Machine-learned uncertainty quantification is not magic: Lessons learned from emulating radiative transfer with ML

Ryan Lagerquist^{1,2}, Imme Ebert-Uphoff^{1,2,3}, David D Turner², and Jebb Q. Stewart²

¹Cooperative Institute for Research in the Atmosphere (CIRA), Colorado State University ²National Oceanic and Atmospheric Administration (NOAA) Global Systems Laboratory (GSL)

³Department of Electrical and Computer Engineering, Colorado State University

November 14, 2023

Abstract

Machine-learned uncertainty quantification (ML-UQ) has become a hot topic in environmental science, especially for neural networks. Scientists foresee the use of ML-UQ to make better decisions and assess the trustworthiness of the ML model. However, because ML-UQ is a new tool, its limitations are not yet fully appreciated. For example, some types of uncertainty are fundamentally unresolvable, including uncertainty that arises from data being out of sample, *i.e.*, outside the distribution of the training data. While it is generally recognized that ML-based point predictions (predictions without UQ) do not extrapolate well out of sample, this awareness does not exist for ML-based uncertainty. When point predictions have a large error, instead of accounting for this error by producing a wider confidence interval, ML-UQ often fails just as spectacularly. We demonstrate this problem by training ML with five different UQ methods to predict shortwave radiative transfer. The ML-UQ models are trained with real data but then tasked with generalizing to perturbed data containing, *e.g.*, fictitious cloud and ozone layers. We show that ML-UQ completely fails on the perturbed data, which are far outside the training distribution. We also show that when the training data are lightly perturbed – so that each basis vector of perturbation has a *little* variation in the training data – ML-UQ can extrapolate along the basis vectors with some success, leading to much better (but still somewhat concerning) performance on the validation and testing data. Overall, we wish to discourage overreliance on ML-UQ, especially in operational environments.

Machine-learned uncertainty quantification is not magic: Lessons learned from emulating radiative transfer with ML

1

2

3

6

7

13

Key Points:

This article was submitted to the Journal for Advances in Modeling Earth Systems (JAMES) on Nov 8 2023.

Ryan Lagerquist^{1,2*}, Imme Ebert-Uphoff^{1,3}, David D. Turner², and Jebb Q. Stewart²

8	1 Cooperative Institute for Research in the Atmosphere (CIRA), Colorado State University, Fort Collins,
9 10	$\begin{tabular}{lllllllllllllllllllllllllllllllllll$
11 12	Colorado ³ Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, Colorado

14	• Machine-learned uncertainty quantification (ML-UQ) is a new tool, and its lim-
15	itations are not yet fully appreciated.
16	• Just like deterministic ML, ML-UQ often does not extrapolate well outside of the
17	training domain.
18	• However, this problem can be somewhat mitigated by perturbing the training data
19	along important basis vectors.

 $^{^*325}$ Broadway, R/GSL6, Boulder, CO 80305

 $Corresponding \ author: \ Ryan \ Lagerquist, \verb"ralager@colostate.edu"$

20 Abstract

Machine-learned uncertainty quantification (ML-UQ) has become a hot topic in envi-21 ronmental science, especially for neural networks. Scientists foresee the use of ML-UQ 22 to make better decisions and assess the trustworthiness of the ML model. However, be-23 cause ML-UQ is a new tool, its limitations are not yet fully appreciated. For example, 24 some types of uncertainty are fundamentally unresolvable, including uncertainty that arises 25 from data being out of sample, *i.e.*, outside the distribution of the training data. While 26 it is generally recognized that ML-based point predictions (predictions without UQ) do 27 not extrapolate well out of sample, this awareness does not exist for ML-based uncer-28 tainty. When point predictions have a large error, instead of accounting for this error 29 by producing a wider confidence interval, ML-UQ often fails just as spectacularly. We 30 demonstrate this problem by training ML with five different UQ methods to predict short-31 wave radiative transfer. The ML-UQ models are trained with real data but then tasked 32 with generalizing to perturbed data containing, e.g., fictitious cloud and ozone layers. 33 We show that ML-UQ completely fails on the perturbed data, which are far outside the 34 training distribution. We also show that when the training data are lightly perturbed 35 - so that each basis vector of perturbation has a *little* variation in the training data -36 ML-UQ can extrapolate along the basis vectors with some success, leading to much bet-37 ter (but still somewhat concerning) performance on the validation and testing data. Over-38 all, we wish to discourage overreliance on ML-UQ, especially in operational environments. 39

⁴⁰ Plain-language summary

Machine-learned uncertainty quantification (ML-UQ) - i.e., ML models that re-41 turn both a point prediction and an estimate of their own uncertainty – is a hot topic 42 in environmental science. Recent developments in ML-UQ have generated much excite-43 ment, but this excitement should be tempered by an awareness of its limitations. For 44 example, just like basic ML (with only point predictions) extrapolates poorly outside the 45 distribution of its training data, so do uncertainty estimates from ML-UQ. This can lead 46 to catastrophic errors, *i.e.*, very wrong predictions made with high confidence (low un-47 certainty). We demonstrate this problem across a range of ML-UQ methods and address 48 a way to alleviate the problem. 49

50 1 Introduction

51

1.1 Machine-learned uncertainty in environmental science

For as long as machine learning (ML) has been used in environmental science (ES), 52 both developers and users have been interested in how *uncertain* the predictions are. This 53 uncertainty quantification (UQ) is especially important for high-impact applications, such 54 as severe weather, where an incorrect prediction can cost property and human lives. The 55 computer-science literature has recently made breakthroughs in ML models that quan-56 tify their own uncertainty (ML-UQ), which could be a game-changer for high-impact ML 57 applications in ES. The next step is for the ES community to familiarize itself with these 58 new ML-UQ tools and modify them to best suit the unique needs of ES applications. This 59 work is already in progress (Rasp et al., 2018; Wimmers et al., 2019; Scheuerer et al., 60 2020; Baran & Baran, 2021; Bihlo, 2021; Barnes et al., 2021; Clare et al., 2021; Ghazvinian 61 et al., 2021; Orescanin et al., 2021; Scher & Messori, 2021; Veldkamp et al., 2021; Chap-62 man et al., 2022; Garg et al., 2022; Klotz et al., 2022; Ortiz et al., 2022; Schulz & Lerch, 63 2022). Specifically, the ES community is asking the following questions: 64

- 1. Which ML-UQ methods are reasonably easy to learn and implement?
- 65 66 67
 - 2. Which methods are available for classification (predicting a category) vs. regres-
 - sion (predicting a continuous real number)? The computer-science community of-

ten develops methods for classification tasks, whereas many ES problems are regression tasks.

- 703. What is the best ML-UQ method for a given application? Just like the quality of71point predictions (*i.e.*, single predictions without UQ, often called "determinis-72tic") can be evaluated with objective tools, so can the quality of uncertainty es-73timates. See Haynes et al. (2023, henceforth H23) for an overview of UQ evalu-74ation.
- We are particularly interested in the last question. Specifically, we take one step
 further back and ask:
- 77 78

79

80

81

82

68

69

- 1. Are there fundamentally *unresolvable* types of uncertainty (*i.e.*, that cannot be captured with any UQ method)?
- 2. If so, what real-world scenarios create unresolvable uncertainty? What are the implications for using ML-UQ in operations? For example, how should the uncertainty estimates be interpreted, knowing that they might completely miss some types of uncertainty?

We are interested in these questions because we foresee a danger of scientists relying too heavily on ML-UQ.

85

1.2 ML-UQ is not magic

To understand why we are concerned, let us briefly recap the state of ML in ES. 86 ML has shown great promise in terms of improved accuracy and faster execution, rel-87 ative to process-based models. Although these advantages have been demonstrated for 88 many ES applications, users, such as operational weather-forecasters, have been slow to 89 accept ML into operations (Gil et al., 2019; Reichstein et al., 2019). The primary rea-90 son is that ML is not guaranteed to generalize well to out-of-sample data (Buiten, 2019) 91 e.q., locations, seasons, or physical regimes that were not included in the training data. 92 In contrast, process-based models, which employ known laws of physics, typically gen-93 eralize much better out-of-sample. Also, where process-based models make an approx-0/ imation, users generally understand the potential impacts - e.g., situations where the model is thereby inappropriate. This understanding is much harder to build for ML mod-96 els. 97

One hope is that recently developed ML-UQ methods can help indicate situations 98 where an ML model is inappropriate, *i.e.*, where it will produce unacceptable errors. Howqq ever, this hope rests on the assumption that the model's estimates of its own uncertainty 100 are highly correlated with its error -i.e., that the model "knows when it is wrong". It 101 is our subjective experience that scientists do not question an ML model's uncertainty 102 estimates as much as they question its predictions. In particular, scientists do not con-103 sider the possibility of *catastrophic errors*: extremely wrong predictions made with high 104 confidence (low uncertainty). The concept of catastrophic errors, especially arising due 105 to unresolvable uncertainty, is at the core of this manuscript. 106

107

1.3 A few examples of unresolvable uncertainty

An ML-based UQ method (*e.g.*, Bayesian neural network) must ground its uncertainty estimates in the training data, just like the base ML model (*e.g.*, neural network) must ground its predictions in the training data. No other information is provided to the model. Thus, if a physical relationship exists in the real world but is not represented in the training data, it will not be learned by the base model or ML-UQ method. From this insight, we construct three scenarios that *any* ML-UQ method would struggle with.

