ_b is a random variable drawn from a normal distribution with mean = 0 and standard deviation = 0.5, denoted as $\mathcal{N}(0, 0.5)$. In other words, the standard deviation of bulk noise is 50% of the value generated by Eq. (1). For structure noise, we multiply the $r_{\text{eff}}^{\text{ice}}$ value at every height by $1 + \epsilon_s$, where ϵ_s is drawn anew at every height from $\mathcal{N}(0, 0.05)$. After adding noise, we bound $r_{\text{eff}}^{\text{ice}}$ values to the range [17.18, 65.33] μm , which is the same as bounding temperature to [-65, -20]°C, the range of validity for Eq. (1). See Fig. 1c.

To create a synthetic profile of $r_{\text{eff}}^{\text{liq}}$, we start with the distribution discovered by Miles et al. (2000). They found that $r_{\text{eff}}^{\text{liq}}$ roughly follows the distribution $\mathcal{N}(6, 1) \mu\text{m}$ over land and $\mathcal{N}(9.5, 1.2) \mu\text{m}$ over ocean. See Fig. 1d. However, using this information alone would lead to constant $r_{\text{eff}}^{\text{liq}}$ profiles, which are unrealistic. Thus, we add structure noise to each profile, using the same method as for $r_{\text{eff}}^{\text{ice}}$. See Fig. 1e.

TABLE 2. Description of synthetic predictor variables.

Variable	Units	Predictor for shortwave RT?	Predictor for longwave RT?	Vector?
Aerosol single-scattering albedo	—	✓		
Aerosol asymmetry parameter	—	✓		
Aerosol extinction coefficient	m ⁻¹	✓		
Liquid effective radius	m	✓	✓	✓
Ice effective radius	m	✓	✓	✓
N ₂ O concentration	ppmv	✓	✓	✓
CH ₄ concentration	ppmv	✓	✓	✓
CO ₂ concentration	ppmv	✓	✓	✓

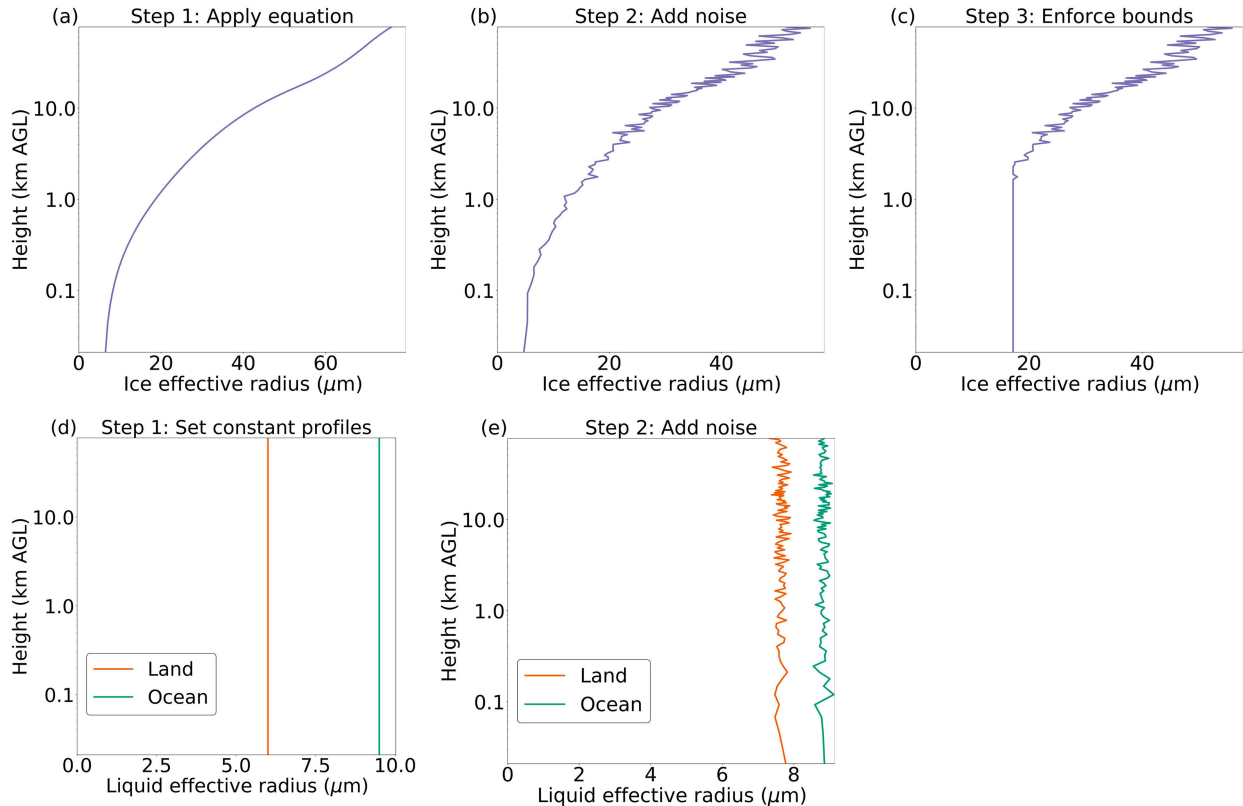


FIG. 1. Procedure for creating synthetic profiles of (a)–(c) ice effective radius and (d),(e) liquid effective radius.

2) TRACE GASES

For trace gases not in the GFSv16 output files— N_2O , CH_4 , and CO_2 —we use canonical profiles provided by Anderson et al. (1986). There is one canonical profile for each gas and each standard atmosphere, the latter defined in Table 3. For example, the five canonical N_2O profiles are shown in Fig. 2a. As for r_{eff}^{ice} , we add both bulk and structure noise to each profile of trace gas concentrations. We use the same noise distributions as for r_{eff}^{ice} . See Fig. 2b.

Note that the values provided in Anderson et al. (1986) are outdated, corresponding to a past climate. However, by

adding noise we sample a wide range of atmospheric conditions, corresponding to both present-day and hypothetical future climates. For example, online supplemental Fig. S3 shows that our dataset includes many CO_2 concentrations well above the present-day value of ~ 412 ppm.

3) AEROSOLS

Due to its complexity, we have relegated our method for creating synthetic aerosol variables—single-scattering albedo (SSA), asymmetry parameter, and extinction coefficient—to supplemental section 1.

c. Target variables

We run the shortwave and longwave RRTM separately for each profile. The target variables are those needed by an NWP model from its embedded RT model: a profile of heating rates (HR), surface downwelling flux (F_{down}^{sw}), top-of-atmosphere upwelling flux (F_{up}^{TOA}), and net flux (F_{net}). All four of these variables have both a shortwave and a longwave version. In machine learning the goal is often to improve accuracy, but our goal is to improve efficiency—i.e., to accelerate the RRTM—while emulating it as faithfully as possible. This means that we treat the RRTM as a perfect model, considering its HRs and fluxes to be the correct answers. Although the RRTM is imperfect, its errors are quite small, at less than 0.1 K day^{-1} for HRs and less than 1 W m^{-2} for fluxes (Iacono et al. 2008).

TABLE 3. Definition of standard atmospheres. The categorization is mutually exclusive and collectively exhaustive, i.e., every profile is assigned to exactly one of the five standard atmospheres.

Standard atmosphere	Months	Latitudes
Midlatitude summer	May–October	[20, 65]°N
Midlatitude summer	November–April	[20, 65]°S
Midlatitude winter	November–April	[20, 65]°N
Midlatitude winter	May–October	[20, 65]°S
Polar summer	May–October	[65, 90]°N
Polar summer	November–April	[65, 90]°S
Polar winter	November–April	[65, 90]°N
Polar winter	May–October	[65, 90]°S
Tropical	All	[−20, 20]°N

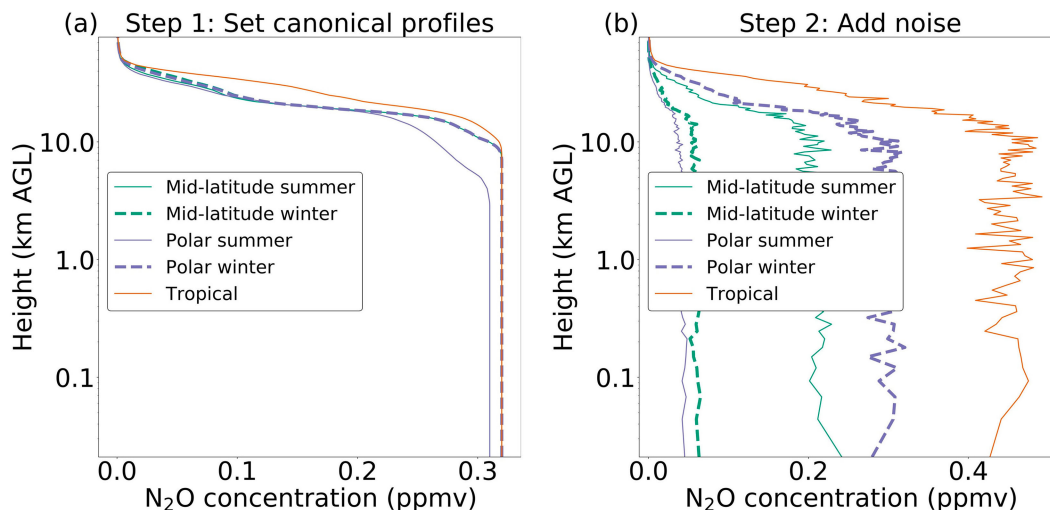


FIG. 2. Procedure for creating synthetic profiles of trace gas concentration—in this example, N_2O concentration.

d. Preprocessing

We apply two types of preprocessing to the data: splitting and normalization. As in L21, we use isotonic regression (IR) to bias correct the NNs, which requires a separate training set. Thus, we split the data into four temporally independent subsets: NN training, IR training, validation, and testing (Table 4). Each subset covers locations around the globe during all seasons. For normalization, we use the same methods described in section 3b of L21, except that we do not normalize any target variables. In L21 we normalized the flux variables, but we have since found that this has a deleterious effect on the quality of NN predictions.

3. Deep learning methods

This section provides a minimal background on the NN architectures used in L21, followed by a more extensive background on the architectures new to the current work, and finally information on the loss functions used to train NNs.

a. U-net and U-net++ without deep supervision

L21 considered two NN architectures, namely, the U-net and U-net++, for shortwave RT. They found that the U-net++ outperforms the U-net in situations with multilayer cloud (their supplemental section Cd), which are the most complex situations for RT and also vitally important for weather/climate prediction. In this article we consider the U-net++ architecture and a new architecture called U-net3+. L21 contains a detailed background on the U-net and U-net++ (their section 2), and we attempt to reproduce as little of this background as possible—only that which is necessary for understanding the current article.

The U-net (Ronneberger et al. 2015) is a type of NN designed for making predictions on a spatial grid, often called “image-to-image translation” in the ML literature. U-nets are typically applied to images with two or three spatial dimensions, but in our case the “images” are vertical profiles, containing only one spatial dimension. The task is to translate a $127 \times M$ image of predictors (M , the number of variables, is

TABLE 4. Partitioning of data into temporally independent subsets. “SW” = shortwave; “LW” = longwave; and “sample size” = number of profiles. SW and LW sample sizes are different because the SW radiation scheme (RRTM or NN-based emulator) is not run when the sun is below the horizon, i.e., when solar zenith angle $> 90^\circ$. Also, “number of days” \neq length of “time period,” because some days are missing from the archive.

Data subset	Time period	No. of days	SW sample size	LW sample size
NN training	1 Sep 2018–21 Dec 2019	237	873 086	3 503 226
IR training	24–30 Dec 2019, 3–9 Feb 2020, 15–21 Mar 2020, 26 Apr–2 May 2020, 7–13 Jun 2020, 18–24 Jul 2020, 28 Aug–3 Sep 2020, 10–16 Oct 2020, 21–27 Nov 2020	63	213 275	939 181
Validation	2–15 Jan 2020, 12–25 Feb 2020, 24 Mar–6 Apr 2020, 5–18 May 2020, 16–29 Jun 2020, 27 Jul–9 Aug 2020, 6–19 Sep 2020, 19 Oct–2 Nov 2020, 30 Nov–13 Dec 2020	126	479 806	1 934 460
Testing	18–31 Jan 2020, 28 Feb–12 Mar 2020, 9–22 Apr 2020, 22 May–4 Jun 2020, 2–15 Jul 2020, 12–25 Aug 2020, 22 Sep–7 Oct 2020, 5–18 Nov 2020, 16–23 Dec 2020	120	474 726	1 929 078

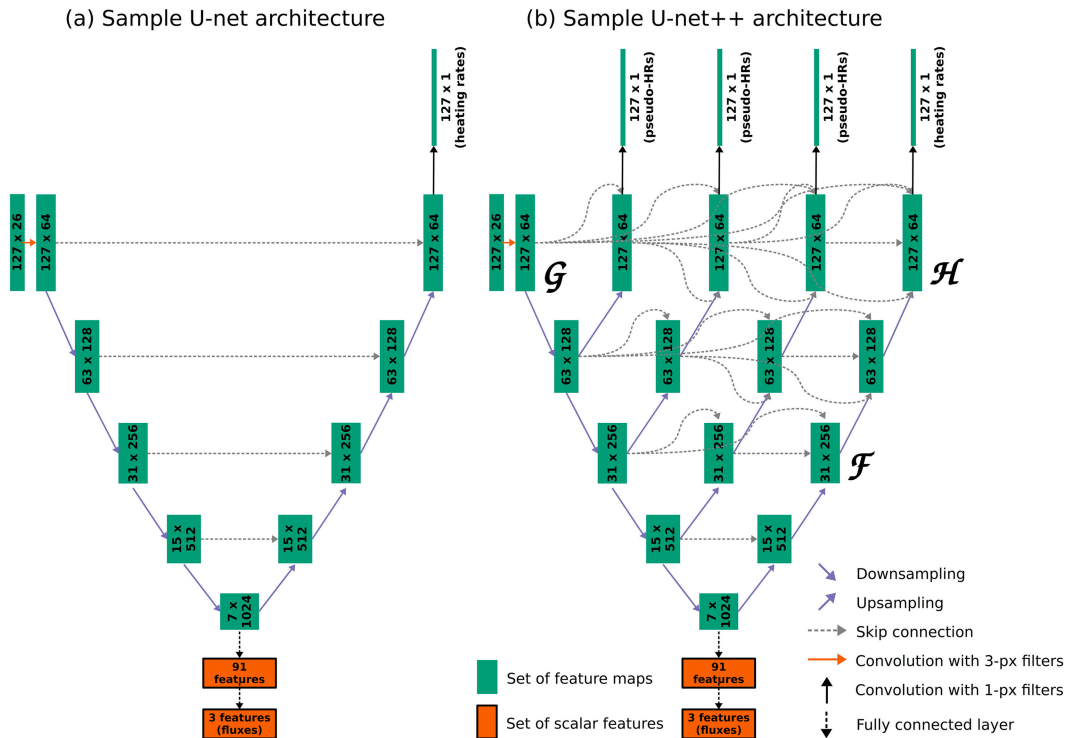


FIG. 3. Sample architectures for (a) U-net and (b) U-net++. Labels \mathcal{F} , \mathcal{G} , and \mathcal{H} are referred to in the main text. Actual models used in this study differ in the number of channels and depth (number of encoder/decoder layers, i.e., number of horizontal rows in this figure). For each set of feature maps (green box), the two dimensions are number of heights and channels, respectively. When the U-net++ is trained without deep supervision, all feature maps labeled “pseudo-HRs” go away, along with the arrows pointing to them. In the remaining discussion, let K be the number of convolutional layers per block, a user-chosen hyperparameter. Each orange “convolution” arrow corresponds to K convolutional layers with three-pixel filters; each “downsampling” arrow corresponds to K convolutional layers with three-pixel filters, followed by a maximum-pooling layer with a two-pixel window; each “upsampling” arrow corresponds to an up-sampling layer with a two-pixel window, followed by a convolutional layer with three-pixel filters; each “skip connection” arrow includes K convolutional layers with three-pixel filters; each black “convolution” arrow corresponds to one convolutional layer with one-pixel filters; and last, each “fully connected layer” arrow corresponds to one fully connected layer.

different for longwave versus shortwave RT) into a 127×1 image of HRs.⁷

U-nets contain four key components (Fig. 3a): convolutional layers, pooling (downsampling) layers, upsampling layers, and skip connections. The role of the convolutional layers is to detect spatial and multivariate features—i.e., features including many pixels and predictor variables—using convolutional filters with weights optimized during training to detect the most useful features for prediction. The role of the pooling and upsampling layers is to change the resolution of the feature maps—a “feature map” being either the original or a transformed version of the predictors—so that convolutional layers at different depths in the network can detect features at different spatial scales. The role of the skip connections is to preserve high-resolution information—i.e., to carry through the network high-resolution information that

has not been degraded by downsampling, a lossy operation that cannot be fully reversed by upsampling. The left side of the U-shaped network (Fig. 3a) is the encoder side, where the predictors are converted to feature maps with decreasing spatial resolution (fewer height levels) and increasing spectral resolution (more channels). The right side is the decoder side, where feature maps are upsampled and converted to the final prediction—an image of HRs. To make our networks also predict the three flux variables, which are scalars and not images, we attach fully connected layers to the deepest encoder layer, as shown in Fig. 3a. These are the layers used in fully connected NNs (chapter 6 of Goodfellow et al. 2016), which are still a popular choice for scalar data.