Scenario 1: Missing variable. The target variable y depends strongly on a vari-114 able not included in the predictors. This scenario is common, as some variables cannot 115 be measured and a limited number of variables can be included in an ML model. Specif-116 ically, consider an example where y is a function of two variables, x_{known} and x_{unknown} , 117 but only x_{known} is included in the predictors. Let the model be f with a probabilistic 118 output vector \vec{y}_{pred} , which represents the full predicted distribution. (For example, if f 119 is an ensemble model, each element of $\vec{y}_{\rm pred}$ is one member of the ensemble; if \hat{f} is a para-120 metric model assuming the normal distribution, the two elements of \vec{y}_{pred} are mean and 121 variance; ...; etc.) 122

$$\begin{cases} y_{\text{true}}(x_{\text{known}}, x_{\text{unknown}}) = x_{\text{known}} \cdot x_{\text{unknown}}; \\ \vec{y}_{\text{pred}}(x_{\text{known}}) = \hat{f}(x_{\text{known}}). \end{cases}$$
(1)

Since y_{true} depends strongly on x_{unknown} but the model does not have access to x_{unknown} , the distribution \vec{y}_{pred} – including any point prediction (*e.g.*, the mean) and any measure of uncertainty (*e.g.*, the variance) – will lack skill. In other words, the model's point predictions will be poor, and the tool designed to alert us when point predictions are poor – namely UQ – will fail as well. Although this example is extreme, unresolvable uncertainty can arise in other ways. The point of this example is to illustrate that *any* UQ method will fail to alert us to the model's poor predictions.

130 Scenario 2: Variable constant in training data. y depends strongly on a vari-131 able x_c that, although it is included in the predictors, takes a constant value over all the 132 training data. In general, though, x_c is not constant. Specifically, consider an example 133 with one other predictor, x:

$$\begin{cases} y_{\text{true}}(x, x_c) = x \cdot x_c; \\ \vec{y}_{\text{pred}}(x, x_c) = \hat{f}(x, x_c) = \hat{f}(x). \end{cases}$$
(2)

Although both x and x_c are provided to the model \hat{f} , it cannot learn anything from a 134 variable that does not actually vary in the training data. Replacing x with x_{known} and 135 x_c with $x_{unknown}$, Equation 2 becomes Equation 1, so uncertainty arising from x_c is un-136 resolvable. For a more intuitive example, consider a climate model trained to predict global-137 annual-averaged temperature (GAAT), with one of the predictors being CO_2 concentra-138 tion (q). All training samples contain the year-2000 value, q = 370 ppm; but the model 139 is then applied to year-2100 data, with q = 600 ppm. The year-2100 data are out of sam-140 ple with respect to q, leading to unresolvable uncertainty and catastrophic errors. Specif-141 ically, the model will severely underpredict GAAT with high confidence. 142

Scenario 3: Missing basis vector in training data. y depends strongly on vari-143 ations along a basis vector b of the predictor space, but the training data contain no vari-144 ation along this direction. Scenario 3 is a more general example of scenario 2, where \hat{b} 145 $= \hat{x}_c$. Scenario 2 is unlikely because it is easy to spot (e.g., by plotting a histogram of 146 every predictor variable), whereas scenario 3 is hard to spot, because it is hard to know 147 all the important basis vectors in a dataset, especially for high-dimensional data. As our 148 experiments in Sections 5 and 6 show, if an important basis vector (e.q., thickness of the 149 ozone layer) is not well sampled in the training data, this can lead to catastrophic out-150 of-sample errors. (We use the term *basis vector* loosely; in the strict definition all ba-151 sis vectors of a space are orthogonal, which is not necessarily true in our data. The term 152 "latent variable" or "latent-space vector" would be more accurate (Van et al., 2020), since 153 latent spaces do not imply orthogonality, but we feel that "basis vector" is more famil-154 iar to an ES audience.) 155

156

1.4 Our sample application: Shortwave radiative transfer

Simulating radiative transfer (RT) - i.e., heating of the atmosphere due to the scattering and absorption of radiation by particles such as hydrometeors, water vapour, aerosols, and trace gases – is a key part of numerical weather prediction (NWP). However, existing RT models are computationally expensive, which slows down NWP. In previous work (Lagerquist et al., 2023, henceforth L23) we demonstrated that one of these models, namely the Rapid Radiative Transfer Model (RRTM; Iacono et al., 2008), can be emulated accurately and quickly with neural networks. The current study builds on L23 and focuses entirely on UQ, which is not included in L23 or the RRTM. We focus on shortwave radiation (wavelengths of 0.2-12.2 μ m), which is largely of solar origin.

Note that the goal of this paper is *not* to generate new insights for the application of emulating RT. Rather, we use this application because it is an ideal setup to experiment with the scenarios of unresolvable uncertainty discussed above. Because we are using ML to emulate another model (the RRTM), we can freely modify the inputs (predictors) and use the RRTM to compute the corresponding correct outputs (targets). We use this setup to explore the following questions:

- 172 1. Can we observe the theoretical scenarios from Section 1.3 in practice? *i.e.*, Can we cause ML-UQ to fail catastrophically for such scenarios?
- How drastic are the failures in practice? Do some UQ methods fare better than
 others? Are there simple ways to address the failures?
- 1.5 Organization of this manuscript

First, we create out-of-sample data that should confound *any* ML-UQ method. Specifically, we perturb the validation and testing data along several basis vectors (*e.g.*, thickness of the ozone layer) that have very little variability in the training data. Second, we train models with five different ML-UQ methods and apply them to the perturbed validation and testing data, verifying that all five methods produce catastrophic errors. Third, we explore whether we can reduce these catastrophic errors by perturbing the training data *just a little* along each basis vector.

¹⁸⁴ 2 Input data

This section is a brief overview of the predictor and target variables, referring to L23 for details. The RRTM and ML-based emulators have the same target variables and mostly the same predictor variables; the emulators have two extra predictors, for reasons discussed in Section 2a of L23. For the target variables, values produced by the RRTM are considered ground truth, or "labels" in ML terminology.

¹⁹⁰ 2.1 Predictor variables

We use 26 predictor variables, summarized in Table 1. Most of these variables are 191 available in output files from version 16 of the Global Forecast System model (GFSv16; 192 see 2021 update at https://www.emc.ncep.noaa.gov/emc/pages/numerical_forecast 193 _systems/gfs/documentation.php), but a few are not. For these synthetic variables, 194 we create fictitious data, following Section 2b of L23. For the GFSv16 variables, we ex-195 tract forecast profiles at locations around the globe from 0000 UTC model runs on dates 196 from Sep 1 2018 to Dec 23 2020. Thus, our dataset is global in terms of both geographic 197 location and seasonality -i.e., covers all times of year at all locations. For more details 198 on the GFSv16 variables, see Section 2a of L23. 199

Table 1: Description of predictor variables. "Scalar?" indicates whether the variable is scalar, versus a full profile. "Synthetic?" indicates whether the values are synthesized from fake data, versus taken from GFSv16 output. "ML only?" indicates whether the variable is used only in the ML-based emulators, versus both ML and RRTM. "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. Downward IWP, upward IWP, downward WVP, and upward WVP have analogous definitions.

Variable	Units	Scalar?	Synthetic?	ML only?
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-vapour 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			
Solar zenith angle	0	\checkmark		
Surface albedo	_	\checkmark		
Height thickness	m			\checkmark
Pressure thickness	Pa			\checkmark
Aerosol single-scattering albedo		\checkmark	\checkmark	
Aerosol asymmetry parameter	_	\checkmark	\checkmark	
Aerosol extinction coefficient	m ⁻¹		\checkmark	
Liquid effective radius	m		\checkmark	
Ice effective radius	m		\checkmark	
N_2O concentration	ppmv		\checkmark	
CH_4 concentration	ppmv		\checkmark	
CO_2 concentration	ppmv		\checkmark	

2.2 Target variables

200

The RRTM performs 1-dimensional RT, assuming that RT occurs only in the vertical. Thus, both the RRTM and emulators are applied to each profile separately. The target variables are those required by an NWP model from its RT parameterization: a full profile of heating rates (HR), surface downwelling flux ($F_{\rm down}^{\rm sfc}$), top-of-atmosphere upwelling flux ($F_{\rm up}^{\rm TOA}$), and net flux ($F_{\rm net}$). See Figure 1 for an example.



Figure 1: RRTM outputs for one data sample. We emulate the full profile of heating rates, $F_{\text{down}}^{\text{sfc}}$ (the bottom value in the green curve), $F_{\text{up}}^{\text{TOA}}$ (the top value in the purple curve), and F_{net} (the difference between the last two values).

2.3 Preparing the data for ML

206

Our data preparation includes three steps. First, we split the data into three tem-207 porally independent partitions: training, validation, and testing (Table 2). We use the 208 training data to optimize parameters (weights and biases) for each ML model, the val-209 idation data to select the best ML model (e.g., best UQ method), and the testing data 210 for a final assessment of the selected model. Second, we perturb the data in each par-211 tition to a different extent: the training data not at all (for the first experiment) or lightly 212 (for the second experiment), the validation data moderately, and the testing data heav-213 ily. In other words, the ML models are trained with clean or lightly perturbed data, se-214 lected based on moderately perturbed data, and then tasked with generalizing to heav-215 ily perturbed data. Third, we normalize each predictor variable from physical units to 216 z-scores, following Section 3b of Lagerquist et al. (2021). 217

Data subset	Time period		Number of days	Sample size
Training	Sep 1 2018 – Dec 21 2019		237	873 086
Validation	Jan 2-15 2020, Mar 24 – Apr 6 2020, Jun 16-29 2020, Sep 6-19 2020, Nov 30 – Dec 13 2020	Feb 12-25 2020, May 5-18 2020, Jul 27 – Aug 9 2020, Oct 19 – Nov 2 2020,	126	479 806
Testing	Jan 18-31 2020, Apr 9-22 2020, Jul 2-15 2020, Sep 22 - Oct 7 2020,	Feb 28 – Mar 12 2020, May 22 – Jun 4 2020, Aug 12-25 2020, Nov 5-18 2020,	120	474 726

Table 2: Partitioning of data into temporally independent subsets. "Sample size" = number of profiles. Also, "Number of days" \neq length of "Time period," because some days are missing from the archive.

2.4 Perturbing to create out-of-sample data

218

Dec 16-23 2020

We create out-of-sample data by perturbing five atmospheric properties represented 219 in the predictor variables. The five properties are near-surface temperature, near-surface 220 humidity, liquid cloud, ice cloud, and ozone. Loosely, each property may be seen as cor-221 responding to one or more basis vectors of the predictor space. Some of our perturba-222 tions – e.g., increasing near-surface temperature and humidity – mimic impacts that are 223 expected from climate change, a real process that creates out-of-sample data. Some re-224 searchers have developed methods to make ML more robust to climate change (Beucler 225 et al., 2021), albeit with a focus on point predictions rather than uncertainty estimates. 226 However, some of our perturbations - e.g., those involving the ozone layer - are unlike 227 anything seen in the Earth's atmosphere or expected with climate change. Supplemen-228 tal Figures S5-S9 show the distribution of each variable before and after perturbation; 229 here it is evident, for example, that the changes to ozone are much more extreme than 230 the changes to temperature and humidity. These *extreme* perturbations allow us to ob-231 serve the behaviour of the UQ methods when tasked with generalizing to *extremely* out-232 of-sample data. In other words, the more extreme perturbations allow us to stress-test 233 the UQ methods in a way that more realistic data would not. 234

The target values -i.e., heating rates and fluxes - must change in response to the new predictors. To obtain the new target values $(\vec{y'})$ for a given profile, we simply feed the new predictors $(\vec{x'})$ to the RRTM.

Two details remain to be specified: [1] Which atmospheric properties are perturbed for which profiles? [2] What are the specific perturbation methods? For each profile Pand each property χ , there is a 50% chance that χ will be perturbed in P, based on drawing a random integer from $\{0, 1\}$. For a given profile P, if all five random numbers evaluate to 0, one of the five is randomly changed to 1, so that at least one property is perturbed for every profile. The subsections below explain the specific perturbation method for each atmospheric property.

245 Near-surface temperature

256

257

258

259

Our motivation for this procedure is to mimic the lower-tropospheric warming expected with climate change. The procedure has two parameters: maximum depth of the warm layer (D_{max}) and maximum surface-temperature increase $(\Delta T_{\text{sfc}}^{\text{max}})$. For the lightly perturbed training data, we set $D_{\text{max}} = 1.25$ km and $\Delta T_{\text{sfc}}^{\text{max}} = 2$ K; for the moderately perturbed validation data, $D_{\text{max}} = 2.5$ km and $\Delta T_{\text{sfc}}^{\text{max}} = 4$ K; for the heavily perturbed testing data, $D_{\text{max}} = 5$ km and $\Delta T_{\text{sfc}}^{\text{max}} = 8$ K. The procedure is shown schematically in Figure 2. After the numbered procedure below, we recompute relative humidity, based on the new temperature and untouched specific humidity.

- 1. Determine the depth of the warm layer by sampling from a uniform distribution
- over $[0, D_{\max}]$. Symbolically, $D \in \mathcal{U}[0, D_{\max}]$.
 - 2. Sample to determine the surface-temperature increase: $\Delta T_{\rm sfc} \in \mathcal{U}[0, \Delta T_{\rm sfc}^{\rm max}]$.
 - 3. At each height in the warm layer, scale the temperature increase linearly from $\Delta T_{\rm sfc}$ at the surface to 0 at height *D* above the surface. See Figures 2a-c.
 - 4. If step 3 led to any temperature above 60 °C, reduce to 60 °C. See Figure 2d.



Figure 2: Procedure for perturbing near-surface temperature. Panel c = a + b. In this example, the warm-layer depth D is 3 km and the surface-temperature increase $\Delta T_{\rm sfc}$ is 8 K.

260 Near-surface humidity

Our motivation is to mimic the lower-tropospheric moistening expected with climate change. We first generate a disturbance for the relative humidity (RH) profile, then recompute the other moisture variable (specific humidity) from the new RH. See Supplemental Section 1 for details.

265 Liquid cloud

²⁶⁶ Our motivation is to create more complex cloud arrangements, as well as denser ²⁶⁷ and deeper clouds, than seen in the real atmosphere. We completely replace the liquid-²⁶⁸ water content (LWC) profile, generating a number of cloud layers from 0 up to N_{max} . ²⁶⁹ N_{max} varies from 2 for the lightly perturbed training data to 5 for the heavily perturbed ²⁷⁰ testing data. See Supplemental Section 1 for details.

271 Ice cloud

The motivation for perturbing ice cloud is the same as for perturbing liquid cloud; the two procedures are nearly identical. See Supplemental Section 1 for details.

274 **Ozone**

Our motivation is to create more complex ozone layers – over a wider range of locations, depths, and mixing ratios – than seen in the real atmosphere. We completely replace the ozone mixing ratio (w) profile, generating an ozone layer with a random location, depth, and structure. See Supplemental Section 1 for details.

²⁷⁹ **3** Methods

280

3.1 The base model: U-net++

The field of deep learning has produced many specialized neural network (NN) ar-281 chitectures for handling spatial data. We have chosen the U-net++ architecture, which 282 L23 found to be the best for shortwave RT. The U-net++ is a slight generalization of 283 the U-net (Ronneberger et al., 2015), which is designed for image-to-image translation, *i.e.*, to output predictions on the same spatial grid as the predictors. The U-net contains 285 four key components: convolutional layers, pooling (downsampling) layers, upsampling 286 layers, and skip connections. Convolutional layers use learned image filters to detect spa-287 tial and multivariate features in the predictor data, producing abstract representations 288 of the predictor data, called "feature maps". Pooling and upsampling layers scale fea-289 ture maps to coarser and finer spatial resolutions, respectively, allowing convolutional 290 layers to detect features at different scales. Skip connections carry high-resolution fea-291 ture maps directly across the network, by passing the series of downsampling and upsam-292 pling layers, which is a lossy operation that degrades high-resolution information. The 293 U-net++ (Zhou et al., 2019) is a U-net with more skip connections, which more effec-294 tively preserve small-scale features, such as cloud boundaries, that are important for short-295 wave RT. Our specific U-net++ setup for point prediction is shown in Figure 3. Our main 296 learning task is to translate a 127-by-26 image of predictor variables into a 127-by-1 im-297 age of heating rates. (There are 127 heights in the GFS grid and 26 predictor variables; 298 see Table 1. We duplicate the 4 scalar variables over all 127 heights, so that they fit into 299 the matrix.) There is also a second learning task: to predict the three flux variables $(F_{\text{down}}^{\text{sfc}},$ 300 $F_{\rm up}^{\rm TOA}$, and $F_{\rm net}$), which are scalars rather than images. For this we attach fully connected 301 layers – which are used in traditional (non-convolutional) NNs (Chapter 6 of Goodfel-302 low et al., 2016) and still a popular choice for scalar data – to the U-net++. 303



Figure 3: [adapted from Figure 3a of L23] Our specific U-net++ setup for point prediction. For each set of feature maps (green box), the label is number_of_heights × number_of_channels. In the remaining discussion, let K be the number of convolutional layers per block. We call this hyperparameter "width" in L23; the chosen value in this study, based on L23, is K = 1. Each orange "convolution" arrow represents K convolutional layers with 3-pixel filters; each "downsampling" arrow represents K convolutional layers with 3-pixel filters, followed by maximum-pooling with a 2-pixel window; each "upsampling" arrow represents upsampling with a 2-pixel window, followed by a convolutional layer with 3-pixel filters; each "skip connection" arrow includes K convolutional layers with 3-pixel filters; each black "convolution" arrow represents one convolutional layer with 1-pixel filters; and finally, each "fully connected layer" arrow represents one fully connected layer.

304

3.2 The ML-UQ methods

The total uncertainty in an ML model is the sum of two components: aleatory and 305 epistemic. The Appendix provides definitions of these terms – which, interestingly, dif-306 fer across disciplines – and shows that the examples of unresolvable uncertainty from 307 Section 1.3 can show up in both components. Thus, our analysis must include UQ meth-308 ods that can capture both the aleatory and epistemic components of uncertainty. Specif-309 ically, we use the three UQ methods discussed in the subsections. The first method (CRPS-310 LF) was found by H23 to perform well, but it can capture only aleatory uncertainty. Thus, 311 we also use the multi-model ensemble (MME) and Bayesian neural networks (BNN). On 312 their own MME and BNN can capture only epistemic uncertainty, but either method can 313 be combined with CRPS-LF to capture both types of uncertainty. 314

How to read this section: Our purpose for testing multiple UQ methods is to show that our results generalize across UQ methods. The interested reader may continue with this section and see H23 for even more details; readers less interested in the inner workings of UQ methods may skip ahead to Section 4.

319 3.2.1 CRPS-LF

This approach involves training a NN with the continuous ranked probability score (CRPS) as the loss function (LF). The CRPS-LF approach can be used with both parametric prediction and ensemble prediction. In ensemble prediction, the NN approximates y_{pred} by generating an ensemble, and its loss function is the ensemble formulation of the CRPS:

$$CRPS = \frac{1}{N} \sum_{i=1}^{N} |y_{true} - y_{pred}^{i}| - \frac{1}{2} \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} |y_{pred}^{i} - y_{pred}^{j}|, \qquad (3)$$

where N is the ensemble size; y_{true} is the correct value; and y_{pred}^k is the k^{th} prediction in the ensemble. The first term is the mean absolute error (MAE), and the second is the mean absolute pairwise difference (MAPD) between ensemble members, a measure of spread. The CRPS ranges from $[0, \infty)$; the optimal value is 0.

The CRPS is an uncertainty-oriented generalization of the MAE, which is a standard loss function for point prediction. However, for point prediction we use a custom loss function to emphasize large heating rates (Section 3d of L23), which the NN predicts poorly when trained with standard loss functions. Specifically, we use the following loss function for point prediction:

$$\mathcal{L} = \frac{1}{H} \sum_{h=1}^{H} \max\left\{ |r_h|, |\hat{r}_h| \right\} (r_h - \hat{r}_h)^2 + \frac{1}{L} \sum_{l=1}^{L} (F_l - \hat{F}_l)^2,$$
(4)

where H = 127 is the number of heights; r_h is the actual heating rate at the h^{th} height; \hat{r}_h is the corresponding prediction; L = 3 is the number of flux variables; F_l is the actual value of the l^{th} flux variable; and \hat{F}_l is the corresponding prediction. The second term is the standard MSE for flux variables, but the first term is a weighted MSE for heating rates, the weight being max $\left\{ |r_h|, |\hat{r}_h| \right\}$. We call this term the dual-weighted MSE (DWMSE).