The U-net++ (Zhou et al. 2020) contains more skip connections than the U-net, which more effectively preserve small-scale features such as cloud boundaries, leading to better predictions for multilayer cloud in L21. The U-net3+ (Huang et al. 2020) contains even more skip connections than the U-net++, so we hypothesize that the U-net3+ will perform even better in situations with multilayer cloud and

⁷ There is a second learning task, which involves image-to-scalar translation—namely, to translate the same $127 \times M$ image of predictors into three flux components.

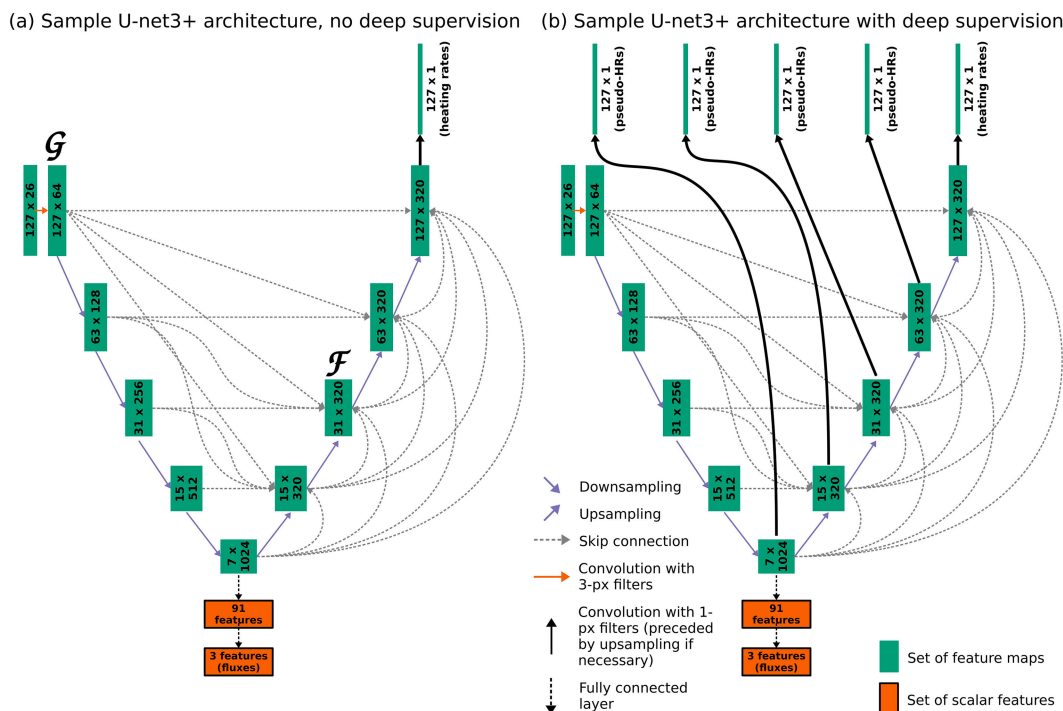


FIG. 4. Sample architectures for U-net3+ (a) without and (b) with deep supervision. Labels \mathcal{F} and \mathcal{G} are referred to in the main text. Actual models used in this study differ in the number of channels and depth. Formatting is explained in the caption of Fig. 3, except that the solid black arrows are slightly different in this figure. The solid black arrow pointing to actual HRs (top right) corresponds to one convolutional layer with one-pixel filters, while a solid black arrow pointing to pseudo-HRs corresponds to an upsampling layer followed by a convolutional layer with one-pixel filters.

perhaps overall. Also, the U-net++ and U-net3+ may be trained with deep supervision, which was not used in L21.

b. U-net++ with deep supervision

When an NN is trained without deep supervision, the loss function optimized by the NN compares the ground truth (here, a length-127 profile of HRs) only to the final prediction, i.e., output from the last NN layer. With deep supervision, the ground truth is also compared to intermediate representations, i.e., layer outputs that are ultimately transformed to the final prediction. Zhou et al. (2020) found that deep supervision improves image segmentation for phenomena that occur at different scales, such as lung nodules. We hypothesize that deep supervision will also improve RT estimation, since relevant features for RT estimation also occur at different scales—e.g., cloud depths range from $\mathcal{O}(10)$ m to $\mathcal{O}(10)$ km.

Figure 3b shows a sample U-net++ architecture with and without deep supervision. The only difference is that deep supervision requires extra convolutional layers—those producing pseudo-HRs—to transform the intermediate representations from many channels to one channel. With deep supervision, all four outputs (the three pseudo-HR profiles and the actual-HR profile) are produced; without deep supervision, only one output (the actual-HR profile) is produced. For details on the loss function, which compares both pseudo-HRs and actual HRs to

the ground truth, see section 3d. Note that deep supervision is applied only to the spatial outputs (HRs) and not the scalar outputs (fluxes). Deep supervision was invented for spatial data, and there is no clear analog for scalars.

c. U-net3+ with and without deep supervision

The U-net3+ has one property that distinguishes it from the U-net++, namely, full-scale skip connections. Full-scale skip connections pass information from all scales to each decoder layer, whereas skip connections in the U-net++ pass information from only two scales to each decoder layer. For example, in the U-net++ shown in Fig. 3b, the feature maps labeled \mathcal{F} combine information from the same scale (other feature maps with 31 heights) and the next-largest scale (feature maps with 15 heights). But in the U-net3+ shown in Fig. 4a, the feature maps labeled \mathcal{F} combine information from equal and smaller scales (feature maps with ≥ 31 heights) on the encoder side, as well as information from larger scales (feature maps with < 31 heights) on the decoder side.

Stated differently, full-scale skip connections more effectively carry high-resolution information through the network. For example, the feature maps labeled \mathcal{G} (in both Figs. 3b and 4a) contain information at the smallest scale that has not been degraded by downsampling. In the U-net++ (Fig. 3b), skip connections carry this information to only one level on the

decoder side, namely, the feature maps labeled \mathcal{H} . Other levels on the decoder side cannot access the undegraded high-resolution information in \mathcal{G} . But in the U-net3+ (Fig. 4a), full-scale skip connections carry the information in \mathcal{G} to all levels on the decoder side, allowing this information to be used in decoded feature maps at all resolutions.

Figures 4a and 4b show how to add deep supervision to the U-net3+ architecture. For the U-net3+, deep supervision requires two architecture changes. The first is extra convolutional layers to reduce the number of channels to one (pseudo-HR), as in the U-net++. The second is extra upsampling layers to increase the number of heights to 127.

d. Loss function

In machine learning, the standard loss function for regression tasks—where the model predicts a continuous value instead of a category—is the mean squared error (MSE). However, in L21 we found that using the MSE causes two problems. First, the MSE does not adequately emphasize large HRs, which are rare but important for weather/climate prediction, causing the NN to dramatically underpredict large HRs. Second, the MSE does not ensure that the following conservation law is respected:

$$F_{\text{net}}^{(b)} = F_{\text{down}}^{\text{sc}(b)} - F_{\text{up}}^{\text{TOA}(b)}, \quad (2)$$

where the superscript (b) denotes that all three variables must come from the same band, either shortwave or longwave. To remedy the first problem, we used the dual-weighted MSE (DWMSE) for HRs, which emphasizes cases with a large actual or predicted HR, “nudging” the NN to predict these cases correctly. See section 3c(2) of L21. To remedy the second problem, we used the basic MSE for flux variables *but* enforced the law of Eq. (2) inside the NN. See section 3c(1) of L21.

Because L21 is concerned with shortwave RT only, the present work requires two updates to the loss function. First, the weight in the DWMSE becomes the maximum of the *absolute* actual and predicted HRs, because although shortwave HR is always ≥ 0 , longwave HR may be negative (i.e., longwave cooling). Second, the flux law must be applied to both shortwave and longwave RT. The total loss function becomes the following:

$$\begin{aligned} \mathcal{L}^{(b)} = & \frac{1}{NH} \sum_{i=1}^N \sum_{j=1}^H \max\{|r_{ij}^{(b)}|, |\hat{r}_{ij}^{(b)}|\} [r_{ij}^{(b)} - \hat{r}_{ij}^{(b)}]^2 \\ & + \frac{1}{NM} \sum_{i=1}^N \sum_{k=1}^M [F_{ik}^{(b)} - \hat{F}_{ik}^{(b)}]^2, \end{aligned} \quad (3)$$

where N is the number of examples, $H = 127$ is the number of heights per example, $r_{ij}^{(b)}$ is the actual HR for the j th height in the i th example, $\hat{r}_{ij}^{(b)}$ is the corresponding prediction, $M = 3$ is the number of flux components, $F_{ik}^{(b)}$ is the actual value of the k th flux component in the i th example, and $\hat{F}_{ik}^{(b)}$ is the corresponding prediction. There is one version of Eq. (3) for the shortwave, where the superscript (b) is SW, and one version for the longwave.

TABLE 5. Experimental hyperparameters.

Hyperparameter	Values attempted
NN type	U-net++ without deep supervision, U-net++ with deep supervision, U-net3+ without deep supervision, U-net3+ with deep supervision,
NN depth	3, 4, 5
NN width	1, 2, 3, 4
Spectral complexity	4, 8, 16, 32, 64, 128

For NNs without deep supervision, Eq. (3) is the whole story. However, for NNs with deep supervision, the loss function includes extra terms for the pseudo-HRs. Specifically, the loss function becomes

$$\mathcal{L}_{\text{deep-sup}}^{(b)} = \mathcal{L}^{(b)} + \frac{1}{PNH} \sum_{p=1}^P \sum_{i=1}^N \sum_{j=1}^H \max\{|r_{ij}^{(b)}|, |\hat{r}_{pji}^{(b)}|\} [r_{ij}^{(b)} - \hat{r}_{pji}^{(b)}]^2, \quad (4)$$

where P is the number of layers with deep supervision and thus the number of pseudo-HR profiles, and $\hat{r}_{pji}^{(b)}$ is the pseudo-HR produced by the p th layer with deep supervision for the j th height in the i th example.

4. Experiment with neural networks of varying complexity

This section describes a hyperparameter-tuning experiment used to find the optimal level of NN complexity for estimating RT. We tune four hyperparameters: the NN type (U-net++ or U-net3+ with or without deep supervision), NN depth, NN width, and spectral complexity. NN depth is the number of encoder/decoder levels (e.g., all architectures shown in Figs. 3 and 4 have a depth of 4), NN width is the number of convolutional layers per set (K in the caption of Fig. 3), and spectral complexity is the number of feature maps produced by the first set of convolutional layers (e.g., all architectures shown in Figs. 3 and 4 have a spectral complexity of 64). Following common practice, we always double the number of feature maps with each downsampling operation. For example, Fig. 3 shows that with a depth of 4 and spectral complexity of 64, the deepest set of feature maps (i.e., that with the coarsest spatial resolution, designed to capture the largest-scale features) has 1024 feature maps. We chose to experiment with NN type so that we could try new methods (deep supervision and U-net3+) from the ML literature. We chose to experiment with the other three hyperparameters because they strongly control overall NN complexity, i.e., the number of trainable weights. As shown in supplemental Figs. S10 and S18, the number of trainable weights varies from $\mathcal{O}(10^5)$ to $\mathcal{O}(10^{8.5})$.

Table 5 lists the exact values attempted for each hyperparameter. We perform a grid search (section 11.4.3 of Goodfellow et al. 2016), training one NN for every combination of values, which leads to $4 \times 3 \times 4 \times 6 = 288$ NNs for each band (shortwave and longwave). Most constant hyperparameters (those

TABLE 6. Metrics used for model selection. “Column-averaged” = averaged over all 127 heights; “near-surface” = at the lowest grid level, which averages 21 m AGL; and “all-flux RMSE” is the square root of the MSE averaged over all three flux variables. Metrics computed on fog profiles are used only to evaluate longwave models, not shortwave models.

Set of profiles	Metrics used
All	Column-averaged HR DWMSE, column-averaged HR bias, near-surface HR DWMSE, near-surface HR bias, all-flux RMSE, net-flux RMSE, net-flux bias
Profiles with multilayer liquid-only cloud	Column-averaged HR DWMSE, column-averaged HR bias, near-surface HR DWMSE, near-surface HR bias, all-flux RMSE, net-flux RMSE, net-flux bias
Profiles with fog (longwave only)	Near-surface HR DWMSE, near-surface HR bias, all-flux RMSE, net-flux RMSE, net-flux bias

not varied during the experiment) are illustrated in Figs. 3 and 4. Constants not included in these figures are documented in supplemental Table S3.

a. Evaluation methods used for model selection

Model evaluation is a multifaceted problem, and there are many possible ways to choose the best model. Most hyperparameter experiments optimize one evaluation metric, often the loss function used for training. However, we care about several aspects of model performance. In previous work we have noticed that even when overall performance is acceptable, the following regime-based errors are unacceptably high:

- HR errors near the surface, especially in the longwave;
- flux and HR errors in profiles with multilayer liquid-only cloud, in both the shortwave and longwave;
- longwave HR errors near the surface in profiles with fog, i.e., cloud reaching the lowest grid level.

Thus, we use the metrics listed in Table 6, computed on validation data only, for model selection. Our choice of the best model is based on a subjective combination of these metrics.

b. Evaluation methods used for best models

As in L21, we evaluate the best models (shortwave and longwave) on the testing dataset as a whole and on meaningful subsets of the testing data. We split the testing data in four ways.

First, we split by cloud regime, because clouds add immense complexity to RT, making the process difficult to emulate, and can result in extreme HRs (large absolute values in both the shortwave and longwave), which are important for weather and climate. For a more detailed explanation of these effects, see section 5a of L21. We focus on liquid-only cloud, which we have found to have a greater effect on RT than ice-only, mixed-phase, or any-phase cloud. We define a liquid-only cloud layer as a contiguous set of model heights with liquid water content (LWC) $> 0 \text{ g m}^{-3}$, total liquid water path $\geq 25 \text{ g m}^{-2}$, and total ice water path = 0 g m^{-2} . As in L21, we define three cloud regimes, which are mutually exclusive and collectively exhaustive (MECE): no cloud, single-layer cloud, and multilayer cloud. For the longwave we add a fourth cloud regime—fog—defined as a cloud reaching the surface (i.e., LWC $> 0 \text{ g m}^{-3}$ at the lowest model height). Thus, cloud regimes for the longwave are not MECE, as every profile with fog is also a profile with single- or multilayer cloud. We include fog because it causes large longwave errors near the surface.

Second, we split the testing data by geographic location, specifically, on a global latitude–longitude grid with 5° spacing. This spacing highlights large RT errors due to features such as high terrain and persistent stratocumulus cloud. Third, for the shortwave model only, we split the testing data by aerosol optical depth (AOD) and solar zenith angle (SZA). In earlier work we found that shortwave errors increase with higher AOD, which adds complexity to RT, and lower SZA,⁸ which increases HRs and the frequency of extreme HRs. Fourth, for the longwave model only, we split the testing data by near-surface thermodynamics, specifically, temperature lapse rate (Γ_T^{sfc}) and humidity lapse rate (Γ_q^{sfc}). These are defined as

$$\begin{cases} \Gamma_T^{\text{sfc}} = \frac{T_1 - T_2}{z_2 - z_1}, \\ \Gamma_q^{\text{sfc}} = \frac{q_1 - q_2}{z_2 - z_1}, \end{cases} \quad (5)$$

where T_1 and T_2 are temperature (K) at the lowest and second-lowest model heights (sigma levels), respectively; q_1 and q_2 are specific humidity (kg kg^{-1}) at the same heights; and z_1 and z_2 are the corresponding physical heights (m AGL). Longwave RT near the surface is highly sensitive to the near-surface temperature and moisture profiles (Schmetz 1989). We also experimented with splitting by surface temperature and humidity, instead of their near-surface lapse rates, but found that lapse rates have a greater impact on longwave-RT errors.