To generalize the above loss function for UQ, we hybridize Equations 3 and 4, yielding the dual-weighted CRPS (DWCRPS):

$$DWCRPS = \frac{1}{H} \frac{1}{N} \sum_{h=1}^{H} \sum_{i=1}^{N} \max\left\{ |r_h|, |\hat{r}_{hi}| \right\} |r_h - \hat{r}_{hi}| - \frac{1}{2} \frac{1}{H} \frac{1}{N^2} \sum_{h=1}^{H} \sum_{i=1}^{N} \sum_{j=1}^{N} \max\left\{ |\hat{r}_{hi}|, |\hat{r}_{hj}| \right\} |\hat{r}_{hi} - \hat{r}_{hj}|.$$
(5)

³⁴² $H, N, \text{ and } r_h \text{ are as defined previously; } \hat{r}_{hk} \text{ is the } k^{\text{th}} \text{ predicted heating rate at the } h^{\text{th}}$ ³⁴³ height; max $\left\{ |r_h|, |\hat{r}_{hi}| \right\}$ is the maximum absolute value of the actual and i^{th} predicted ³⁴⁴ heating rate at the h^{th} height; and max $\left\{ |\hat{r}_{hi}|, |\hat{r}_{hj}| \right\}$ is the maximum absolute value of ³⁴⁵ the i^{th} and j^{th} predicted heating rates at the h^{th} height. Both max terms are weights ³⁴⁶ that emphasize large heating rates.

The DWCRPS is used only for heating rates; the standard CRPS is used for fluxes, since the distribution of fluxes is less skewed (Figure 5 of Lagerquist et al., 2021) and therefore does not necessitate a custom loss function to ensure good prediction of extreme values. Thus, the total loss function we use for the CRPS-LF approach is:

$$\mathcal{L} = \text{DWCRPS} + \frac{1}{L} \frac{1}{N} \sum_{l=1}^{L} \sum_{i=1}^{N} |F_l - \hat{F}_{li}| - \frac{1}{2} \frac{1}{L} \frac{1}{N} \sum_{l=1}^{L} \sum_{i=1}^{N} \sum_{j=1}^{N} |\hat{F}_{li} - \hat{F}_{lj}|.$$
(6)

The second and third terms, collectively, are the CRPS for flux variables.

To make the CRPS-LF approach work, in addition to changing the loss function, one must change the NN architecture to output N estimates per target variable (Figure 4).



Figure 4: [adapted from Figure 3a of L23] Our specific U-net++ setup for UQ. If using the CRPS-LF approach, N (the ensemble size) > 1 and the loss function is Equation 6; otherwise, N = 1 and the loss function is Equation 4. If using the BNN approach, one or more of the double arrows contain Bayesian layers. The arrows marked "fully connected layer" and "convolution with 1-px filters" each represent a single layer; the corresponding layer may or may not be Bayesian. Meanwhile, recall from the caption of Figure 3 that

the arrows marked "upsampling" and "skip connection" each contain K convolutional layers. If an upsampling or skip connection is made Bayesian, then all convolutional layers therein are Bayesian.

355 3.2.2 Multi-model ensemble

360

The idea behind the multi-model ensemble (MME) is simple: train many pointprediction NNs, each with a different random seed, then ensemble the predictions. The random seed determines how the NN weights are initialized, and different initializations lead to different solutions, the "solution" being the final set of weights.

3.2.3 Bayesian neural networks

Any NN can be made Bayesian by replacing traditional (point-prediction) layers 361 with Bayesian layers; thus, BNNs are a highly flexible approach to UQ. A point-prediction 362 NN learns a single value for each weight, but a BNN learns a full *distribution* for some 363 weights, determined by fitting the parameters of a user-chosen canonical distribution. 364 In this work and in common practice, the normal distribution is chosen, so the BNN must 365 learn two values for each Bayesian weight: the mean and variance. It is unnecessary to 366 make all layers Bayesian. For example, a popular approach is to make only the last few 367 layers Bayesian, which often achieves the same performance (*i.e.*, quality of mean pre-368 dictions and uncertainty estimates) at a fraction of the computing cost (Jospin et al., 369 2022; Hertel et al., 2023). 370

While simple in theory, "making a layer Bayesian" is non-trivial in practice. To update a Bayesian weight w, one must compute the posterior distribution $p(w \mid D)$, where

 $\mathcal D$ represents the training data. Solving the posterior exactly involves a computationally 373 intractable integral, so in practice, variational inference is often used as an approxima-374 tion (Hoffman et al., 2013; Ranganath et al., 2014; Rezende et al., 2014). Furthermore, 375 weights in a NN are updated via gradient descent with backpropagation, which involves 376 the gradient of the loss with respect to every weight, $\frac{\partial \mathcal{L}}{\partial w}$. However, a Bayesian weight 377 is a random variable, and it is impossible to backpropagate the gradient through ran-378 dom variables. There are two popular solutions to this problem: the reparameterization 379 trick (Kingma & Welling, 2013), which involves loss gradients with respect to only the 380 *parameters* of the weight distribution (e.q., the mean and variance of a normal distri-381 bution), and flipout (Wen et al., 2018), which involves sampling weight perturbations. 382 The advantage of reparameterization is speed – per training epoch, it is faster than flipout 383 - while the advantage of flipout is more accurate gradient estimates. This accuracy of-384 ten translates to needing fewer training epochs, which can make flipout faster per net-385 work even though it is slower per epoch. 386

³⁸⁷ 4 Experimental setup

39 39

396

397

398

399

400

We attempt five UQ methods with the U-net++ base model: CRPS-only, MMEonly, MME/CRPS, BNN-only, and BNN/CRPS. Each of these methods generates an ensemble; for fair comparison across UQ methods, we set the ensemble size to 50. (Larger ensemble sizes lead to memory issues for training the BNN/CRPS models.) Specifically, we use the following techniques:

3	1. CRPS-only. Train a single U-net++ with the probabilistic loss function (Equa-
4	tion 6) and 50 output neurons per target variable $(N = 50 \text{ in Figure 4})$.
5	2. MME-only. Train 50 U-net++ models, each with the deterministic loss function

- (Equation 4) and 1 output neuron per target variable.
- 3. MME/CRPS. Train 50 U-net++ models, each with the probabilistic loss function and 25 output neurons per target variable. Hence, the inner ensemble size is 25 and the outer ensemble size is 50 – leading to a total ensemble size of 1250, from which we randomly select 50 members.
- 4. BNN-only. Train a single Bayesian U-net++ with the deterministic loss function and 1 output neuron per target variable. At inference time, run the Bayesian Unet++ 50 times to get 50 predictions per target variable.
- 5. BNN/CRPS. Train a single Bayesian U-net++ with the probabilistic loss function and 50 output neurons per target variable. At inference time, run the Bayesian U-net++ 10 times, so that each probabilistic weight is sampled 10 times. Hence, the inner ensemble size is 50 and the outer ensemble size is 10 leading to a total ensemble size of 500, from which we randomly select 50 members.
- For methods involving a BNN, this leaves the question of which layers are Bayesian (probabilistic weights) and which are not (deterministic weights), as well as which method to use for training Bayesian layers (reparameterization or flipout). For both the BNNonly and BNN/CRPS methods, we optimize these hyperparameters with an experiment documented in Supplemental Section 2.

Models trained with a single UQ method can capture only one type of uncertainty (aleatory for CRPS-only, epistemic for MME-only and BNN-only), while those trained with a hybrid method can capture both types of uncertainty. Since uncertainty arising from out-of-sample data is partly aleatory and partly epistemic, we expect the hybrid UQ methods to perform best on the perturbed (validation and testing) data.

419 Experiment 1: Clean training data

In this experiment we train the models with clean (unperturbed) data, then task them with generalizing – both point predictions and uncertainty estimates – to perturbed validation and testing data. We expect all UQ methods to perform poorly on the perturbed data, because the perturbations are made along basis vectors with much less variability in the clean training data (Supplemental Figures S5-S9).

Experiment 2: Lightly perturbed training data

The confirmation of the above expectation (see Section 5) motivates another science question: what happens if the models are trained with *lightly perturbed*, instead of clean, data? Said differently, what happens if the models "see" a light version of the perturbations occurring in the validation and testing data? On the out-of-sample validation and testing data, we expect models trained with lightly perturbed data to perform better than models trained with clean data, but *how much* better is an open question.

432 Tools for evaluating UQ results

425

This section provides a light background on UQ-evaluation tools (both graphics and single-number metrics), which should be sufficient for readers to understand the ensuing results and discussion. See H23 for more details.

Figure 5 demonstrates our evaluation tools for two synthetic datasets. The first dataset 436 (Figure 5a) represents a model with good mean predictions but too much spread (i.e.,437 ensemble ranges are wider than necessary); we call this Model A. The second dataset (Fig-438 ure 5b) represents a model with poor mean predictions and too little spread (i.e., the439 observation often falls completely outside the ensemble range); we call this Model B. The 440 attributes diagram – which is a reliability curve with extra information (Hsu & Mur-441 phy, 1986) – is used to evaluate point predictions, showing the mean observation $y_{\rm true}$ 442 as a function of the ensemble-mean prediction $\overline{y_{\text{pred}}}$. This graphic is used to assess con-443 ditional bias, *i.e.*, bias as a function of $\overline{y_{\text{pred}}}$. Model A has no conditional biases (Fig-444 ure 5c), leading to a reliability curve that follows the 1:1 line and a reliability (REL, the 445 mean squared distance between the curve and 1:1 line) of $0.00 \text{ K}^2 \text{ day}^{-2}$. Meanwhile, Model 446 447 B completely misses the extremes, *i.e.*, the lowest (highest) predictions are far too high (low). This leads to the classic sigmoid-shaped reliability curve and a large REL (Fig-448 ure 5d). 449



Figure 5: Demonstration of evaluation tools on two synthetic datasets. [a-b] The two synthetic datasets, representing "Model A" and "Model B". [c] Attributes diagram for Model A. Of the dashed grey lines: the diagonal (1:1) line represents the perfect reliability curve; the vertical line is the climatology line, representing the mean observation in the dataset (16 K day⁻¹); and the horizontal line is the no-resolution line, representing the reliability curve for a completely uninformative model. The blue shading is the positive-skill area, where the model's MSE is better than that yielded by always predicting climatology (here, 16 K day⁻¹). [d] Same but for Model B. [e] Spread-skill plot for Model A. The diagonal (1:1) line represents a perfect spread-skill curve; the grey histogram shows how often each spread value occurs; and the inset plot shows any biases as a function of model spread. [f] Same but for Model B. [g] Discard test for Model A. The inset plot shows any biases as a function of discard fraction. [h] Same but for Model B. [i] PIT histogram for Model A. The dashed line represents a perfect (uniform) PIT histogram. [j] Same but for Model B.