We use several evaluation metrics and plotting tools, most of which are familiar to atmospheric scientists, such as the mean absolute error and bias (mean signed error). We also use the attributes diagram, which is a reliability curve with added reference lines (Hsu and Murphy 1986). However, we have adapted this plot for regression (predicting a continuous value, like flux in watts per square meter) instead of their typical use, which is binary classification (predicting the probability of an event). For readers interested in the details, see section 5a of L21. You can interpret the regression- and classification-based version of the attributes diagram in roughly the same way: the curve should be close to the

⁸ Lower SZA means that the sun is higher above the horizon. Specifically, SZA is 0° when the sun is directly overhead, and 90° when the sun is on the horizon.

diagonal reference line, indicating perfect reliability, and inside the shaded area, indicating a positive skill score. For the regression-based attributes diagram, this is the MSE skill score. A positive MSE skill score means that the NN model has a better MSE than the climatological model. The climatological model is a simple model that always predicts the climatological mean, estimated as the average in the training data. For example, if the mean net flux in the training data is 100 W m^{-2} , the climatological model will predict a net flux of 100 W m^{-2} for every case.

5. Results and discussion

We start with a brief discussion of the hyperparameter experiment (used to determine the best models), followed by a comparison of computing time between the RRTM and our NN-based emulators, then an in-depth discussion of the best shortwave model and best longwave model.

a. Hyperparameter experiment

Results are discussed briefly here and at length in supplemental section 3. For both shortwave and longwave RT, the most important hyperparameter is spectral complexity, while NN depth and width are of secondary importance. The better NNs have large spectral complexity, large depth, and small width. In other words, the better NNs are deep and narrow with many feature maps. For the other hyperparameter—NN type—we hypothesized that the U-net3+ architecture would outperform U-net++ (section 3a) and that NNs trained with deep supervision would outperform those with no deep supervision (section 3b). We are unable to confirm either hypothesis—deep supervision leads to *worse* performance, and architecture has little effect on performance. The best shortwave model—based on our subjective assessment of the metrics listed in Table 6—is a U-net++ with no deep supervision, depth of 3, width of 1, and spectral complexity of 128, leading to $10^{7.52}$ trainable weights. The best longwave model—again based on Table 6—is a U-net3+ with no deep supervision, depth of 5, width of 1, and spectral complexity of 64, leading to $10^{7.28}$ trainable weights. Therefore, the best models are on the high end of the overall-complexity range in our experiment, with number of weights ranging from $\mathcal{O}(10^5)$ to $\mathcal{O}(10^{8.5})$. This is because spectral complexity is the main control on both performance (allowing the models to represent and leverage many features of the input data) and number of weights (see supplemental Figs. S10 and S18).

b. Computing time

The original motivation for NNs was to decrease computing time. To this point, we have compared the wall-clock time of the RRTM and best NNs when run on the same hardware—i.e., one node with 24 CPUs and no GPUs—in stand-alone mode. See Table 7 for details. In summary, the shortwave RRTM (NN) processes 0.11 (843) profiles per second, resulting in a speedup factor of 7510; while the longwave RRTM (NN) processes 5.13 (460) profiles per second, resulting in a speedup factor of 90. Thus, we have accelerated the RRTM by orders of magnitude.

TABLE 7. Timing tests for the RRTM and NN-based emulators, based on the testing dataset. All computing times are given in wall-clock time. Because the RRTM is slower for cloudy profiles and faster for cloud-free profiles, the “time per profile” reported is an average over all atmospheric conditions represented in the dataset. Meanwhile, the NNs have constant computing time for each profile, regardless of atmospheric conditions.

Model	No. of profiles	Total time (s)	Time per profile (s)
Shortwave RRTM	472 412	4 207 793	0.11
Shortwave NN	474 726	563	843
Longwave RRTM	1 894 239	369 363	5.13
Longwave NN	1 929 078	4194	460

c. Best shortwave model

Figure 5 shows the overall performance—i.e., averaged over the whole testing set—of the best shortwave model. For all flux variables (Figs. 5a–c), the model is almost perfectly reliable (see overlap between reliability curve and diagonal reference line) and almost perfectly reproduces the observed distribution (see similarity between the two histograms). However, the model has slight conditional biases, namely, an overprediction of $\sim 10 \text{ W m}^{-2}$ for the highest $F_{\text{down}}^{\text{sc}}$ and $F_{\text{up}}^{\text{TOA}}$ predictions. In other words, when the model predicts an extremely large downwelling or upwelling flux, the prediction is slightly too extreme. However, these two biases offset in the calculation of F_{net} [Eq. (2)], resulting in near-zero bias for all predicted F_{net} values. The model has an absolute bias $< 0.1 \text{ K day}^{-1}$ for HR at every height (Fig. 5d); this suggests that it could be stably integrated into an NWP system such as the GFS (Iacono et al. 2008), as systematic errors for an RT parameterization are much more important than random errors (Pincus et al. 2003). The model has a substantially larger MAE than bias for HR at every height (Figs. 5d,e), which indicates that most of the model’s HR error is random instead of systematic. Both bias and MAE are largest in the upper stratosphere, where shortwave RT is dominated by O_3 absorption. The bias and MAE profiles in L21 were similar—even with a dataset that used a constant profile for trace gases such as O_3 —which suggests that O_3 absorption is a fundamentally difficult process to emulate. Since the average HR in the upper stratosphere is large (e.g., 21.6 K day^{-1} at 47 km AGL), the climatological model also has a large MAE here, so the NN’s spike in MAE translates to only a small dip in its MAE skill score (Fig. 5f). Last, the attributes diagram for HR (Fig. 5g) tells a similar story to those for the flux variables: the model is almost perfectly reliable and almost perfectly reproduces the observed distribution. However, the model has a slight positive bias ($\ll 1 \text{ K day}^{-1}$) for the highest predicted HR values.

Supplemental Figs. S22 and S23 are analogous to Fig. 5 but only for extreme cases—i.e., the 3% of testing profiles with the greatest height-maximum and height-averaged HR, respectively. Although errors are expectedly higher for the extreme cases, HR and flux predictions are still almost perfectly reliable and absolute HR bias is well below 0.1 K day^{-1} at almost every height.

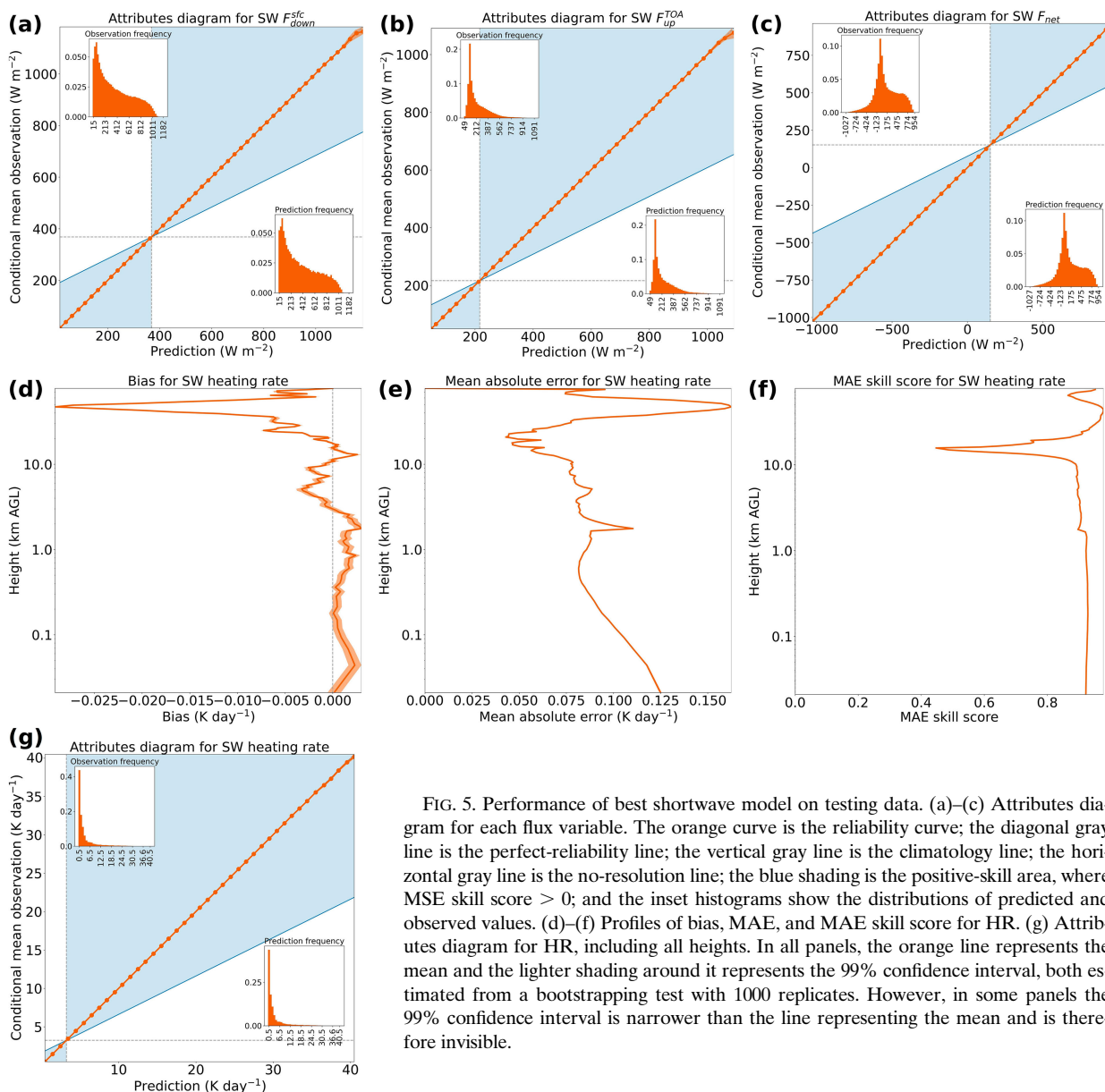


FIG. 5. Performance of best shortwave model on testing data. (a)–(c) Attributes diagram for each flux variable. The orange curve is the reliability curve; the diagonal gray line is the perfect-reliability line; the vertical gray line is the climatology line; the horizontal gray line is the no-resolution line; the blue shading is the positive-skill area, where MSE skill score > 0 ; and the inset histograms show the distributions of predicted and observed values. (d)–(f) Profiles of bias, MAE, and MAE skill score for HR. (g) Attributes diagram for HR, including all heights. In all panels, the orange line represents the mean and the lighter shading around it represents the 99% confidence interval, both estimated from a bootstrapping test with 1000 replicates. However, in some panels the 99% confidence interval is narrower than the line representing the mean and is therefore invisible.

Figure 6 shows the model's performance as a function of liquid-only cloud regime. Performance for other cloud phases (ice only, mixed phase, and any phase) is shown in supplemental Figs. S19–S21. The attributes diagram for each flux variable (Figs. 6a–c) tells a similar story to its cloud-agnostic analog (Figs. 5a–c): slight conditional bias for extreme predictions of F_{down}^{sfc} and F_{up}^{TOA} but with no absolute bias exceeding 20 W m^{-2} . The following discussion of error profiles for HR (Figs. 6d–f) focuses on the troposphere (below $\sim 15 \text{ km AGL}$), where shortwave heating is dominated by cloud rather than O_3 . In the bottom few 100 m, errors are largest for clear-sky profiles and smallest for cloudy profiles, because in cloudy profiles most of the incoming solar radiation has already been absorbed by clouds above, which leaves little shortwave

radiation in the bottom few 100 m, thus making shortwave RT an easier problem here. Meanwhile, in the troposphere above $\sim 1 \text{ km}$, errors are smallest for clear-sky profiles and largest for cloudy profiles, because this is the region where most clouds and their associated extreme HRs occur. Also, errors for multilayer cloud are greater than for single-layer cloud, because multilayer cloud produces nonlocal effects that are difficult to emulate. For example, consider a profile with two clouds of equal thickness and structure (i.e., equal series of LWC values), one based at 10 km AGL and the other based at 1 km AGL . The upper cloud will absorb most of the incoming solar radiation, leaving little shortwave radiation to be absorbed by the lower cloud; thus, the upper cloud will cause much larger HRs, even though the two clouds are identical except for

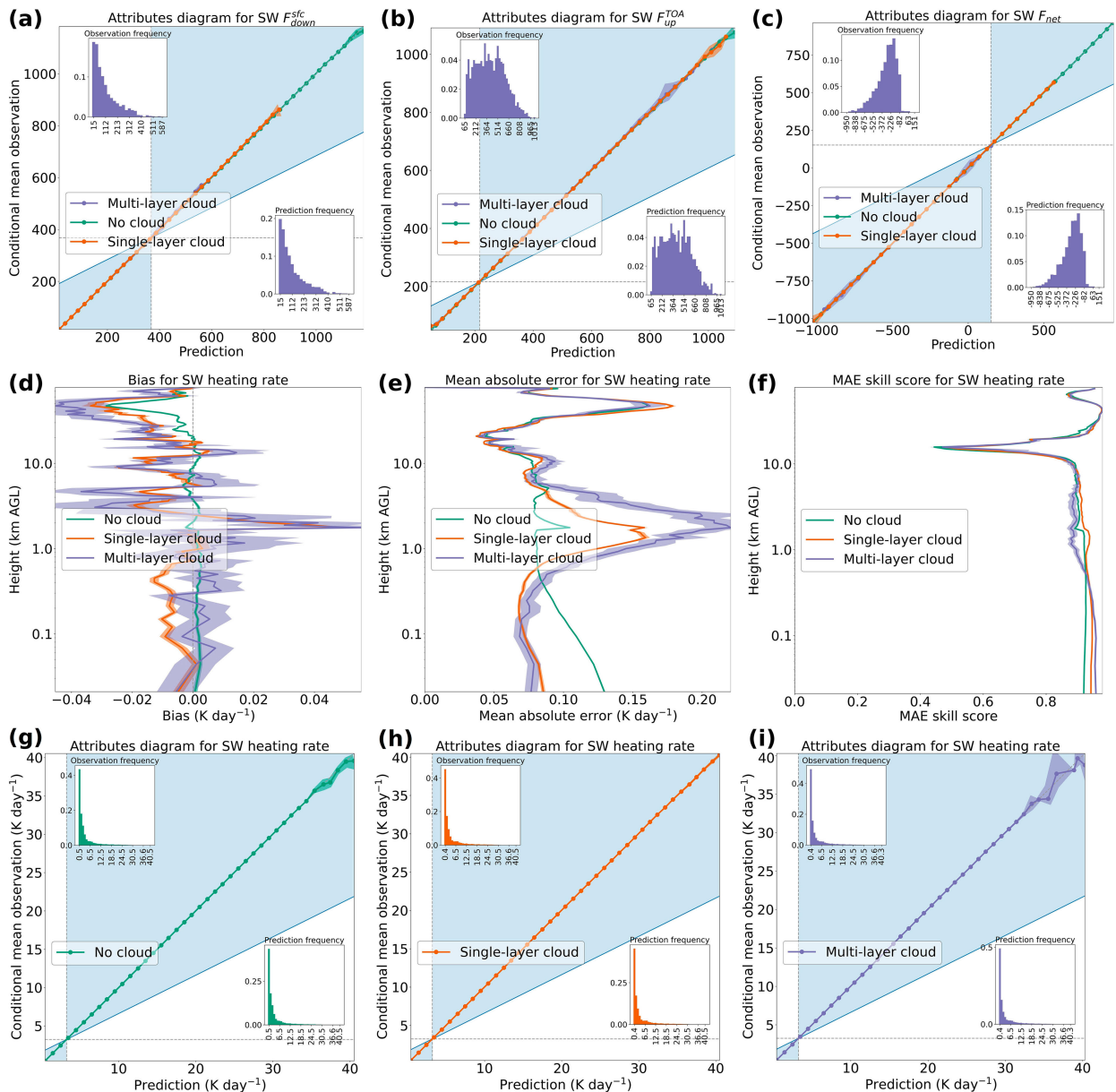


FIG. 6. Performance of best shortwave model on testing data, separated by liquid-only cloud regime. (a)–(c) Attributes diagram (formatting explained in the caption of Fig. 5) for each flux variable. The inset histograms are based only on cases with multilayer cloud. (d)–(f) Profiles of bias, MAE, and MAE skill score for HR. (g) Attributes diagram for HR, including all heights, only for cases with no cloud (89.67% of the testing data). (h) As in (g), but for single-layer cloud (9.98% of the testing data). (i) As in (g), but for multilayer cloud (0.35% of the testing data). In all panels, the green, orange, and purple lines represent the mean and the lighter shading around them represents the 99% confidence interval, both estimated from a bootstrapping test with 1000 replicates.

location. This is a nonlocal effect, as the two clouds are far (more than a few grid cells) apart. Last, the attributes diagrams for HR (Figs. 6g–i) tell a similar story to their cloud-agnostic analog (Fig. 5g): an overall positive bias for the highest predicted HR values and near-zero bias for all other values. However, this positive bias is largest for multilayer cloud (up to ~ 2 K day⁻¹)—likely due to a small sample size for the highest predicted HR values, indicated by the wide confidence intervals in Fig. 6i.

Figure 7 shows the model’s performance as a function of location. The column-averaged MAE for HR (Fig. 7a) is mostly between⁹ 0.07 and 0.11 K day⁻¹; it exceeds 0.11 K day⁻¹ at a few locations, notably the Tibetan Plateau and east Antarctica.

⁹ Henceforth, “mostly between” corresponds to the middle 95% of the distribution, i.e., the 2.5th–97.5th percentiles. However, note that the color bar in each panel shows 100% of the distribution, ranging from the minimum to the maximum.

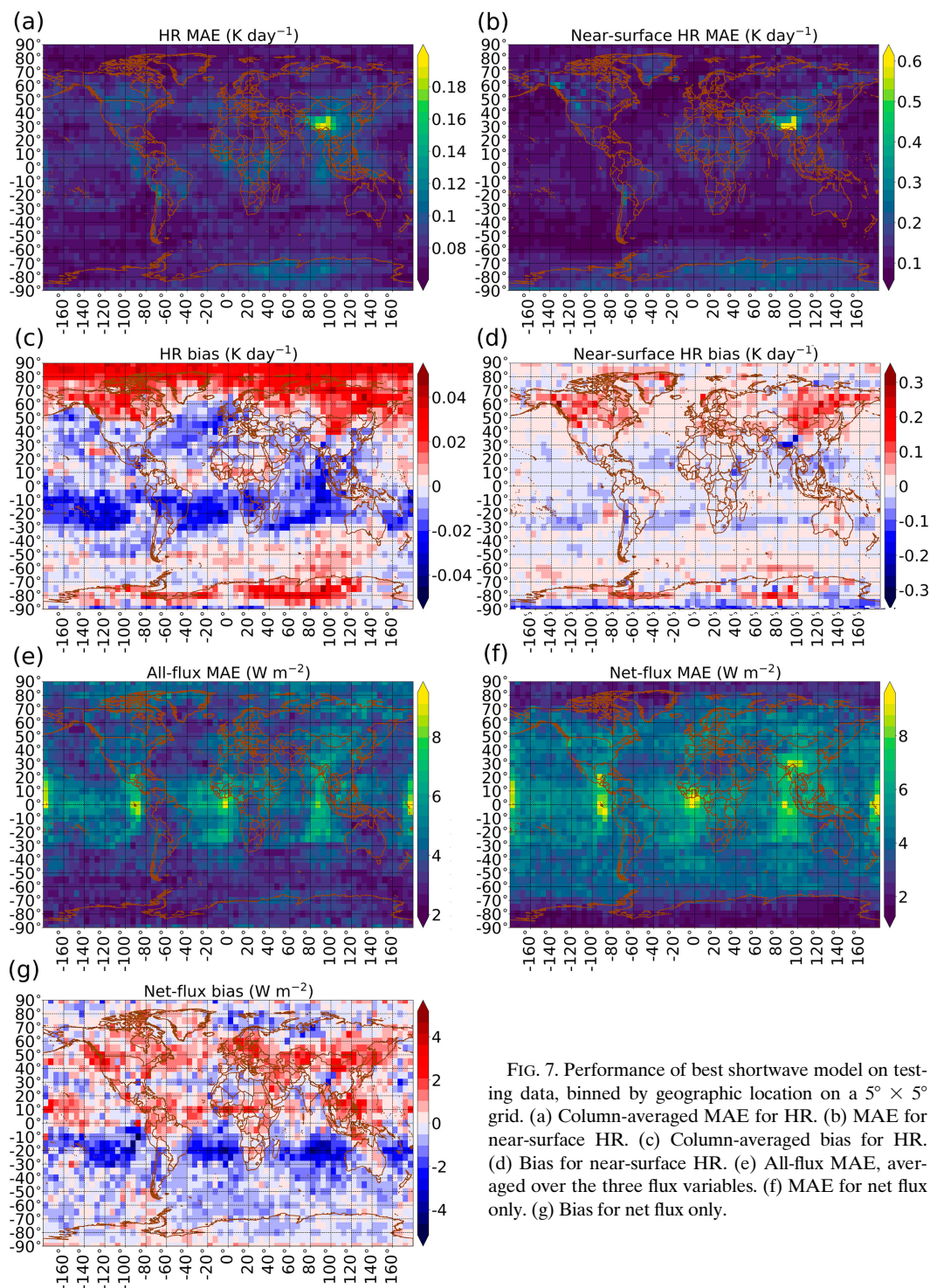


FIG. 7. Performance of best shortwave model on testing data, binned by geographic location on a $5^\circ \times 5^\circ$ grid. (a) Column-averaged MAE for HR. (b) MAE for near-surface HR. (c) Column-averaged bias for HR. (d) Bias for near-surface HR. (e) All-flux MAE, averaged over the three flux variables. (f) MAE for net flux only. (g) Bias for net flux only.

The MAE for near-surface HR (Fig. 7b) is larger—mostly between 0.07 and 0.23 K day^{-1} , exceeding 0.23 K day^{-1} at a few locations, again notably Tibet and east Antarctica. The two locations have very high surface elevation and albedo, the latter due to snow/ice cover. High elevation decreases atmospheric

thickness and therefore increases near-surface HR; high albedo decreases near-surface HR; and both extremes are globally rare, causing high model error under these extremes. Many error metrics (Figs. 7a,b,d,f) are especially large over the Tibetan Plateau, as it is the largest and highest plateau in

the world, thus exacerbating both the thickness and albedo effects. The column-averaged bias for HR (Fig. 7c) is mostly between -0.02 and $+0.03$ K day^{-1} , with absolute bias not exceeding 0.05 K day^{-1} at any location. The bias for near-surface HR (Fig. 7d) is larger—mostly between -0.09 and $+0.09$ K day^{-1} , with absolute value exceeding 0.09 K day^{-1} over high-latitude continents such as Canada, Siberia, and Antarctica. The all-flux MAE (Fig. 7e) is mostly between 2.5 and 6.4 W m^{-2} , exceeding 6.4 W m^{-2} mainly in the Southern Hemisphere stratocumulus regions. These are regions of semipersistent stratocumulus cloud in the subtropics off the west coast of a continent—including South America, southern Africa, and Australia (Fig. 6 of Neubauer et al. 2014). The net-flux MAE (Fig. 7f) follows a similar pattern to the all-flux MAE. Last, the net-flux bias (Fig. 7g) is mostly between -2.2 and $+2.0$ W m^{-2} , with mostly negative bias in the Southern Hemisphere and positive bias in the Northern Hemisphere.

Supplemental Fig. S24 is analogous to Fig. 7 but shows relative, instead of raw, errors. For example, “relative net-flux MAE” at grid point P is (raw net-flux MAE at P)/(mean observed net flux at P). We make two observations from the two figures. First, for column-averaged HR MAE (Fig. S24a and Fig. 7a), the highest relative errors are collocated with the highest raw errors—in Tibet and east Antarctica. This indicates that shortwave HR is *fundamentally* harder to predict at said locations—i.e., these maxima in HR error are not just caused by maxima in HR itself. Second, for all other error metrics (Figs. S24b–g and Figs. 7b–g), the largest relative errors occur at polar latitudes, where raw errors are small. Polar latitudes receive little solar radiation, leading to small shortwave HRs and fluxes, so a small raw error translates to a large relative error. Supplemental Fig. S25 is another variant of Fig. 7, but showing errors for individual flux variables instead of averaging to produce all-flux quantities. The main conclusion from this figure is that $F_{\text{down}}^{\text{swc}}$ errors are worst at the low latitudes, including in the stratocumulus cloud regions, while $F_{\text{up}}^{\text{TOA}}$ errors are worst at the high latitudes.

Figure 8 shows case studies from two regions with high model error: Tibet (Figs. 8a–d) and east Antarctica (Figs. 8e–h). To select these case studies, we first plotted 400 random profiles—200 from each region—and then manually selected 4 profiles that are representative of the original 400. In the following conclusions, although we reference Fig. 8, we have ensured that they represent most of the original 400 profiles as well. First, Tibet experiences a lot of cloud, often complex mixtures of liquid and ice. Second, east Antarctica also experiences a lot of cloud, often ice cloud reaching the surface as fog. Third, although the model matches the shape of the HR profile well, it often misses extreme HRs associated with cloud by >1 K day^{-1} . Sometimes the model underestimates HR maxima (e.g., ~ 3 km in Fig. 8a, ~ 6 km in Fig. 8c), and sometimes it overestimates (e.g., ~ 7 km in Fig. 8a, ~ 3 km in Fig. 8c, ~ 8 km in Fig. 8e). Fourth, Figs. 8e and 8g are manifestations of the model’s positive near-surface HR bias in east Antarctica (Fig. 7d).

Figure 9 shows the model’s performance as a function of SZA and AOD. Supplemental Fig. S26 is analogous but shows relative, instead of raw, errors. We make three observations

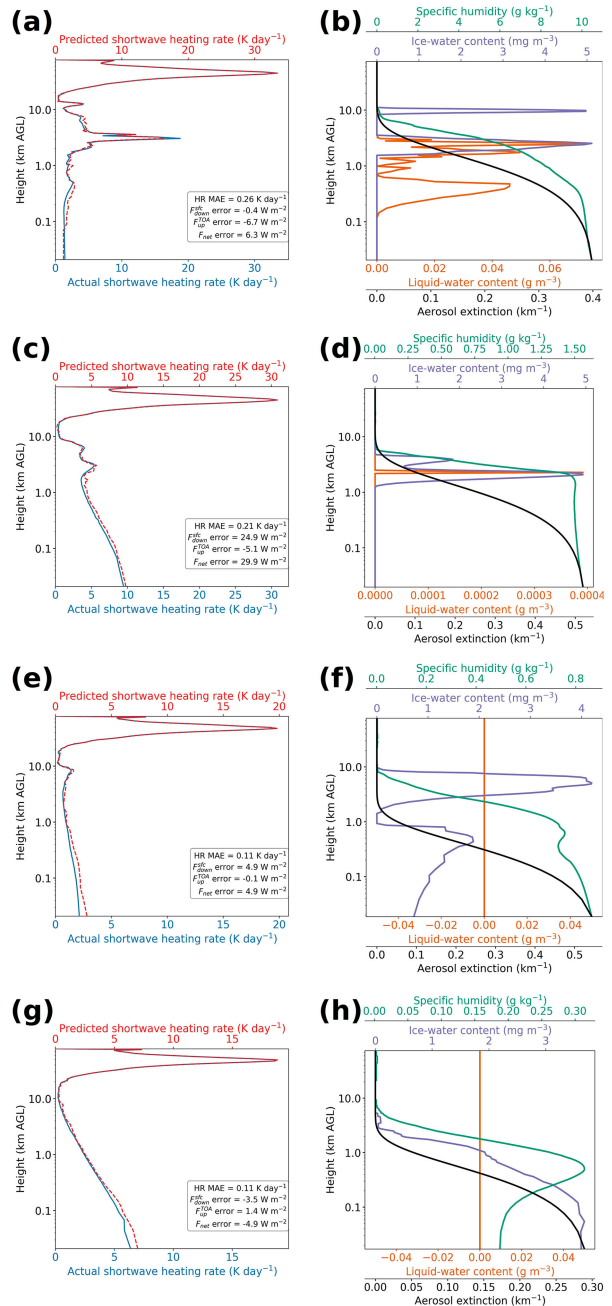


FIG. 8. Geography-based case studies for the best shortwave model. (a),(b) Case study from Tibet, with AOD of 0.61 and SZA of 16.6° ; (c),(d) another case study from Tibet, with AOD of 0.72 and SZA of 11.2° ; (e),(f) case study from east Antarctica, with AOD of 0.23 and SZA of 67.7° ; (g),(h) another case study from east Antarctica, with AOD of 0.17 and SZA of 70.7° . For each case study, (left) actual and predicted RT solutions and (right) four of the most important predictor variables for shortwave RT. In each left panel, the legend shows column-averaged MAE for HR (labeled “HR MAE”) and errors for the three flux variables (predicted minus actual). AOD is a summary of an important predictor variable (the height-integrated aerosol extinction), while SZA is an important predictor variable itself. These scalars are thus reported in the caption above.

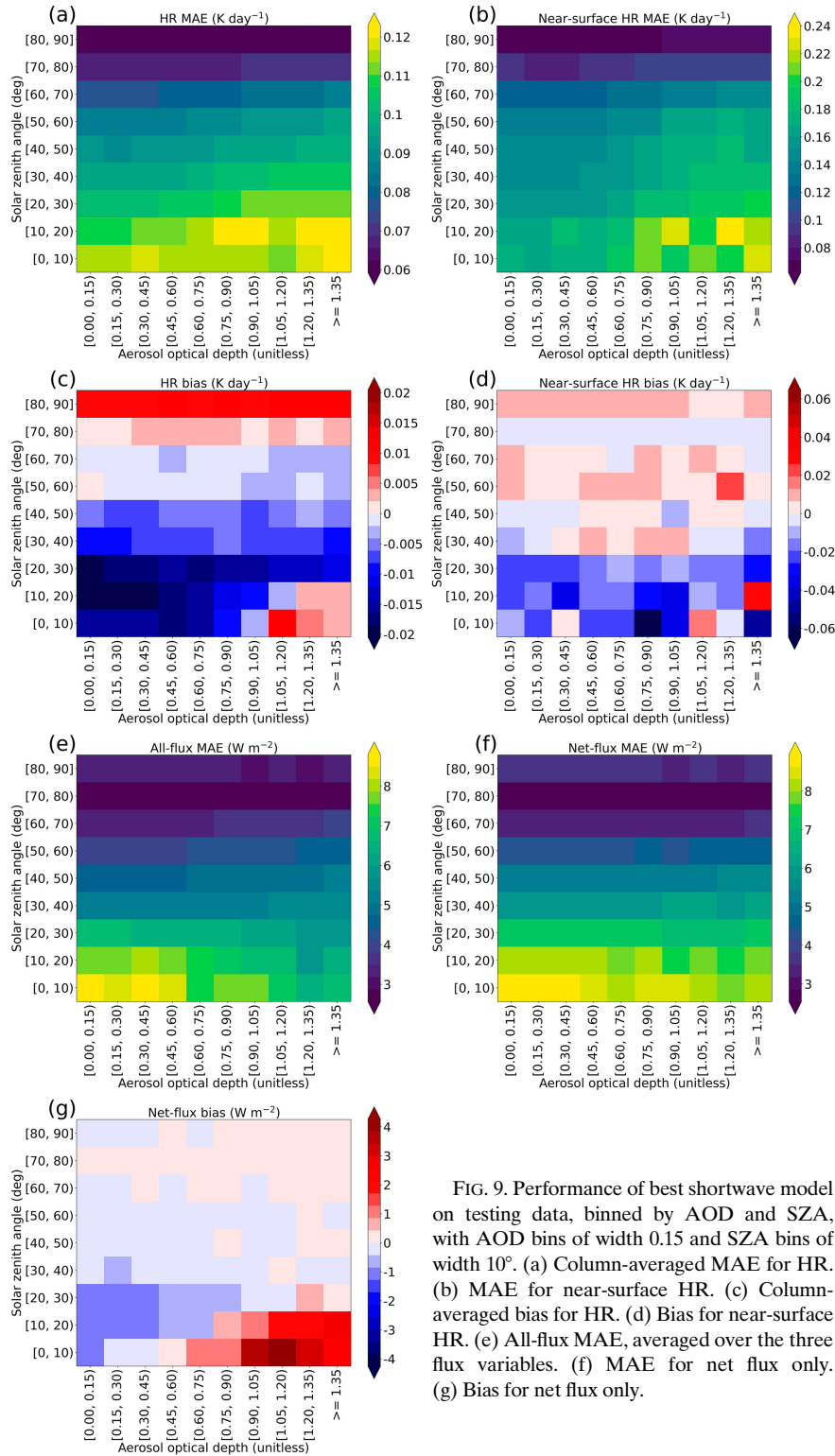


FIG. 9. Performance of best shortwave model on testing data, binned by AOD and SZA, with AOD bins of width 0.15 and SZA bins of width 10°. (a) Column-averaged MAE for HR. (b) MAE for near-surface HR. (c) Column-averaged bias for HR. (d) Bias for near-surface HR. (e) All-flux MAE, averaged over the three flux variables. (f) MAE for net flux only. (g) Bias for net flux only.