Remaining tools shown in Figure 5 are for uncertainty estimates rather than point 450 predictions. The spread-skill plot (Delle Monache et al., 2013) shows the root mean square 451 error (RMSE) achieved by $\overline{y_{\text{pred}}}$, or "skill," as a function of the ensemble standard de-452 viation, or "spread". For a perfect model, the spread-skill ratio is 1.0 across all spread 453 values, so the curve follows the 1:1 line. Model A (Figure 5e) is very overspread or "un-454 derconfident," leading to a curve well below the 1:1 line and a large overall spread-skill 455 ratio (SSRAT). Model B (Figure 5f) has the opposite problem. Spread-skill reliability 456 (SSREL), the mean distance between the curve and 1:1 line, is substantially lower (bet-457 ter) for Model A. 458

In the discard test, data samples are thrown out in descending order of model un-459 certainty (*i.e.*, the highest-uncertainty samples are thrown out first) and the effect on 460 model error is observed. The error should decrease monotonically, *i.e.*, whenever the dis-461 card fraction increases. For all discard tests in this paper, model error is based on the 462 ensemble mean $\overline{y_{\text{pred}}}$ and model uncertainty is the height-averaged variance of HR pre-463 dictions. (Mathematically, this is $\frac{1}{H} \sum_{h=1}^{H} \left[\frac{1}{N-1} \sum_{i=1}^{N} (\hat{r}_{hi} - \overline{\hat{r}_h})^2 \right]$, where $\overline{\hat{r}_h}$ is the ensem-464 ble mean at the h^{th} height; all other variables are defined in Equation 5. There are two 465 reasons that we use only HR, and not flux, to define overall uncertainty. First, most of 466 the model's outputs [127 of every 130] are HRs; second, combining HR and flux uncer-467 tainties into an overall uncertainty is non-trivial, as they have different units.) Model 468 A (Figure 5g) has a perfect discard test, leading to a monotonicity fraction (MF) of 100%. 469 Model B (Figure 5h) has an imperfect discard test; error increases as the discard frac-470 tion increases from 20-25%, from 25-30%, from 30-35%, and from 35-40%. Thus, model 471 error decreases only 15 of 19 times that the discard fraction increases, leading to an MF 472 of $\frac{15}{10} = 78.9\%$. 473

The probability integral transform (PIT), defined for each data sample, is the rank-474 ing of y_{true} in the distribution \vec{y}_{pred} . For example, if y_{true} is less than all \vec{y}_{pred} , its PIT 475 is 0.0; if y_{true} is greater than all \vec{y}_{pred} , its PIT is 1.0; if y_{true} is the median of all \vec{y}_{pred} , 476 its PIT is 0.5; etc. The PIT *histogram* – which is similar to the rank histogram, or "Ta-477 lagrand diagram" (Hamill, 2001), and can be interpreted similarly – then shows the dis-478 tribution of PIT values. For a perfectly calibrated model – which is neither overconfi-479 dent nor underconfident – all PIT values occur equally often, leading to a uniform his-480 togram. For Model A (Figure 5i), nearly all PIT values are between 0.3 and 0.7, mean-481 ing that y_{true} is usually in the middle 40% of the \vec{y}_{pred} distribution and rarely at the ex-482 tremes. In other words, the $\vec{y}_{\rm pred}$ distribution is usually too wide; the model is under-483 confident. For Model B (Figure 5j), nearly all PIT values are below 0.05 or above 0.95, 484 meaning that y_{true} is usually in the bottom or top 5% of the distribution. In other words, 485 the \vec{y}_{pred} distribution is usually too narrow; the model is overconfident. The PIT devi-486 ation (PITD), the mean absolute difference between bar height and the horizontal line, 487 is substantially better (lower) for Model A. (The horizontal line line denotes the bar height for the ideal [uniform] PIT histogram, which is $\frac{1}{\text{number of bins}}$.) 488 489

REL, SSREL, and PITD are negatively oriented with a perfect value of 0.0; MF
is positively oriented with a perfect value of 1.0; and the perfect SSRAT is 1.0, with higher
(lower) values indicating underconfidence (overconfidence).

Lastly, we define "large point error" as a data sample where $\overline{y_{\text{pred}}}$ has absolute error $\geq 1 \text{ K day}^{-1}$; we define "catastrophic error" as a large point error where the observation also falls outside the 95% confidence interval (in other words, PIT is either < 0.025 or > 0.975).

497 How to read the results

Sections 5 and 6, which discuss the results of the two experiments, contain many specific terms from the RT application and the field of UQ. All these terms have been

- explained briefly heretofore, with longer explanations found in L23 (for RT) and H23 (for
- ⁵⁰¹ UQ). However, we do not expect readers to be fluent in these terms, so we highlight "key
- points" throughout Sections 5 and 6. Readers with less interest in the details can jump
- ⁵⁰³ directly to the key points.

⁵⁰⁴ 5 Results for Experiment 1: Clean training data

We start with overall diagnostics (metrics computed from the entire validation or testing set), which allow us to understand the UQ methods' performance and choose the best method. Then we present a small number of case studies, which allow us to understand the UQ methods' performance in a way that overall diagnostics cannot.

509 5.1 Overall diagnostics

Figure 6 compares all five UQ methods on the *moderately perturbed* validation data. Error metrics pertaining to heating rates (HR) are averaged over the 127 heights; those pertaining to fluxes are averaged over the three variables $(F_{\text{down}}^{\text{sfc}}, F_{\text{up}}^{\text{TOA}}, \text{ and } F_{\text{net}})$.



Figure 6: Comparison of UQ methods on validation data, for models trained with clean data. In panel g, higher is better; in panel e, closer to 1.0 is better; in all other panels, lower is better. "CEF" in panel h is catastrophic-error frequency.

513	We highlight several observations from Figure 6:
514	1. UQ methods that can capture only epistemic uncertainty $-i.e.$, MME-only and
515	BNN-only – produce too little spread for both HR and fluxes (panel e). This sug-
516	gests that much of the uncertainty in the validation data is aleatory.
517	2. Methods that can capture aleatory uncertainty $-i.e.$, those involving the CRPS
518	– produce too much spread for fluxes (panel e). However, out of these three meth-
519	ods, BNN/CRPS is the least overspread.

 127 heights (not shown). 4. All UQ methods produce catastrophic errors at least ~10% of the time (panel h); the non-hybrid methods (CRPS-only, MME-only, and BNN-only) produce catas- trophic errors substantially more often. 5. The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	520	3.	All UQ methods produce far too little spread for HR (panel e); this is true at all
 4. All UQ methods produce catastrophic errors at least ~10% of the time (panel h); the non-hybrid methods (CRPS-only, MME-only, and BNN-only) produce catas- trophic errors substantially more often. 5. The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	521		127 heights (not shown).
 the non-hybrid methods (CRPS-only, MME-only, and BNN-only) produce catas- trophic errors substantially more often. 5. The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	522	4.	All UQ methods produce catastrophic errors at least $\sim 10\%$ of the time (panel h);
 trophic errors substantially more often. 5. The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	523		the non-hybrid methods (CRPS-only, MME-only, and BNN-only) produce catas-
 5. The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	524		trophic errors substantially more often.
 (flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error fre quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	525	5.	The BNN/CRPS method performs best on 4 of the 10 uncertainty-based metrics
 performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT, HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second- worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met- rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	526		(flux SSREL, flux SSRAT, flux PITD, and HR MF). The MME/CRPS method
 HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error free quency. However, MME/CRPS performs worst on flux MF (panel g) and second-worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based metrics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS produces better point predictions than BNN/CRPS; however, the differences here are small. 	527		performs best on 6 of the 10 uncertainty-based metrics (HR SSREL, HR SSRAT,
 quency. However, MME/CRPS performs worst on flux MF (panel g) and second-worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based metrics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS produces better point predictions than BNN/CRPS; however, the differences here are small. 	528		HR PITD, HR MF, HR CEF, and flux CEF), where CEF is catastrophic-error fre-
 worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based metrics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	529		quency. However, MME/CRPS performs worst on flux MF (panel g) and second-
 for any uncertainty-based metric. 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based metrics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS produces better point predictions than BNN/CRPS; however, the differences here are small. 	530		worst on flux SSRAT (panel e), while BNN/CRPS does not perform this badly
 6. MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based metrics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS produces better point predictions than BNN/CRPS; however, the differences here are small. 	531		for any uncertainty-based metric.
 rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro- duces better point predictions than BNN/CRPS; however, the differences here are small. 	532	6.	MME/CRPS outperforms BNN/CRPS on 3 of the 4 point-prediction-based met-
duces better point predictions than BNN/CRPS; however, the differences here are small.	533		rics (HR MAE, flux MAE, and flux REL; not HR REL). Thus, MME/CRPS pro-
535 small.	534		duces better point predictions than BNN/CRPS; however, the differences here are
	535		small.

Key points: Points 3 and 4 exemplify that, when trained with clean data, all UQ
methods fail dramatically. Based on points 5 and 6, we judge that the best UQ method
is BNN/CRPS, followed by MME/CRPS. Both are hybrid methods, which can capture
both aleatory and epistemic uncertainty. The remaining analysis will focus largely on
BNN/CRPS.