from the two figures. First, for all error metrics except net-flux bias (Figs. 9a–f and Figs. S26a–f), raw error decreases strongly with SZA and increases weakly with AOD. In other words, raw errors are worst when there is a lot of incoming solar

radiation and a lot of interaction with aerosols. Second, for the same error metrics, relative error increases strongly with SZA (the opposite relationship to raw error) and has no apparent relationship with AOD. Thus, higher solar radiation and

aerosol content do not make shortwave RT *fundamentally* harder to predict; raw errors increase because the actual values (HRs and fluxes) increase. Third, for net-flux bias (Fig. 9g and Fig. S26g), when $SZA < 20^\circ$, both raw and relative error increase with decreasing SZA and increasing AOD. In other words, when $SZA < 20^\circ$, higher solar radiation and aerosol content make it fundamentally harder to predict net flux without bias. Supplemental Fig. S27—with errors for individual flux variables rather than all-flux errors—shows that this last relationship is driven primarily by biases in $F_{\text{down}}^{\text{sc}}$, which are larger than biases in $F_{\text{up}}^{\text{TOA}}$.

Figure 10 shows case studies from the low-SZA/high-AOD regime (defined as $SZA \leq 20^\circ$ and $AOD \geq 0.75$), where raw errors are highest. The following observations aim to represent 200 random profiles, a superset of the four shown in Fig. 10. First, many low-SZA/high-AOD cases feature ice cloud near the tropopause, including the first three in Fig. 10. This is a known climatological feature of the tropics (Jensen et al. 2013), where the vast majority of low-SZA/high-AOD cases occur. Second, low-SZA/high-AOD cases without liquid cloud (Figs. 10e–h) feature large HRs in the bottom ~ 1 km of the atmosphere, where the model sometimes overestimates (Fig. 10e) but generally underestimates (Fig. 10g)—consistent with the bottom grid row in Fig. 9d. Third, the model generally overestimates net flux for these cases (by a large amount in Fig. 10e). This is due mainly to overestimating $F_{\text{down}}^{\text{sc}}$ in the low-SZA/high-AOD regime (supplemental Fig. S27).

d. Best longwave model

Figure 11 shows the overall performance of the best longwave model. For all flux variables (Figs. 11a–c), the model is almost perfectly reliable and almost perfectly reproduces the observed distribution. The model has only one perceptible conditional bias, namely, an underprediction of $\sim 10 \text{ W m}^{-2}$ for the lowest $F_{\text{up}}^{\text{TOA}}$ predictions. In other words, when the model predicts an extremely low $F_{\text{up}}^{\text{TOA}}$, the prediction is slightly too extreme. The model has an absolute bias $\ll 0.1 \text{ K day}^{-1}$ for HR at every height (Fig. 11d) but much larger MAEs (Fig. 11e), reaching 0.55 and 0.24 K day^{-1} at the bottom two grid levels (~ 21 and $\sim 44 \text{ m AGL}$). As will be shown, longwave RT near the surface is sensitive to fine-scale details of the thermodynamic profile, which the model struggles to capture. Because the climatological model also has its largest HR MAE at the surface, the NN model's local maximum in MAE does not translate to a local minimum in MAE skill score (Fig. 11f). Last, the attributes diagram for HR (Fig. 11g) tells a similar story to those for the flux variables: the model is almost perfectly reliable and almost perfectly reproduces the observed distribution. Supplemental Figs. S31 and S32 are analogous to Fig. 11 but only for extreme cases—i.e., the 3% of testing profiles with the greatest height-maximum and height-averaged absolute HR, respectively. As for the shortwave model, we find that although errors are higher for the extreme cases, HRs and fluxes still have almost perfect reliability and absolute HR bias is well below 0.1 K day^{-1} throughout the profile.

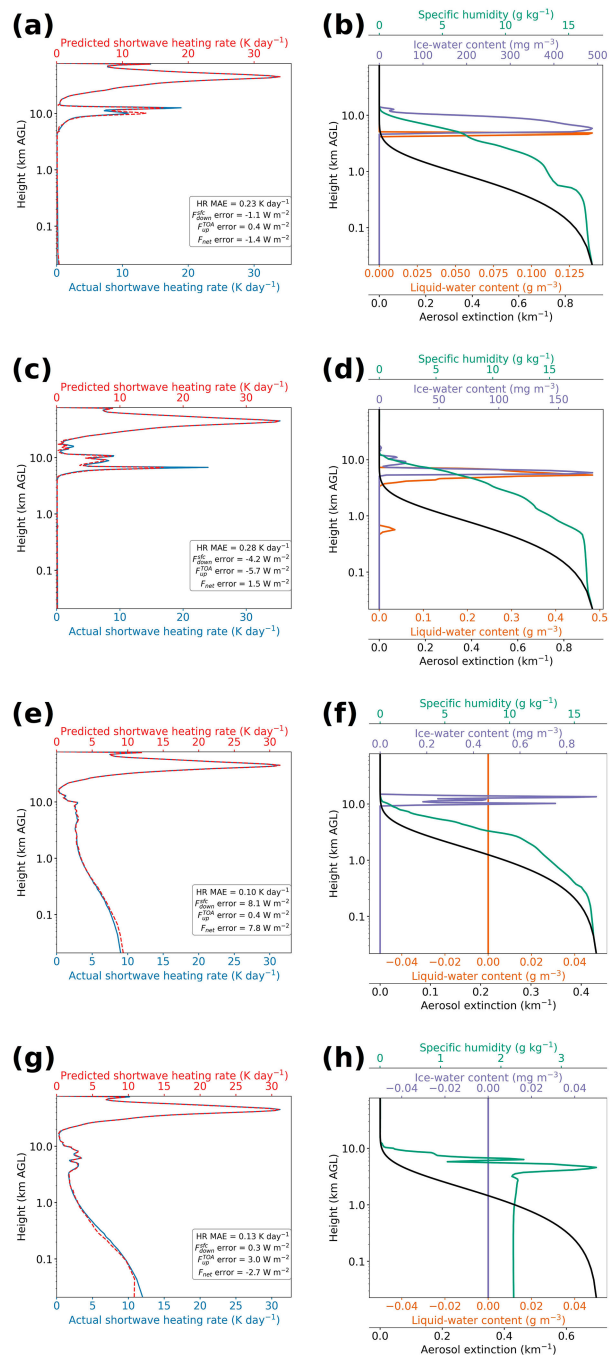


FIG. 10. Regime-based case studies for the best shortwave model, specifically from the low-SZA–high-AOD regime, defined as $SZA \leq 20^\circ$ and $AOD \geq 0.75$. Formatting is explained in the caption of Fig. 8. (a),(b) $AOD = 0.85$ and $SZA = 10.7^\circ$; (c),(d) $AOD = 0.81$ and $SZA = 7.6^\circ$; (e),(f), $AOD = 0.76$ and $SZA = 7.9^\circ$; (g),(h), $AOD = 1.44$ and $SZA = 5.5^\circ$.

Figure 12 shows the model's performance as a function of liquid-only cloud regime. Performance for other cloud phases (ice only, mixed phase, and any phase) is shown in supplemental Figs. S28–S30. The attributes diagrams for flux variables

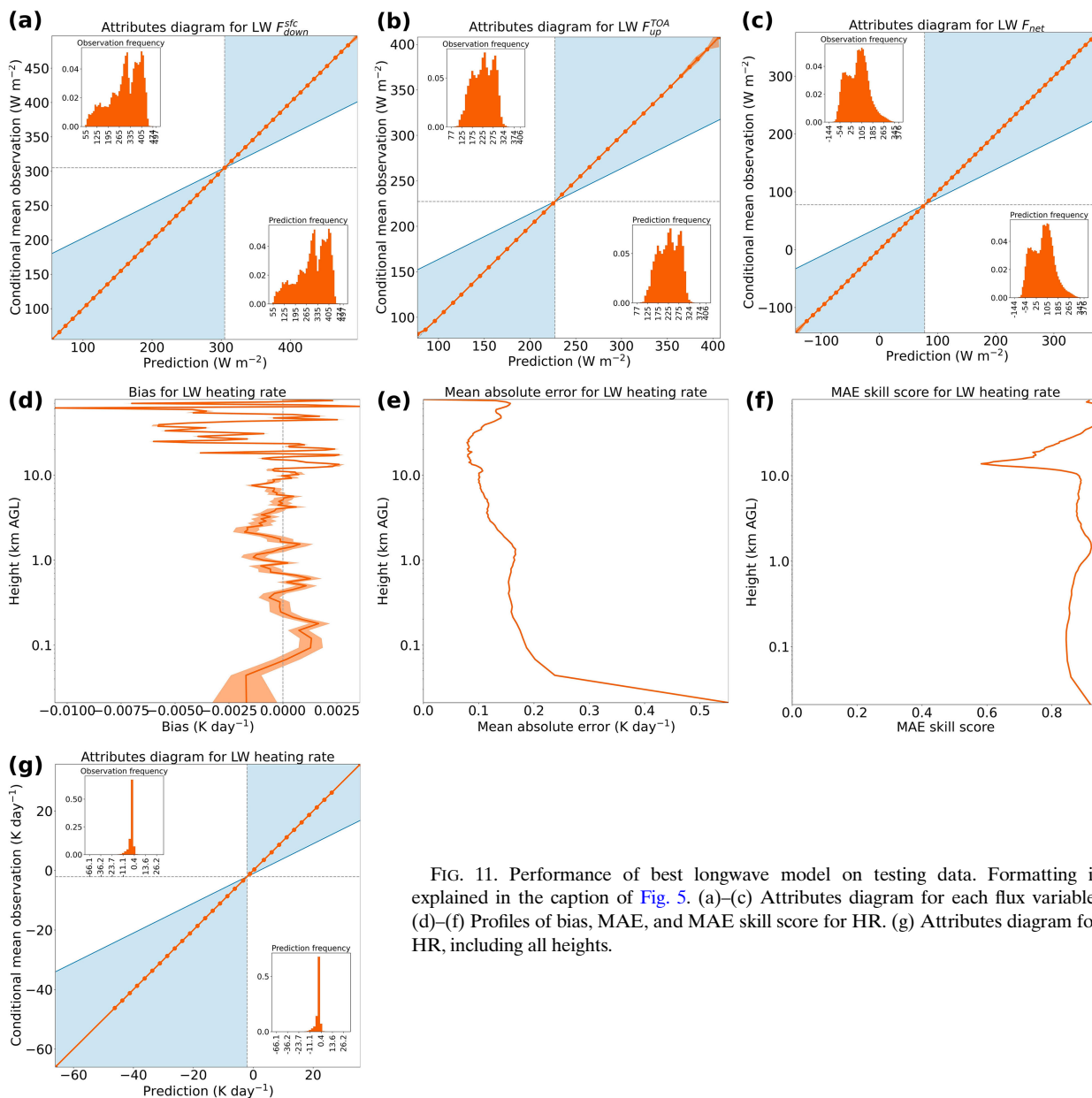


FIG. 11. Performance of best longwave model on testing data. Formatting is explained in the caption of Fig. 5. (a)–(c) Attributes diagram for each flux variable. (d)–(f) Profiles of bias, MAE, and MAE skill score for HR. (g) Attributes diagram for HR, including all heights.

(Figs. 11a–c) tell a similar story to the cloud-agnostic versions (Figs. 11a–c): a few slight conditional biases but no absolute bias exceeding 20 W m^{-2} . In the bottom few 100 m of the troposphere, HR errors (Figs. 12d–f) are best for clear-sky profiles, followed by single- and multilayer cloud, and worst for foggy profiles. In other words, the largest HR errors in the bottom few 100 m are caused by clouds, especially clouds that reach the surface. Meanwhile, in the troposphere above $\sim 1 \text{ km}$, HR errors (Figs. 12d–f) are best for clear-sky profiles and worst for those with multilayer cloud. Errors for foggy profiles above $\sim 1 \text{ km}$ are intermediate, because many surface-based clouds are not thick enough to reach these heights. Last, the attributes diagram for HR (Figs. 12g–j) is nearly perfect in

all cloud regimes except fog. The model has a considerable negative bias (as large as 1 K day^{-1}) when predicting HR above 20 K day^{-1} in foggy profiles, but as shown by the confidence interval—which overlaps the 1:1 line—this apparent defect could be an artifact of small sample size.

Figure 13 shows the model's performance as a function of location. The column-averaged MAE for HR (Fig. 13a) is mostly between 0.10 and 0.15 K day^{-1} ; it exceeds 0.15 K day^{-1} at a few locations, notably Tibet, southern Peru, and the northwestern Rocky Mountains. The MAE for near-surface HR (Fig. 13b) is much larger—mostly between 0.35 and 0.94 K day^{-1} , exceeding 0.94 K day^{-1} at the same locations. The column-averaged bias for HR (Fig. 13c) is mostly between -0.01 and $+0.01 \text{ K day}^{-1}$,

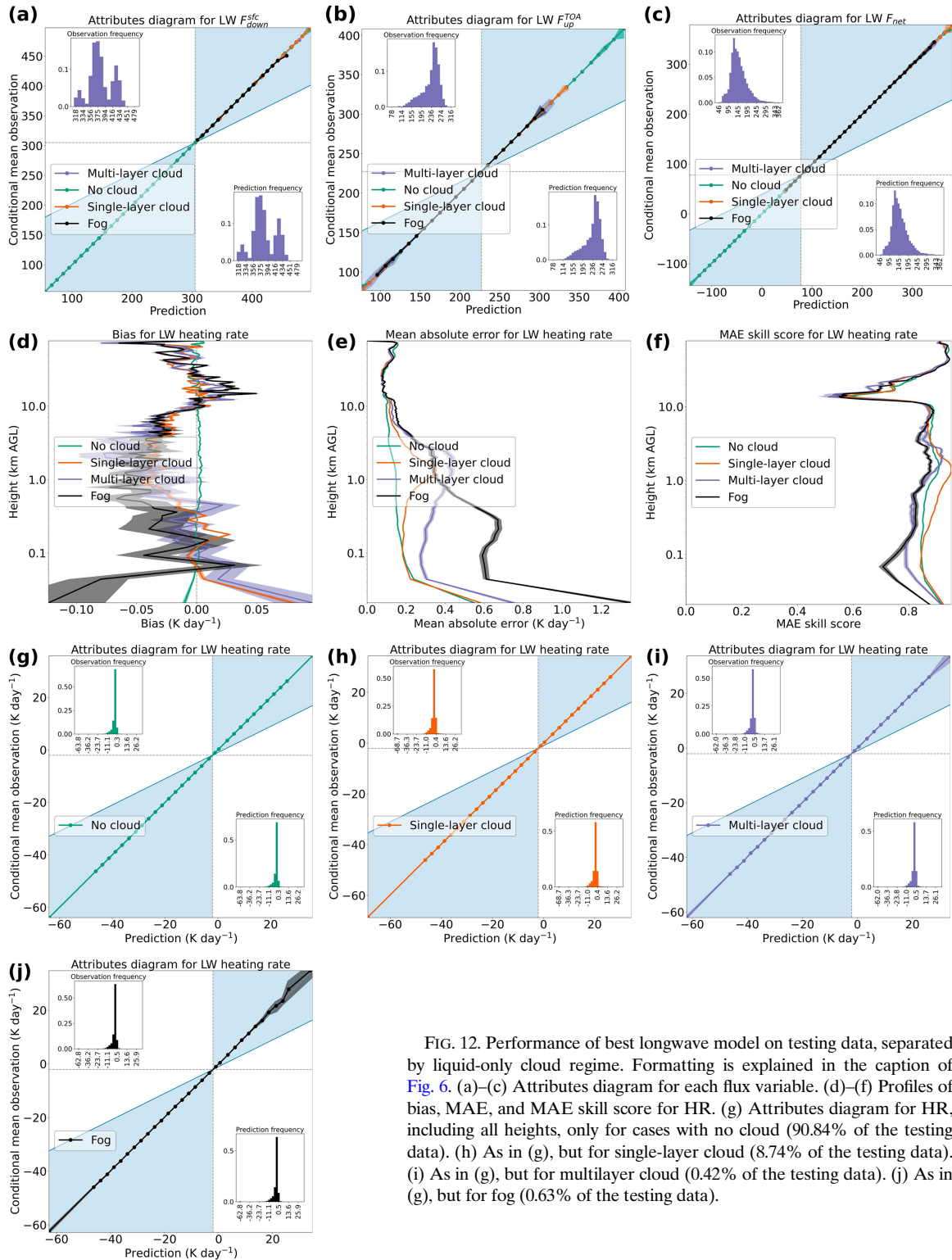


FIG. 12. Performance of best longwave model on testing data, separated by liquid-only cloud regime. Formatting is explained in the caption of Fig. 6. (a)–(c) Attributes diagram for each flux variable. (d)–(f) Profiles of bias, MAE, and MAE skill score for HR. (g) Attributes diagram for HR, including all heights, only for cases with no cloud (90.84% of the testing data). (h) As in (g), but for single-layer cloud (8.74% of the testing data). (i) As in (g), but for multilayer cloud (0.42% of the testing data). (j) As in (g), but for fog (0.63% of the testing data).

with absolute bias not exceeding 0.02 K day^{-1} at any location. The bias for near-surface HR (Fig. 13d) is larger—mostly between -0.24 and $+0.22 \text{ K day}^{-1}$, with absolute value exceeding 0.24 K day^{-1} in Tibet, northern South America, and the

northwestern Rockies. The all-flux MAE (Fig. 13e) is mostly between 0.24 and 0.63 W m^{-2} , exceeding 0.63 W m^{-2} mainly in Tibet. The net-flux MAE (Fig. 13f) follows a similar pattern to the all-flux MAE. The net-flux bias (Fig. 13g) is mostly

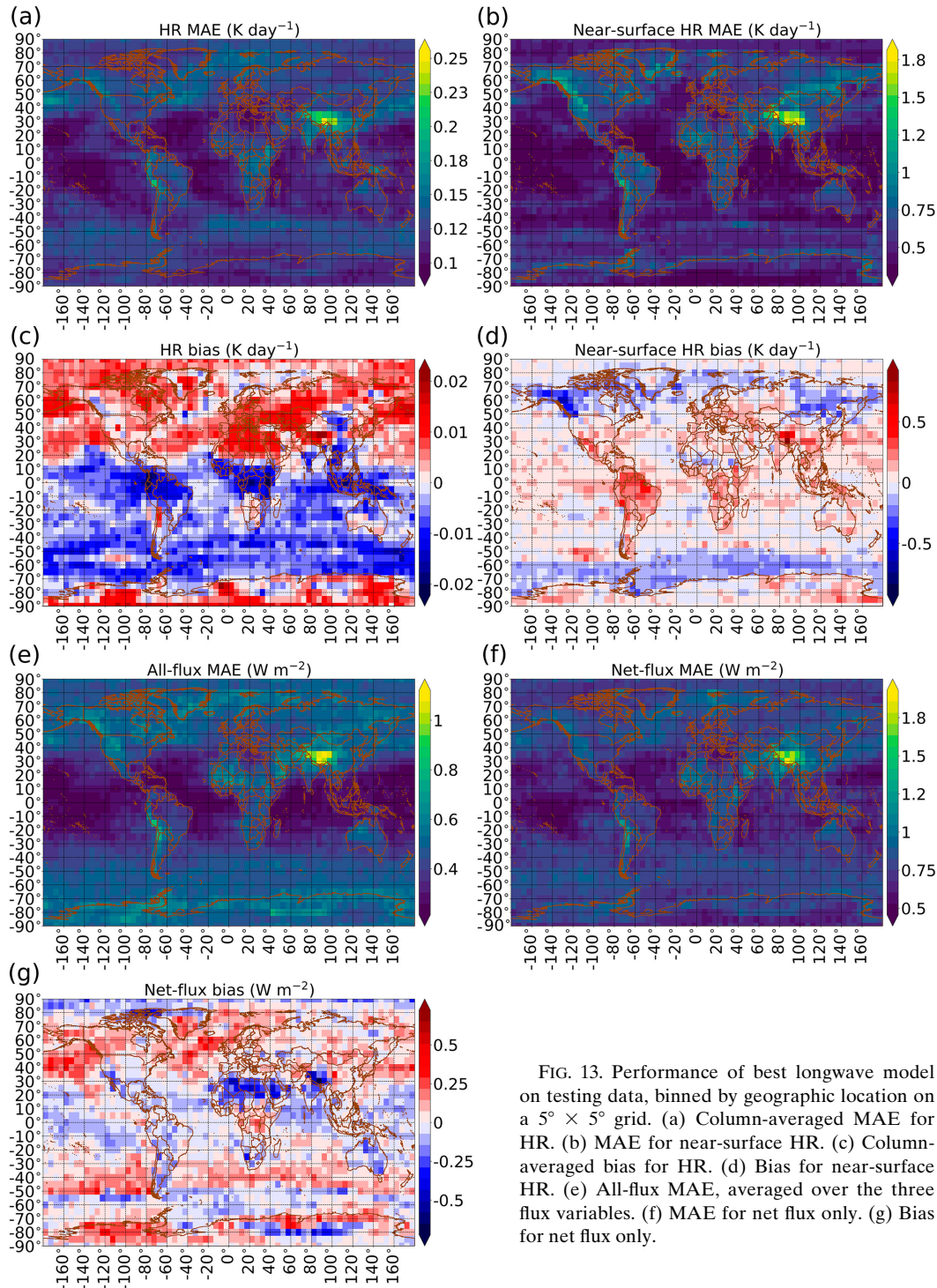


FIG. 13. Performance of best longwave model on testing data, binned by geographic location on a $5^\circ \times 5^\circ$ grid. (a) Column-averaged MAE for HR. (b) MAE for near-surface HR. (c) Column-averaged bias for HR. (d) Bias for near-surface HR. (e) All-flux MAE, averaged over the three flux variables. (f) MAE for net flux only. (g) Bias for net flux only.

between -0.23 and $+0.24 \text{ W m}^{-2}$, with absolute bias not exceeding 0.72 K day^{-1} at any location. Maxima in raw error mostly correspond to maxima in relative error (supplemental Fig. S33), which indicates that longwave RT is fundamentally harder to predict in these regions. Last, supplemental Fig. S34

shows that, while $F_{\text{down}}^{\text{sc}}$ and $F_{\text{up}}^{\text{TOA}}$ have similar MAE values over most of the globe, $F_{\text{down}}^{\text{sc}}$ bias is worse than $F_{\text{up}}^{\text{TOA}}$ bias at most locations. Thus, at most locations, net-flux bias (which equals $F_{\text{down}}^{\text{sc}}$ bias minus $F_{\text{up}}^{\text{TOA}}$ bias) primarily reflects $F_{\text{down}}^{\text{sc}}$ bias, with a small contribution from $F_{\text{up}}^{\text{TOA}}$.

Figure 14 shows case studies from regions with high model error: Tibet (Figs. 14a–d), the northwestern Rockies (Figs. 14e,f), and southern Peru (Figs. 14g,h). The following observations aim to represent 800 random profiles (200 per region), a superset of the four shown in Fig. 14. First, most of the 800 profiles feature liquid and/or ice cloud. Like the shortwave model, the longwave model matches the shape of the HR profile well but often misses extreme HRs associated with cloud by $>1 \text{ K day}^{-1}$. Sometimes the model overestimates longwave cooling above clouds (e.g., ~ 2.5 and $\sim 10 \text{ km}$ in Fig. 14a, $\sim 8 \text{ km}$ in Fig. 14c), and sometimes it underestimates cooling (e.g., ~ 0.4 and $\sim 4 \text{ km}$ in Fig. 14g). Second, as for shortwave RT, regions with high longwave error have very high surface elevations, which are globally rare. Third, sometimes longwave HR error near the surface is large even for profiles that appear uncomplicated near the surface (e.g., Figs. 14e,f), because near-surface longwave RT is sensitive to fine details of the near-surface thermodynamic profile.

Figure 15 shows the model’s performance as a function of near-surface thermodynamics, specifically, the temperature lapse rate [Γ_T^{sfc} in Eq. (5)] and humidity lapse rate [Γ_q^{sfc} in Eq. (5)]. First, we note that all error metrics (Figs. 15a–g) are worst in two regimes, which we call the positive/positive and negative/negative regimes. The positive/positive regime has large positive Γ_T^{sfc} and Γ_q^{sfc} —i.e., both temperature and humidity decrease strongly with height. The negative/negative regime has large negative lapse rates—i.e., both temperature and humidity exhibit a strong inversion, increasing with height. Second, both the positive/positive and negative/negative regimes are quite rare in our dataset, as shown in Fig. 15h. Most profiles have a small positive Γ_T^{sfc} and small positive Γ_q^{sfc} , the “common” regime labeled in Fig. 15. Third, while all error metrics are worst in the positive/positive and negative/negative regimes, the most egregious errors are for near-surface HR, where both MAE (Fig. 15b) and absolute bias (Fig. 15d) can be $\gg 1 \text{ K day}^{-1}$. Fourth, relative error (supplemental Fig. S35) is also maximized in the positive/positive and negative/negative regimes, which indicates that extreme near-surface thermodynamics make longwave RT fundamentally harder to predict. Last, supplemental Fig. S36 shows that $F_{\text{down}}^{\text{sfc}}$ errors are worse than $F_{\text{up}}^{\text{TOA}}$ errors in both regimes.

Figure 16 shows case studies from the negative/negative regime (Figs. 16a–d) and positive/positive regime (Figs. 16e–h). The following observations aim to represent 400 random profiles (200 per regime), a superset of the four shown in Fig. 16. First, we note that most of these profiles feature extreme near-surface heating or cooling. Second, like the geography-based case studies (Fig. 14), the model generally performs well for these regime-based case studies, except for near-surface HR and a few extremes associated with cloud (e.g., $\sim 1.5 \text{ km}$ in Fig. 16e). Third, the model’s fractional error for near-surface HR is generally quite low; cases like Fig. 16a do not occur very often.

6. Summary and future work

We have developed neural networks (NN) to emulate the full RRTM, i.e., the shortwave and longwave RRTM with all

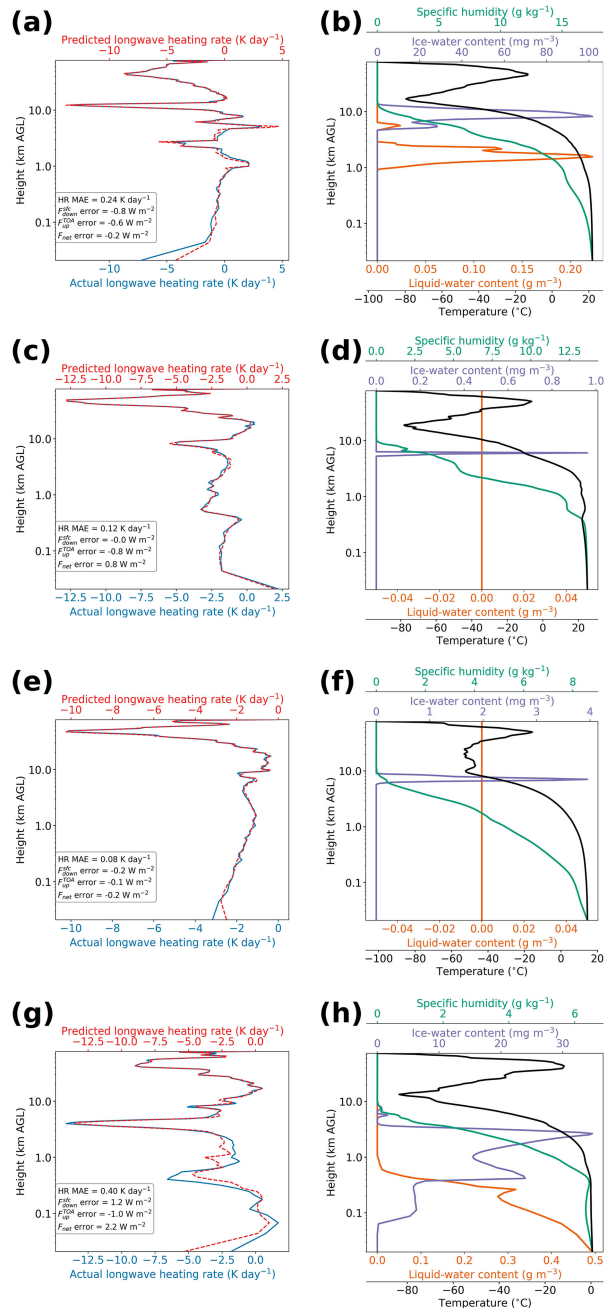


FIG. 14. Geography-based case studies for the best longwave model. (a),(b) Case study from Tibet, with $\Gamma_T^{\text{sfc}} = 4.40 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 5.9 \text{ g kg}^{-1} \text{ km}^{-1}$; (c),(d) another case study from Tibet, with $\Gamma_T^{\text{sfc}} = 11.75 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 0.7 \text{ g kg}^{-1} \text{ km}^{-1}$; (e),(f) case study from northwestern Rockies, with $\Gamma_T^{\text{sfc}} = 4.62 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 11.9 \text{ g kg}^{-1} \text{ km}^{-1}$; (g),(h) case study from southern Peru, with $\Gamma_T^{\text{sfc}} = 10.91 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 4.3 \text{ g kg}^{-1} \text{ km}^{-1}$. For each case study, (left) actual and predicted RT solutions and (right) four of the most important predictor variables for longwave RT. In each left panel, the legend shows column-averaged MAE for HR (labeled “HR MAE”) and errors for the three flux variables (predicted minus actual). Γ_T^{sfc} and Γ_q^{sfc} [Eq. (5)] are summaries of important predictor variables (the thermodynamic profiles). These scalars are thus reported in the caption above.