Figure 7: Detailed results of the BNN/CRPS method on validation data, for a model trained with clean data. [a-e] Evaluation of point predictions (ensemble means). [a]
Attributes diagram for F_{net}; [b] attributes diagram for HR, aggregated over all heights; [c] profile of mean absolute errors for HR; [d] profile of mean signed errors (biases) for HR;
[e] profile of large-point-error frequencies for HR. [f-h] Evaluation of uncertainty estimates

for F_{net} . [f] Spread-skill plot; [g] discard test; [h] PIT histogram. [i-l] Evaluation of uncertainty estimates for HR. [i] Profile of catastrophic-error frequencies for HR; [j-l] as in panels f-h but for HR.

Figure 7 shows detailed results for the BNN/CRPS model, based on the *moder*-541 ately perturbed validation data. For both $F_{\rm net}$ (panel a) and HR (panel b), the attributes 542 diagram is nearly perfect, except a large positive bias ($\sim 5 \text{ K day}^{-1}$) when $\overline{y_{\text{pred}}} \gtrsim 38 \text{ K}$ 543 day⁻¹. In other words, the highest HR predictions are too high. Panels c-d show MAE 544 and bias at each height for the ensemble-mean HR prediction. For shortwave RT, errors 545 on the order of 0.1 K day⁻¹ are generally considered acceptable - e.g., Table 2 of Krasnopolsky 546 et al. (2012), page 7 of Song and Roh (2021), Figure 1 of Kim and Song (2022). At most 547 heights the errors are on this order, except in the upper stratosphere, where MAE jumps 548 to 3.91 K day^{-1} and bias jumps to 2.71 K day^{-1} . Panel e shows the frequency of large 549 point errors for HR, which is below 5% throughout the troposphere but jumps to 45%550 in the upper stratosphere. Error maxima in the upper stratosphere are associated with 551 perturbed ozone layers; some examples will be shown in case studies. 552

Figures 7f-h show the quality of uncertainty estimates for $F_{\rm net}$. The spread-skill 553 plot (panel f) shows that $F_{\rm net}$ predictions are almost perfectly calibrated when spread 554 $\lesssim 40 \text{ W m}^{-2}$; for higher spread values, the model is slightly underconfident. The discard 555 test (panel g) shows that, despite the underconfidence at higher spread values, the model's 556 overall uncertainty is strongly correlated with its error for $F_{\rm net}$. In other words, one can 557 trust that lower uncertainty means lower expected $F_{\rm net}$ error. The PIT histogram (panel 558 h) shows that the model's $F_{\rm net}$ predictions are quite well calibrated, except slightly too 559 many PIT values below 0.5. In other words, $y_{\rm true}$ falls in the bottom half of the $\vec{y}_{\rm pred}$ 560 distribution more often than it should. Meanwhile, Figures 7i-l show the quality of un-561 certainty estimates for HR. Panel i shows the profile of CEFs, which are similar to large-562 error frequencies (panel e). In other words, most large errors are also catastrophic er-563 rors, because the confidence interval (CI) cannot account for errors $> 1 \text{ K day}^{-1}$. The 564 spread-skill plot (panel j) shows that the model is extremely overconfident, producing 565 only 14% as much spread as it should. The discard test (panel k) shows that, despite 566 this overconfidence, the model's overall uncertainty is strongly correlated with its HR 567 error. Finally, the PIT histogram (panel l) shows that the model's HR predictions are 568 poorly calibrated, with extreme PIT values (the first and last bars) occurring for 30%569 of data samples. In other words, y_{true} falls near the bottom or top of the \vec{y}_{pred} distribu-570 tion 30% of the time, three times as often as it should. 571

Having used the *moderately perturbed* validation data to select the best UQ method (BNN/CRPS), we now evaluate BNN/CRPS on the *heavily perturbed* testing data. By comparison of Figure 6 and Supplemental Figure S23, the overall ranking of UQ methods is similar between the validation and testing data. Most importantly, BNN/CRPS appears to be the best method for both datasets.



Figure 8: Detailed results of the BNN/CRPS method on testing data, for a model trained with clean data. Formatting is explained in the caption of Figure 7.

Figure 8 shows detailed results for the BNN/CRPS model on the testing data. Here we highlight differences from the validation results (Figure 7). The attributes diagrams (panels a-b) show that point predictions of F_{net} and HR are worse on the testing data (note the higher REL values), but these plots still indicate good skill except for the highest HR predictions. The MAE and bias profiles (panels c-d) show that point predictions of HR are still mostly acceptable in the troposphere, but problems in the stratosphere are worse in the testing data. Large-error frequencies for HR (panel e) are similar to the validation data but with a slightly better maximum (7% decrease) in the upper stratosphere. Results for $F_{\rm net}$ uncertainty (panels f-h) are similar to the validation data, except with a worse SSREL (55% increase) and slightly better PITD (12% decrease). Results for HR uncertainty (panels i-l) are also similar to the validation data – showing very poor skill – except with a slightly better maximum for CEF (7% decrease), worse SSREL (24% increase), and worse PITD (19% increase).

Key points: For models trained with clean data, even the best model produces unacceptable errors, as expected. The most notable errors are poor HR predictions in the stratosphere, with CEF > 40%; poor HR predictions at the highest values, with a bias of $\gg 1$ K day⁻¹; and poor HR uncertainty estimates throughout the atmosphere, with SSRAT < 14%. Results for the testing data are worse than for the validation data, as expected.

5.2 Case studies

596

Case study 1: Validation data. Figure 9 shows a case with the following per-597 turbations: a two-layer ice cloud (panel a), a multi-layer liquid cloud with large/noisy 598 LWC values (panel a), and an ozone layer with noisy mixing ratios (panel b). For the 599 HR spike due to ice cloud (around 10 km), all point predictions are too low and most 600 CIs (for all models except BNN/CRPS; panel f) miss the observation. For the HR spikes 601 due to liquid cloud (from 1-3 km), point predictions have a large error (> 1 K day⁻¹) but 602 observations generally fall within the CI, especially for the non-BNN models (panels c, 603 d, g). For the HR spike due to ozone (around 45 km), all point predictions and CIs are 604 far too low -i.e., all models produce a catastrophic error. The models also fail to cap-605 ture other aspects of ozone-related heating (from 15-60 km), including the HR minimum 606 around 60 km. 607

Key points: This case study exemplifies that while perturbed cloud sometimes causes larger point errors than perturbed ozone, perturbed ozone causes catastrophic errors more often. This conclusion is supported more rigorously by comparing panel i across Supplemental Figures S26-S28. The reason is that ozone varies much less in the training data than LWC/IWC – *e.g.*, all training samples have exactly one ozone layer with a maximum mixing ratio between 5.5 and 18.5 mg kg⁻¹, while different training samples have very different LWC/IWC profiles. Thus, perturbed ozone layers are more alien to the training data than perturbed cloud layers.



Figure 9: Case study for models trained with clean data and applied to a validation sample: 0000 UTC 9 Dec 2020, 1.58 °S, 94.80 °E. [a-b] Key predictor variables, *i.e.*, those subject to perturbation. [c] Actual HR profile (blue), along with ensemble-mean prediction (dashed red line) and 95% confidence interval (shaded red envelope), from the MME-only model. [d] Same but for MME/CRPS model. [e] Same but for BNN-only model. [f] Same but for BNN/CRPS model. [g] Same but for CRPS-only model. In the legends, "HR MAE" is the MAE of ensemble-mean HR predictions over the 127 heights; " $F_{\rm net}$ error" is the ensemble-mean $F_{\rm net}$ prediction minus actual; and " $F_{\rm down}^{\rm sfc}$ error" and " $F_{\rm up}^{\rm TOA}$ error" are defined analogously.

Case study 2: Cloudy testing data. Figure 10 shows a case with two follow-616 ing perturbations: a shallow ozone layer and near-surface moist layer (panel b). All mod-617 els struggle with ozone-related heating (from 15-30 km), as in the validation case but worse. 618 The near-surface moist layer (bottom 0.6 km) causes an HR maximum of ~ 10 K day⁻¹, 619 for which all UQ methods fail completely. The best models in this region are MME-only 620 (panel c) and MME/CRPS (panel d), but the HR maximum is still $\sim 3 \text{ K day}^{-1}$ above 621 both ensemble means and $\sim 1 \text{ K day}^{-1}$ above both CIs, so these errors are considered catas-622 trophic. 623

Key points: This case study exemplifies two conclusions from the broader dataset. First, although perturbed near-surface moisture generally causes smaller errors than perturbed ozone and cloud layers, near-surface moisture can still cause catastrophic errors. These are most common in profiles with little to no cloud, leaving ample solar radiation to reach the near-surface and interact with water vapour there. Second, perturbations in the testing data cause worse errors than perturbations in the validation data (*cf.* Figures 7 and 8), as expected.



Figure 10: Case study for models trained with clean data and applied to a testing sample: 1200 UTC 18 Dec 2020, 27.47 °N, 6.91 °E. All formatting is explained in the caption of Figure 9.

Case study 3: Cloud-free testing data. Figure 11 shows a case with the fol lowing perturbations: a shallow ozone layer with large mixing ratios (panel b), three cloud
 layers with large/noisy LWC values (panel a), a near-surface moist layer with specific
 humidity reaching ~30 g kg⁻¹ (panel b), and a near-surface warm layer with tempera ture reaching ~35 °C (panel b). There is virtually no heating near the surface, because
 all solar radiation is attenuated by the clouds above. The models capture this lack of heat-

ing quite well, except for MME/CRPS (panel d) and BNN/CRPS (panel f), which produce an erroneous HR maximum in the bottom 0.2 km. For the HR spike due to liquid
cloud (around 10 km), all models produce a catastrophic error. This error is worse than
cloud-related errors in the validation case (Figure 9), consistent with the more extreme
LWC values in this, a testing case. For ozone-related heating, all models produce catastrophic errors throughout the stratosphere.