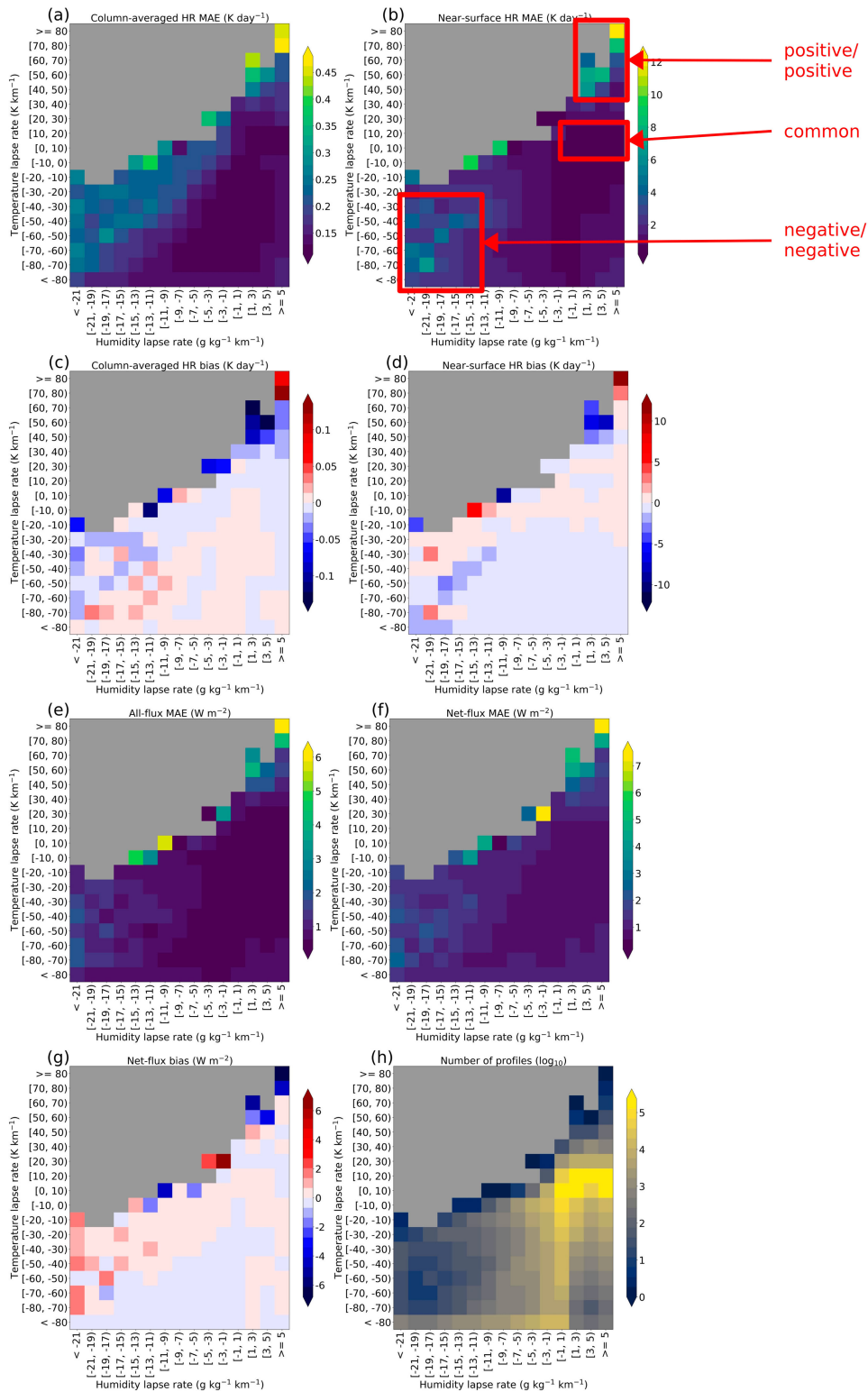


FIG. 15. Performance of best longwave model on testing data, binned by near-surface thermodynamic lapse rates, with Γ_T^{stc} bins of width 10 K km^{-1} and Γ_q^{stc} bins of width $2 \text{ g kg}^{-1} \text{ km}^{-1}$. The three labeled regimes (positive/positive, negative/negative, and common) are explained in the main text. (a) Column-averaged MAE for HR. (b) MAE for near-surface HR. (c) Column-averaged bias for HR. (d) Bias for near-surface HR. (e) All-flux MAE, averaged over the three flux variables. (f) MAE for net flux only. (g) Bias for net flux only. (h) Number of testing samples per bin, in logarithmic scale.

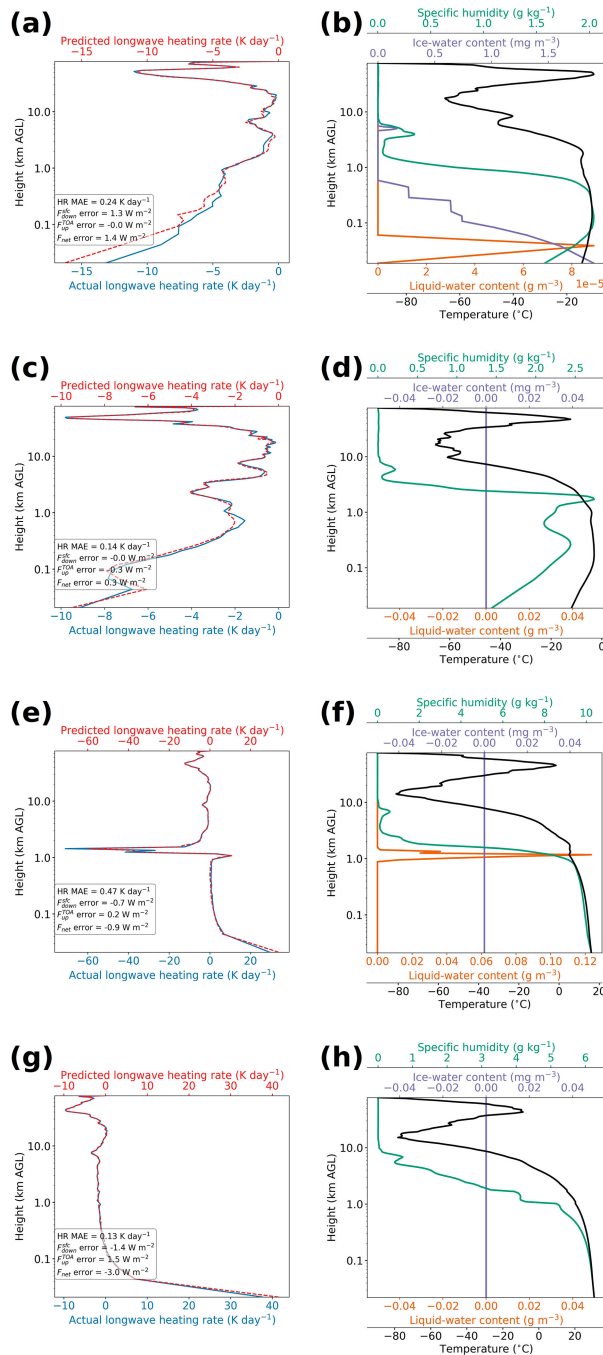


FIG. 16. Regime-based case studies for the best longwave model. (a),(b) Case study from the negative/negative regime, defined as $\Gamma_T^{\text{sfc}} < -30 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} < -13 \text{ g kg}^{-1} \text{ km}^{-1}$. Exact values here are $\Gamma_T^{\text{sfc}} = -94.56 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = -14.0 \text{ g kg}^{-1} \text{ km}^{-1}$. (c),(d) Another case study from the negative/negative regime, with $\Gamma_T^{\text{sfc}} = -164.83 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = -16.2 \text{ g kg}^{-1} \text{ km}^{-1}$. (e),(f) Case study from the positive/positive regime, defined as $\Gamma_T^{\text{sfc}} > 40 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} > 1 \text{ g kg}^{-1} \text{ km}^{-1}$. Exact values here are $\Gamma_T^{\text{sfc}} = 44.06 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 8.1 \text{ g kg}^{-1} \text{ km}^{-1}$. (g),(h) Another case study from the positive/positive regime, with $\Gamma_T^{\text{sfc}} = 40.25 \text{ K km}^{-1}$ and $\Gamma_q^{\text{sfc}} = 2.2 \text{ g kg}^{-1} \text{ km}^{-1}$. Formatting is explained in the caption of Fig. 14.

predictor variables. Both the RRTM and NN-based emulators are driven by forecast profiles from the GFSv16 on the native vertical grid, which uses hybrid pressure–sigma coordinates. We experimented with novel deep learning methods designed to produce realistic and accurate spatial structure in gridded predictions: the U-net++ architecture, U-net3+ architecture, and deep-supervision training method. We hypothesized that the best NNs would be those with the U-net3+ architecture and deep supervision. Contrary to our hypotheses, we found that deep supervision leads to worse performance and architecture has little impact. We also experimented with three other hyperparameters—NN width, depth, and spectral complexity—which strongly control the NN’s overall complexity, causing the number of trainable weights to vary from $\mathcal{O}(10^5)$ to $\mathcal{O}(10^{8.5})$. We found that the best NNs are at the more complex end of the spectrum; the selected shortwave and longwave NNs have $10^{7.52}$ and $10^{7.28}$ trainable weights, respectively. Overall, the better NNs are deep (have encoders and decoders at many spatial resolutions), narrow (have only one convolutional layer per block), and have large spectral complexity (many convolutional filters and thus many feature maps). While NN type (U-net++ or U-net3+) has only a weak effect on performance, the best shortwave NN is a U-net++ model, while the best longwave NN is a U-net3+ model. Our NNs are an example of knowledge-guided machine learning, identified as a major need in ML applications to the geosciences (Gil et al. 2019; Reichstein et al. 2019). Specifically, we enforce energy conservation in the NNs [Eq. (2)]; use a custom loss function to emphasize large heating rates (HR), which are rare but important for weather and climate [Eq. (3)]; and include custom predictors to account for vertically nonlocal effects [section 3c(3) of L21].

The best shortwave NN model performs extremely well in an aggregate sense, i.e., averaged over all the testing data. Highlights include reliable fluxes, with all conditional biases $< 10 \text{ W m}^{-2}$ in absolute value; reliable HRs, with all conditional biases $< 1 \text{ K day}^{-1}$ in absolute value; and absolute HR bias $< 0.1 \text{ K day}^{-1}$ at all heights, suggesting that the NN could be stably integrated into the GFSv16 as a parameterization. The model also performs extremely well in all cloud regimes, at most geographic locations, and in most regimes defined by solar zenith angle (SZA) and aerosol optical depth (AOD). The largest errors occur in Tibet and east Antarctica, which feature high surface elevation/albedo, and in the low-SZA/high-AOD regime, which features a lot of incoming solar radiation and interaction with aerosols. However, even these largest errors are quite small: mean absolute error (MAE) for HR does not exceed 0.6 K day^{-1} , even near the surface; absolute HR bias does not exceed 0.3 K day^{-1} , even near the surface; MAE for flux variables does not exceed 10 W m^{-2} ; and net-flux bias does not exceed 5 W m^{-2} . For regimes that make RT fundamentally harder to predict—e.g., high elevation/albedo, which increase both raw and relative errors—results could potentially be improved by adding training data from these regimes. Table 8 compares our model to NN-based emulators of shortwave RT from three other studies: Krasnopolsky et al. (2012, hereafter K12), SR21, and KS22. Although our model appears to perform best, this

TABLE 8. Comparison of NN-based emulators for shortwave RT. For our model, we use the testing data only. For the comparison studies, we take results from Table 2 of K12, page 7 of SR21 for HR errors, Table 3 (the “WRF15” column) of SR21 for flux errors, and Fig. 1 of KS22 (these values are estimated visually). “Profile RMSE” is defined in Eq. (A1) of K12; “near-surface” means for the lowest model level; and “N/A” means that the statistic is not reported. Although KS22 reports flux errors, the statistic is all-flux RMSE, computed by averaging over three variables: $F_{\text{down}}^{\text{sc}}$, $F_{\text{up}}^{\text{TOA}}$, and $F_{\text{up}}^{\text{sc}}$. We predict a different set of flux variables— F_{net} instead of $F_{\text{up}}^{\text{sc}}$ —and thus do not compare our flux errors with KS22.

Model statistic	Ours	K12	SR21	KS22
Column-averaged HR RMSE (K day ⁻¹)	0.14	0.26	0.17	~0.2
Column-averaged HR bias (K day ⁻¹)	-0.002	-0.007	N/A	N/A
HR profile RMSE (K day ⁻¹)	0.12	0.18	N/A	N/A
Near-surface HR RMSE (K day ⁻¹)	0.20	0.20	N/A	N/A
Near-surface HR bias (K day ⁻¹)	+0.0001	-0.03	N/A	N/A
$F_{\text{down}}^{\text{sc}}$ RMSE (W m ⁻²)	5.85	N/A	43.75	N/A
$F_{\text{up}}^{\text{TOA}}$ RMSE (W m ⁻²)	3.94	N/A	36.20	N/A

comparison is not apples-to-apples, due to different vertical resolutions (127 levels here, 64 in K12, 39 in the other two studies), testing cases (time period and spatial domain), and predictor variables. The three comparison studies omit aerosols, all trace gases other than O₃, LWC and IWC (they use cloud fraction instead, with no distinction between liquid and ice), and the particle size distribution (for which we use liquid and ice effective radii). Last, our shortwave NN runs 7510 times faster than the shortwave RRTM.

The best longwave NN model also performs extremely well in an aggregate sense; highlights include near-perfect reliability for both fluxes and HRs and absolute HR bias \ll 0.1 K day⁻¹ at every height. The model’s main deficiency is a large error in near-surface HR, e.g., an MAE of 0.55 K day⁻¹ at the lowest grid level. However, longwave RT near the surface is complicated, and errors here are often quite large. For example, in Veerman et al. (2020), who emulated only the gas-optics part of the RRTMGp, near-surface HR bias is on the order of 1 K day⁻¹ (their Fig. 2c). The model performs well in all cloud regimes, at most geographic locations, and in most regimes defined by near-surface thermodynamics. The largest errors occur with liquid-only fog, where the bias and MAE for near-surface HR reach -0.12 and 1.3 K day⁻¹, respectively; in Tibet, where near-surface bias and MAE reach almost 1 and 2 K day⁻¹, respectively; and under extreme near-surface thermodynamics, where near-surface absolute bias and MAE are \gg 1 K day⁻¹. However, the extreme thermodynamic regimes are quite rare, so this last number is affected by small sample size. Also, even in the aforementioned regimes with large error in near-surface HR, column-averaged bias for HR does not exceed 0.15 K day⁻¹ in absolute value, column-averaged MAE for HR does not exceed 0.6 K day⁻¹, MAE for flux variables does not exceed 10 W m⁻², and net-flux bias does not exceed 7 W m⁻². Table 9 shows that our

TABLE 9. Comparison of NN-based emulators for longwave RT. For technical notes, see the caption of Table 8.

Model statistic	Ours	K12	SR21	KS22
Column-averaged HR RMSE (K day ⁻¹)	0.22	0.52	0.46	~0.375
Column-averaged HR bias (K day ⁻¹)	-0.0006	+0.008	N/A	N/A
HR profile RMSE (K day ⁻¹)	0.20	0.38	N/A	N/A
Near-surface HR RMSE (K day ⁻¹)	0.83	0.55	N/A	N/A
Near-surface HR bias (K day ⁻¹)	-0.002	+0.02	N/A	N/A
$F_{\text{down}}^{\text{sc}}$ RMSE (W m ⁻²)	0.64	N/A	5.71	N/A
$F_{\text{up}}^{\text{TOA}}$ RMSE (W m ⁻²)	0.81	N/A	7.11	N/A

longwave NN compares very favorably to other studies. Last, our longwave NN runs 90 times faster than the longwave RRTM.

Future work will include three items. First, we will develop grid-agnostic NNs that work on profiles with any vertical resolution. This work may benefit from Fourier neural operators (FNO; Lu et al. 2019; Li et al. 2020), which naturally learn physics in a grid-agnostic manner. Second, we will implement the NNs in online mode, i.e., as a parameterization in the GFSv16. To this end we have converted the NNs to a Fortran-friendly format, using the Infero library (ECMWF 2022), and ensured that the NNs yield the same predictions in Fortran as in Python. Note that the NNs alone cannot handle sub-grid-scale fractional cloudiness, as cloud fraction is a predictor in neither the RRTM nor the NNs. To handle fractional cloudiness in online mode, we will couple the NNs with the Monte Carlo independent-column approximation (Pincus et al. 2003). Third, we will perform thorough testing of the NNs in online mode. Specifically, we will conduct month-long retrospective simulations in both the summer and winter, using a control model (original parameterization) and experimental model (NN parameterization). We will compare the two models against each other and against observations, using methods as in Turner et al. (2012, 2020). Given the accuracy and efficiency of modern deep NNs, we expect them to replace many existing parameterizations in weather and climate models. However, operational use should proceed only after thorough NN evaluation and with the caution that NNs may generalize poorly outside the distribution of their training data, e.g., to future climates.¹⁰ Safeguards against this problem should be built into NN parameterizations, such as continued online learning or out-of-distribution detection.