Key points: The models were trained with clean data and simply have not seen
liquid cloud or ozone layers like the one here (cf. Figure 11b and Supplemental Figure
S9a). Thus, when presented with heavily perturbed testing data, which are far out of sample compared to the training data, the models (including their uncertainty estimates)
completely fail. This confirms our expectation from Section 4 and motivates Experiment

 $_{648}$ 2 – with lightly perturbed, instead of clean, training data.



Figure 11: Case study for models trained with clean data and applied to a testing sample: 1200 UTC 12 Aug 2020, 16.69 °S, 37.14 °E. All formatting is explained in the caption of Figure 9.

⁶⁴⁹ 6 Results for Experiment 2: Lightly perturbed training data

As in the discussion for Experiment 1, we start with overall diagnostics, then dig deeper with case studies.

652 6.1 Overall diagnostics

658

659

660

664

665

666

667

Figure 12 compares all five UQ methods on the validation data. This figure, which shows models trained with lightly perturbed data (henceforth LP-trained models), is analogous to Figure 6, which shows clean-trained models. We highlight several observations from Figure 12. Unless otherwise stated, these observations are true for the clean-trained models as well.

- 1. MME-only and BNN-only produce too little spread for all variables (panel e).
 - 2. Methods involving the CRPS produce too much spread for fluxes, but BNN/CRPS is the least overspread among these methods (panel e).
- 3. The clean-trained models produce far too little spread for HR (no SSRAT > 0.19; Figure 6e). However, among the LP-trained models, all except BNN-only produce an HR SSRAT > 0.74 (Figure 12e).
 - 4. The LP-trained models produce far fewer catastrophic errors than the clean-trained models, especially for HR (panel h). For example, the LP-trained model with the MME/CRPS method produces a CEF of 0.26% and 4.2% for HR and flux, respectively; analogous values for the clean-trained model are 6.7% and 7.6%.
- 5. As for the clean-trained models, BNN/CRPS produces the best uncertainty estimates overall (performing best on 8 of 10 uncertainty-based metrics). Unlike for the clean-trained models, there is no clear second-best method for uncertainty estimates.
- 672 6. BNN/CRPS also produces competitive point predictions (4th-best HR MAE, 2ndbest flux MAE, best HR REL, best flux REL).



Figure 12: Comparison of UQ methods on validation data, for models trained with lightly perturbed data. In panels a-d and f, lower is better; in panel e, closer to 1.0 is better; and in panel g, higher is better.

As for the clean-trained models, we judge that BNN/CRPS is the best UQ method overall. Supplemental Figure S29 shows that this conclusion also holds on the testing data. The remainder of this section focuses on the BNN/CRPS method and, for brevity, focuses on the testing data rather than the validation data. We have already seen for clean-trained models that performance deteriorates from the validation to the testing data, and this is true for LP-trained models as well.



Figure 13: Detailed results of the BNN/CRPS method on testing data, for a model trained with lightly perturbed data. Formatting is explained in the caption of Figure 7.

Figure 13 shows detailed testing results for the LP-trained BNN/CRPS model. We compare these results to the clean-trained BNN/CRPS model (Figure 8). The attributes diagrams (panels a-b) show that point predictions of both F_{net} and HR are extremely well calibrated, except for the highest HR predictions. The attributes diagrams are better for the LP-trained model. The MAE and bias profiles (panels c-d) show that point HR predictions are better for the LP-trained model, much better in the stratosphere. Specifically, the maximum MAE is 85% lower, and the maximum absolute bias is 92%

lower. The frequency of large HR errors (panel e) is also better for the LP-trained model, 687 with a 69% decrease in the upper-stratosphere maximum. Results for $F_{\rm net}$ uncertainty 688 (panels f-h) show that the LP-trained model is well calibrated, although slightly under-689 confident in general (panel f) and producing slightly too many PIT values below 0.5 (panel h). The SSREL and PITD are better than for the clean-trained model, but the SSRAT 691 is worse. Results for HR uncertainty (panels i-l) show that the LP-trained model is poorly 692 calibrated, with overconfidence at nearly all spread values (panel j) and too many PIT 693 values above 0.5 (panel l). However, these results are *much* better than for the clean-694 trained model, with a 69% decrease in maximum CEF, 89% decrease in SSREL, 285% 695 increase in SSRAT, and 38% decrease in PITD. Overall, the comparison of Figures 8 and 696 13 shows that the LP-trained model is better than the clean-trained model, but three 697 results of the LP-trained model are still concerning: large positive bias for HR point pre-698 dictions $\gtrsim 38$ K day⁻¹, a 14% frequency of catastrophic HR errors in the upper strato-699 sphere, and large overconfidence for HR in general. Supplemental Figure S30 – analo-700 gous to Figure 13 but for the validation data – shows that similar concerns exist in the 701 validation data but are much less severe. 702

Key points: Experiment 1 showed that training with clean data, which barely sample important basis vectors in the validation/testing data, leads to catastrophic errors.
 Experiment 2 shows that perturbing the training data *just a little* along these basis vectors leads to much better performance on the validation/testing data, even if the latter are still out of sample. However, catastrophic errors still occur, showing that ML-UQ is not magic.

⁷⁰⁹ 6.2 Case studies

710

Both case studies in this section are from the testing data.

Case study 1: Extreme liquid cloud and humidity. Figure 14 shows a case 711 with the following perturbations: a very dense liquid cloud (panel a), an ozone layer with 712 very large/noisy mixing ratios (panel b), and a near-surface moist layer with very large 713 humidity (panel b). Also, the maximum ozone content occurs at a lower height than usual, 714 around 20 km (in the lower stratosphere). For ozone-related heating around this level, 715 all models produce a catastrophic error. However, for ozone-related heating above this 716 level, most models (all except BNN-only; panel e) perform quite well. This result is in 717 stark contrast to case studies for the clean-trained models, which struggle with perturbed 718 ozone everywhere in the stratosphere. For the heating due to liquid cloud and the moist 719 layer (from 0-0.3 km), only the MME-only and MME/CRPS models (panels c-d) per-720 form well. The BNN-only and BNN/CRPS models (panels e-f) produce catastrophic er-721 rors, and the CRPS-only model (panel g) produces a large point error (~ 10 K day⁻¹) 722 for the HR maximum. 723

Key points: Although the LP-trained models are much better than the clean-trained models, every LP-trained model produces a catastrophic error somewhere in the profile.



Figure 14: Case study for LP-trained models applied to a testing sample: 1800 UTC 22 Apr 2020, 30.99 °N, 99.26 °W. All formatting is explained in the caption of Figure 9.

Case study 2: Extreme ice cloud and ozone. Figure 15 shows a case with the
following perturbations: a very dense two-layer ice cloud (panel a) and an ozone layer
with small/noisy mixing ratios (panel b). There is virtually no heating below 7 km, because all solar radiation is attenuated by the ice cloud above; all models capture this lack
of heating well. For the ice-cloud-related heating (around 8 km), all point predictions
are 1-3 K day⁻¹ too low. However, most CIs (for all models except BNN-only and BNN/CRPS;
panels e-f) capture the observed HR. For the ozone-related heating (above 10 km), the

MME/CRPS model (panel d) performs best. Specifically, at every height except the HR
 maximum around 29 km, the point error is < 1 K day⁻¹ and the CI captures the observation. The other models perform nearly as well, except the BNN-only model, which produces catastrophic errors from 30-45 km.

Key points: Again, the LP-trained models are much better than the clean-trained
 models. The MME/CRPS model even manages to produce no catastrophic errors for this
 case study, which is far out of sample.



Figure 15: Case study in the testing data: 0000 UTC 29 Jan 2020, 41.65 °N, 168.52 °E. All formatting is explained in the caption of Figure 9.
⁷⁴⁰ 7 Summary and future work

For a long time, uncertainty quantification (UQ) has been a key ambition for ma-741 chine learning (ML) in environmental science (ES). The field of computer science has re-742 cently made breakthroughs in ML-UQ, and environmental scientists are just beginning 743 to apply the resulting methods. One hope of the ES community is that new ML-UQ meth-744 ods can be used to assess the trustworthiness of an ML model on a case-by-case basis 745 e.g., to alert users when the model is expected to have a large error. However, ML-746 UQ methods, just like the base ML models with which they are coupled, do not gener-747 748 alize well to out-of-sample data. When a UQ-enhanced ML model encounters out-of-sample data, it is likely to produce a *catastrophic error* -i.e., an extremely wrong prediction 749 with high confidence. While scientists are generally aware that ML generalizes poorly 750 out of sample, in our experience they do not have this awareness for UQ, leaving them 751 prone to catastrophic errors. We wish to discourage overreliance on ML-UQ by show-752 ing that there are fundamentally unresolvable types of uncertainty, including that which 753 arises from out-of-sample data. 754

To this end, we trained neural networks (NN) to predict shortwave radiative trans-755 fer. The NNs predict 130 quantities – a length-127 vector of heating rates and 3 flux com-756 ponents – each with a 50-member ensemble. The ensemble is produced by one of five 757 UQ methods: a multi-model ensemble (MME), Bayesian neural network (BNN), train-758 ing with the continuous ranked probability score (CRPS) loss function, or a hybrid method 759 (MME/CRPS or BNN/CRPS). The validation and testing data are pushed out of sam-760 ple by perturbing several predictor variables: temperature, humidity, liquid cloud, ice 761 cloud, and ozone. 762

In Experiment 1, the NNs are trained with clean (unperturbed) data, then tasked 763 with generalizing to moderately perturbed validation data and heavily perturbed testing 764 data. Irrespective of the UQ method with which they are coupled, the NNs completely 765 fail on the validation and testing data, generating poor point predictions (ensemble means) 766 and uncertainty estimates. Even the best-performing UQ method (BNN/CRPS) is ex-767 tremely overconfident for heating rates, producing only 14% as much spread as it should. 768 Perturbations made to the ozone layer – which are more severe than perturbations made 769 to other atmospheric properties – confound our NNs the most. The models fail on val-770 idation and testing data because the perturbations therein are simply not "seen" in the 771 clean training data. In other words, the perturbations occur along basis vectors with lit-772 tle to no variability in the training data. While it is generally recognized that ML-based 773 predictions will fail in this setting, ML-generated uncertainty estimates fail just as spec-774 tacularly. This has serious implications for operational use of ML, e.g., in weather-forecasting. 775 If a high-impact event occurs in real time that is not represented in the training data 776 (*i.e.*, is out of sample), overreliance on ML – including ML-based uncertainty estimates 777 - could have severe consequences. 778