Acknowledgments. This work was partially supported by the NOAA Global Systems Laboratory, Cooperative Institute for Research in the Atmosphere, and NOAA Award

¹⁰ In our case, motivated by the strong influence of clouds on radiation—including their phase and number of layers—we paid particular attention to the NNs’ ability to emulate the RRTM for all cloud types.

NA19OAR4320073. Author Ebert-Uphoff's work was partially supported by NSF AI Institute Grant 2019758 and NSF Grant 1934668.

Data availability statement. The input data (predictor and target variables for all the time periods: NN training, IR training, validation, and testing) and selected models (best shortwave NN, best longwave NN, and IR model used to bias correct each one) are stored on NOAA's high-performance computing systems and are available from the authors upon request. We used version 2.0.0 of Machine Learning for Radiative Transfer (ML4RT; <https://doi.org/10.5281/zenodo.7378773>)—a Python library managed by author Lagerquist—for all training, evaluation, and analysis.

REFERENCES

- Anderson, G. P., S. A. Clough, F. X. Kneizys, J. H. Chetwynd, and E. P. Shettle, 1986: AFGL atmospheric constituent profiles (0–120 km). Air Force Geophysics Laboratory Tech. Rep. AFGL-TR-86-0110, 46 pp., <https://apps.dtic.mil/sti/citations/ADA175173>.
- Baek, S., 2017: A revised radiation package of G-packed McICA and two-stream approximation: Performance evaluation in a global weather forecasting model. *J. Adv. Model. Earth Syst.*, **9**, 1628–1640, <https://doi.org/10.1002/2017MS000994>.
- Belochitski, A., and V. Krasnopolsky, 2021: Robustness of neural network emulations of radiative transfer parameterizations in a state-of-the-art general circulation model. *Geosci. Model Dev.*, **14**, 7425–7437, <https://doi.org/10.5194/gmd-14-7425-2021>.
- Beucler, T., and Coauthors, 2021: Climate-invariant machine learning. arXiv, 2112.08440v2, <https://doi.org/10.48550/arXiv.2112.08440>.
- Boukabara, S.-A., V. Krasnopolsky, J. Q. Stewart, E. S. Maddy, N. Shahrudi, and R. N. Hoffman, 2019: Leveraging modern artificial intelligence for remote sensing and NWP: Benefits and challenges. *Bull. Amer. Meteor. Soc.*, **100**, ES473–ES491, <https://doi.org/10.1175/BAMS-D-18-0324.1>.
- Chantry, M., P. Dueben, R. Hogan, and P. Ukkonen, 2022: Progress on emulating the radiation scheme via machine learning. *ECMWF Newsletter*, No. 173, ECMWF, Reading, United Kingdom, 9–10, <https://www.ecmwf.int/en/newsletter/173/news/progress-emulating-radiation-scheme-machine-learning>.
- , P. Ukkonen, R. Hogan, and P. Dueben, 2023: Emulating radiative transfer in a numerical weather prediction model. *2023 EGU General Assembly*, Vienna, Austria, European Geoscience Union, Abstract EGU23-3256, <https://doi.org/10.5194/egusphere-egu23-3256>.
- Chevallier, F., F. Chéruy, N. A. Scott, and A. Chédin, 1998: A neural network approach for a fast and accurate computation of a longwave radiative budget. *J. Appl. Meteor.*, **37**, 1385–1397, [https://doi.org/10.1175/1520-0450\(1998\)037<1385:ANNAFA>2.0.CO;2](https://doi.org/10.1175/1520-0450(1998)037<1385:ANNAFA>2.0.CO;2).
- Cotronei, A., and T. Slawig, 2020: Single-precision arithmetic in ECHAM radiation reduces runtime and energy consumption. *Geosci. Model Dev.*, **13**, 2783–2804, <https://doi.org/10.5194/gmd-13-2783-2020>.
- ECMWF, 2022: Inero: A lower-level API for machine learning inference in operations. GitHub, <https://github.com/ecmwf-projects/infero>.
- Geiss, A., P.-L. Ma, B. Singh, and J. C. Hardin, 2022: Emulating aerosol optics with randomly generated neural networks. *Geosci. Model Dev.*, **16**, 2355–2370, <https://doi.org/10.5194/gmd-16-2355-2023>.
- Gil, Y., and Coauthors, 2019: Intelligent systems for geosciences: An essential research agenda. *Commun. ACM*, **62**, 76–84, <https://doi.org/10.1145/3192335>.
- Goodfellow, I., Y. Bengio, and A. Courville, 2016: *Deep Learning*. MIT Press, 800 pp.
- Hsu, W.-R., and A. H. Murphy, 1986: The attributes diagram: A geometrical framework for assessing the quality of probability forecasts. *Int. J. Forecasting*, **2**, 285–293, [https://doi.org/10.1016/0169-2070\(86\)90048-8](https://doi.org/10.1016/0169-2070(86)90048-8).
- Huang, H., and Coauthors, 2020: UNet 3+: A full-scale connected UNet for medical image segmentation. arXiv, 2004.08790v1, <https://doi.org/10.48550/arXiv.2004.08790>.
- Iacono, M. J., J. S. Delamere, E. J. Mlawer, M. W. Shephard, S. A. Clough, and W. D. Collins, 2008: Radiative forcing by long-lived greenhouse gases: Calculations with the AER radiative transfer models. *J. Geophys. Res.*, **113**, D13103, <https://doi.org/10.1029/2008JD009944>.
- Jensen, E. J., and Coauthors, 2013: Ice nucleation and dehydration in the tropical tropopause layer. *Proc. Natl. Acad. Sci. USA*, **110**, 2041–2046, <https://doi.org/10.1073/pnas.1217104110>.
- Kim, P. S., and H.-J. Song, 2022: Usefulness of automatic hyperparameter optimization in developing radiation emulator in a numerical weather prediction model. *Atmosphere*, **13**, 721, <https://doi.org/10.3390/atmos13050721>.
- Krasnopolsky, V., A. Belochitski, Y. Hou, S. Lord, and F. Yang, 2012: Accurate and fast neural network emulations of long and short wave radiation for the NCEP Global Forecast System model. NCEP Office Note 471, 36 pp., <https://repository.library.noaa.gov/view/noaa/6951>.
- Lagerquist, R., D. Turner, I. Ebert-Uphoff, J. Stewart, and V. Hagerty, 2021: Using deep learning to emulate and accelerate a radiative transfer model. *J. Atmos. Oceanic Technol.*, **38**, 1673–1696, <https://doi.org/10.1175/JTECH-D-21-0007.1>.
- Le, T., C. Liu, B. Yao, V. Natraj, and Y. Yung, 2020: Application of machine learning to hyperspectral radiative transfer simulations. *J. Quant. Spectrosc. Radiat. Transfer*, **246**, 106928, <https://doi.org/10.1016/j.jqsrt.2020.106928>.
- Li, Z., N. Kovachki, K. Azizzadenesheli, B. Liu, A. Stuart, K. Bhattacharya, and A. Anandkumar, 2020: Multipole graph neural operator for parametric partial differential equations. *Proc. 34th Int. Conf. on Neural Information Processing Systems*, Online, ACM, 6755–6766, <https://dl.acm.org/doi/abs/10.5555/3495724.3496291>.
- Liu, Y., R. Caballero, and J. M. Monteiro, 2020: RadNet 1.0: Exploring deep learning architectures for longwave radiative transfer. *Geosci. Model Dev.*, **13**, 4399–4412, <https://doi.org/10.5194/gmd-13-4399-2020>.
- Lu, L., P. Jin, and G. E. Karniadakis, 2019: DeepONet: Learning nonlinear operators for identifying differential equations based on the universal approximation theorem of operators. arXiv, 1910.03193v3, <https://doi.org/10.48550/arXiv.1910.03193>.
- Meyer, D., R. J. Hogan, P. D. Dueben, and S. L. Mason, 2022: Machine learning emulation of 3D cloud radiative effects. *J. Adv. Model. Earth Syst.*, **14**, e2021MS002550, <https://doi.org/10.1029/2021MS002550>.
- Miles, N. L., J. Verlinde, and E. E. Clothiaux, 2000: Cloud droplet size distributions in low-level stratiform clouds. *J. Atmos. Sci.*, **57**, 295–311, [https://doi.org/10.1175/1520-0469\(2000\)057<0295:CDSDIL>2.0.CO;2](https://doi.org/10.1175/1520-0469(2000)057<0295:CDSDIL>2.0.CO;2).

- Mishra, S., D. L. Mitchell, D. D. Turner, and R. P. Lawson, 2014: Parameterization of ice fall speeds in midlatitude cirrus: Results from SPARTICUS. *J. Geophys. Res. Atmos.*, **119**, 3857–3876, <https://doi.org/10.1002/2013JD020602>.
- Mitchell, D. L., R. P. Lawson, and B. Baker, 2011: Understanding effective diameter and its application to terrestrial radiation in ice clouds. *Atmos. Chem. Phys.*, **11**, 3417–3429, <https://doi.org/10.5194/acp-11-3417-2011>.
- Mlawer, E. J., and D. D. Turner, 2016: Spectral radiation measurements and analysis in the ARM program. *The Atmospheric Radiation Measurement (ARM) Program: The First 20 Years, Meteor. Monogr.*, Vol. 57, Amer. Meteor. Soc., <https://doi.org/10.1175/AMSMONOGRAPHS-D-15-0027.1>.
- , S. J. Taubman, P. D. Brown, M. J. Iacono, and S. A. Clough, 1997: Radiative transfer for inhomogeneous atmospheres: RRTM, a validated correlated-k model for the longwave. *J. Geophys. Res.*, **102**, 16663–16682, <https://doi.org/10.1029/97JD00237>.
- Muñoz-Esparza, D., C. Becker, J. A. Sauer, D. J. Gagne II, J. Schreck, and B. Kosović, 2022: On the application of an observations-based machine learning parameterization of surface layer fluxes within an atmospheric large-eddy simulation model. *J. Geophys. Res. Atmos.*, **127**, e2021JD036214, <https://doi.org/10.1029/2021JD036214>.
- Myhre, G., C. Myhre, B. Samset, and T. Storelvmo, 2013: Aerosols and their relation to global climate and climate sensitivity. *Nat. Educ. Knowl.*, **4**, 7, <https://www.nature.com/scitable/knowledge/library/aerosols-and-their-relation-to-global-climate-102215345/>.
- Neubauer, D., U. Lohmann, C. Hoese, and M. G. Frontoso, 2014: Impact of the representation of marine stratocumulus clouds on the anthropogenic aerosol effect. *Atmos. Chem. Phys.*, **14**, 11997–12022, <https://doi.org/10.5194/acp-14-11997-2014>.
- Pal, A., S. Mahajan, and M. R. Norman, 2019: Using deep neural networks as cost-effective surrogate models for superparameterized E3SM radiative transfer. *Geophys. Res. Lett.*, **46**, 6069–6079, <https://doi.org/10.1029/2018GL081646>.
- Pincus, R., and B. Stevens, 2013: Paths to accuracy for radiation parameterizations in atmospheric models. *J. Adv. Model. Earth Syst.*, **5**, 225–233, <https://doi.org/10.1002/jame.20027>.
- , H. W. Barker, and J.-J. Morcrette, 2003: A fast, flexible, approximate technique for computing radiative transfer in inhomogeneous cloud fields. *J. Geophys. Res.*, **108**, 4376, <https://doi.org/10.1029/2002JD003322>.
- , E. J. Mlawer, and J. S. Delamere, 2019: Balancing accuracy, efficiency, and flexibility in radiation calculations for dynamical models. *J. Adv. Model. Earth Syst.*, **11**, 3074–3089, <https://doi.org/10.1029/2019MS001621>.
- Rasp, S., 2020: Coupled online learning as a way to tackle instabilities and biases in neural network parameterizations: General algorithms and Lorenz 96 case study (v1.0). *Geosci. Model Dev.*, **13**, 2185–2196, <https://doi.org/10.5194/gmd-13-2185-2020>.
- Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, N. Carvalhais, and Prabhat, 2019: Deep learning and process understanding for data-driven Earth system science. *Nature*, **566**, 195–204, <https://doi.org/10.1038/s41586-019-0912-1>.
- Roh, S., and H.-J. Song, 2020: Evaluation of neural network emulations for radiation parameterization in cloud resolving model. *Geophys. Res. Lett.*, **47**, e2020GL089444, <https://doi.org/10.1029/2020GL089444>.
- Ronneberger, O., P. Fischer, and T. Brox, 2015: U-net: Convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention*, N. Navab et al. Eds., Lecture Notes in Computer Science, Vol. 9351, Springer, 234–241, https://doi.org/10.1007/978-3-319-24574-4_28.
- Schmetz, J., 1989: Towards a surface radiation climatology: Retrieval of downward irradiances from satellites. *Atmos. Res.*, **23**, 287–321, [https://doi.org/10.1016/0169-8095\(89\)90023-9](https://doi.org/10.1016/0169-8095(89)90023-9).
- Song, H.-J., and S. Roh, 2021: Improved weather forecasting using neural network emulation for radiation parameterization. *J. Adv. Model. Earth Syst.*, **13**, e2021MS002609, <https://doi.org/10.1029/2021MS002609>.
- Sonoda, S., and N. Murata, 2017: Neural network with unbounded activation functions is universal approximator. *Appl. Comput. Harmonic Anal.*, **43**, 233–268, <https://doi.org/10.1016/j.acha.2015.12.005>.
- Stegmann, P. G., B. Johnson, I. Moradi, B. Karpowicz, and W. McCarty, 2022: A deep learning approach to fast radiative transfer. *J. Quant. Spectrosc. Radiat. Transfer*, **280**, 108088, <https://doi.org/10.1016/j.jqsrt.2022.108088>.
- Tang, T., and Coauthors, 2020: Response of surface shortwave cloud radiative effect to greenhouse gases and aerosols and its impact on summer maximum temperature. *Atmos. Chem. Phys.*, **20**, 8251–8266, <https://doi.org/10.5194/acp-20-8251-2020>.
- Turner, D. D., and Coauthors, 2004: The QME AERI LBLRTM: A closure experiment for downwelling high spectral resolution infrared radiance. *J. Atmos. Sci.*, **61**, 2657–2675, <https://doi.org/10.1175/JAS3300.1>.
- , A. Merrelli, D. Vimont, and E. J. Mlawer, 2012: Impact of modifying the longwave water vapor continuum absorption model on Community Earth System Model simulations. *J. Geophys. Res.*, **117**, D04106, <https://doi.org/10.1029/2011JD016440>.
- , and Coauthors, 2020: A verification approach used in developing the Rapid Refresh and other numerical weather prediction models. *J. Oper. Meteor.*, **8**, 39–53, <https://doi.org/10.15191/nwajom.2020.0803>.
- Ukkonen, P., 2022: Exploring pathways to more accurate machine learning emulation of atmospheric radiative transfer. *J. Adv. Model. Earth Syst.*, **14**, e2021MS002875, <https://doi.org/10.1029/2021MS002875>.
- , R. Pincus, R. J. Hogan, K. P. Nielsen, and E. Kaas, 2020: Accelerating radiation computations for dynamical models with targeted machine learning and code optimization. *J. Adv. Model. Earth Syst.*, **12**, e2020MS002226, <https://doi.org/10.1029/2020MS002226>.
- Veerman, M., R. Pincus, R. Stoffer, C. V. Leeuwen, D. Podareanu, and C. V. Heerwaarden, 2020: Predicting atmospheric optical properties for radiative transfer computations using neural networks. *2020 EGU General Assembly*, Online, European Geoscience Union, Abstract EGU2020-5574, <https://doi.org/10.5194/egusphere-egu2020-5574>.
- Wallace, J. M., and P. V. Hobbs, 2006: *Atmospheric Science: An Introductory Survey*. 2nd ed. Academic Press, 504 pp.
- Yang, C., J. Chiu, J. Gristey, G. Feingold, and W. Gustafson, 2022: Machine learning emulation of 3D shortwave radiative transfer for shallow cumulus cloud fields. *Conf. on Atmospheric Radiation, Atmospheric Radiative Transfer, and Light-Scattering Theory*, Madison, WI, Amer. Meteor. Soc., 7.4, <https://ams.confex.com/ams/CMM2022/meetingapp.cgi/Paper/406293>.
- Zhou, Z., M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, 2020: UNet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE Trans. Med. Imaging*, **39**, 1856–1867, <https://doi.org/10.1109/TMI.2019.2959609>.