The discouraging results from Experiment 1 motivated another question: what hap-779 pens if the basis vectors represented by the perturbations are represented just a little in 780 the training data? To answer this, in Experiment 2 we trained NNs with *lightly perturbed* 781 data, calling these "LP-trained models" (as opposed to the clean-trained models in Ex-782 periment 1). On the validation and testing data, the LP-trained models performed much 783 better than clean-trained models. This result illustrates the power of triggering each ba-784 sis vector of variability – even just a little – in the training data. Ebert-Uphoff and Deng 785 (2017) found a similar result for causal discovery in ES: if a causal mechanism is not trig-786 gered in the training data, it will not be learned by the model. 787

Despite the obvious advantages of the LP-trained models, evaluation on the testing data revealed some concerning properties. For example, the best-performing UQ method (BNN/CRPS) is still quite overconfident for heating rates, producing only 52% as much spread as it should. Thus, lightly triggering important basis vectors in the training data

allowed the ML-UQ models to extrapolate much better along these basis vectors, but 792 it did not cure all ills. This is especially true for the testing data; the three concerns listed 793 above are quite minor in the validation data. The above results have two important im-794 plications for operational ML. First, while enhancing the training data by triggering im-795 portant basis vectors of the predictor space allows ML to better extrapolate along these 796 vectors, ML will likely struggle when extrapolating far out of sample. Second, in large 797 predictor spaces (which are common in ES), it is hard to know all the important basis 798 vectors, especially those representing high-impact events. Thus, even for ML models with 799 safeguards against poor out-of-sample performance – such as UQ or enhanced training 800 data – we still discourage overreliance on ML and ML-UQ. Also, we encourage users to 801 be familiar with the training data used for an ML model, so that they can identify out-802 of-sample (or poorly sampled) situations and approach the model with a healthy skep-803 ticism. 804

Future work will proceed along two lines. First, we will explore strategies for adapting ML-UQ to more realistic out-of-sample data, such as those caused by climate change (our application was a sandbox for testing the generalization of ML-UQ methods under extreme conditions). Second, we will try combining ML-UQ with tools that automatically detect out-of-sample data (Bulusu et al., 2020). Although these tools cannot improve an ML model's generalization to out-of-sample data, they can alert users when outof-sample data appear. This would automate part of the process of determining an ML model's trustworthiness.

Appendix A Aleatory vs. epistemic uncertainty

Uncertainty can be divided into two components: aleatory and epistemic. Briefly, according to the ML literature, aleatory uncertainty is due to gaps in the (training) dataset, while epistemic uncertainty is due to gaps in model development. Note, however, that the definitions vary across disciplines – see Figure A1 and discussions in Hüllermeier and Waegeman (2021), Bevan (2022), Haynes et al. (2023)). We use the ML definition throughout this manuscript, shown in Figure A1b. Figure A1 is adapted from Figure 3 of Haynes et al. (2023).



Figure A1: The aleatory/epistemic divide, according to different disciplines. [a] The math (original) definition is based only on the mathematical nature of the observed system. Only uncertainty due to stochastic processes, such as the chaotic nature of the atmosphere, is considered aleatory. [b] The ML definition of aleatory uncertainty is much wider, including not only uncertainty due to the stochastic nature of the system, but due to all other shortcomings of the dataset, such as limited observations.

In Section 1.3 we provide examples of unresolvable uncertainty, which cannot be captured by *any* ML-UQ method. One might wonder whether the unresolvable uncertainty in these scenarios is aleatory or epistemic. The answer is: it depends on how the dataset is chosen. We illustrate this below for Scenario 1 from Section 1.3, where uncertainty depends strongly on a variable $x_{unknown}$ not included in the ML model. The key question in distinguishing aleatory from epistemic is: where was the information lost? Let us track our steps:

- 1. The physical system contains both variables: x_{known} and x_{unknown} (Equations 1).
 - 2. The dataset collected by observing the physical system may or may not include
 - x_{unknown} . Letting M be the number of samples, the two possibilities for the dataset are

$$\mathcal{D}_1 = \left\{ (x_{\text{known}}^i, x_{\text{unknown}}^i); \quad i = 1, 2, \dots, M \right\} \text{ or }$$
$$\mathcal{D}_2 = \left\{ (x_{\text{known}}^i); \quad i = 1, 2, \dots, M \right\}.$$

3. Regardless of which dataset was chosen $(\mathcal{D}_1 \text{ or } \mathcal{D}_2)$, the ML model has no access to x_{unknown} and depends only on x_{known} .

In other words, regardless of which dataset was chosen, the ML model is exactly the same. 834 However, uncertainty in the model's output arising from its ignorance of x_{unknown} is con-835 sidered **epistemic** if the dataset chosen is \mathcal{D}_1 (because the problem is deemed to be in 836 the model), versus **aleatory** if the dataset chosen is \mathcal{D}_2 (because the problem is deemed 837 to be in the data). In other words, the distinction of aleatory vs. epistemic depends on 838 whether the relevant information was dropped during the data-collection or model-development 839 step. The key conclusion for this study is that the types of unresolvable uncertainty in 840 Section 1.3 can show up in both the aleatory and epistemic components; thus, we need 841 to employ UQ methods that can capture both types. 842

Appendix B Open research

We used version 3.0.0 of ML4RT (Machine Learning for Radiative Transfer; https:// doi.org/10.5281/zenodo.10086129) – a Python library managed by author Lagerquist – for all tasks in this study. The input data and all models not involved in a MME can be found at https://zenodo.org/doi/10.5281/zenodo.10081204; the MMEs can be found at https://zenodo.org/doi/10.5281/zenodo.10084393, https://zenodo.org/ doi/10.5281/zenodo.10084403, https://zenodo.org/doi/10.5281/zenodo.10084445, and https://zenodo.org/doi/10.5281/zenodo.10084454.

Acknowledgments

This work was partially supported by the NOAA Global Systems Laboratory, Cooper-

- ative Institute for Research in the Atmosphere, and NOAA Award Number NA19OAR4320073.
- Author Ebert-Uphoff's work was partially supported by NSF AI Institute grant #2019758.

References

864

829

830

831

832

833

- Baran, S., & Baran, A. (2021). Calibration of wind speed ensemble forecasts for
 power generation. arXiv e-prints, 2104 (14910). Retrieved from https://arxiv
 .org/abs/2104.14910
- Barnes, E., Barnes, R., & Gordillo, N. (2021). Adding uncertainty to neural network
 regression tasks in the geosciences. arXiv e-prints, 2109(07250). Retrieved
 from https://arxiv.org/abs/2109.07250
- Beucler, T., Pritchard, M., Yuval, J., Gupta, A., Peng, L., Rasp, S., ... Gentine, P. (2021). Climate-invariant machine learning. arXiv e-prints. 2112(08440).
 - P. (2021). Climate-invariant machine learning. arXiv e-prints, 2112(08440). Retrieved from https://arxiv.org/abs/2112.08440

865	Bevan, L. (2022). The ambiguities of uncertainty: A review of uncertainty frame-
866	works relevant to the assessment of environmental change. $Futures, 137$,
867	102919. Retrieved from https://doi.org/10.1016/j.futures.2022.102919
868	Bihlo, A. (2021). A generative adversarial network approach to (ensemble) weather
869	prediction. Neural Networks, 139, 1-16. Retrieved from https://doi.org/10
870	.1016/j.neunet.2021.02.003
871	Buiten, M. (2019). Towards intelligent regulation of artificial intelligence. European
872	Journal of Risk Regulation, $10(1)$, 41-59. Retrieved from https://doi.org/10
873	.1017/err.2019.8
874	Bulusu, S., Kailkhura, B., Li, B., Varshney, P., & Song, D. (2020). Anomalous exam-
875	ple detection in deep learning: A survey. <i>IEEE Access</i> , 8, 132330-132347. Re-
876	trieved from https://doi.org/10.1109/ACCESS.2020.3010274
877	Chapman, W., Monache, L. D., Alessandrini, S., Subramanian, A., Ralph, F., Xie,
878	S., Hayatbini, N. (2022). Probabilistic predictions from deterministic atmo-
879	spheric river forecasts with deep learning. Monthly Weather Review, $150(1)$,
880	215-234. Retrieved from https://doi.org/10.1175/MWR-D-21-0106.1
881	Clare, M., Jamil, O., & Morcrette, C. (2021). Combining distribution-based neu-
882	ral networks to predict weather forecast probabilities. Quarterly Journal of the
883	Royal Meteorological Society, 147(741), 4337-4357. Retrieved from https://
884	doi.org/10.1002/qj.4180
885	Cooney, J., Bowman, K., Homeyer, C., & Fenske, T. (2018). Ten year analy-
886	sis of tropopause-overshooting convection using GridRad data. Journal
887	of Geophysical Research: Atmospheres, 123(1), 329-343. Retrieved from
888	https://doi.org/10.1002/2017JD027718
889	Delle Monache, L., Eckel, F., Rife, D., Nagarajan, B., & Searight, K. (2013). Prob-
890	abilistic weather prediction with an analog ensemble. Monthly Weather
891	<i>Review</i> , 141(10), 3498-3516. Retrieved from https://doi.org/10.1175/
892	MWR-D-12-00281.1
893	Ebert-Uphoff, I., & Deng, Y. (2017). Causal discovery in the geosciences—Using
894	synthetic data to learn how to interpret results. Computers and Geosciences,
895	<i>99</i> , 50-60. Retrieved from https://doi.org/10.1016/j.cageo.2016.10
896	.008
897	Garg, S., Rasp, S., & Thuerey, N. (2022). WeatherBench Probability: A bench-
898	mark dataset for probabilistic medium-range weather forecasting along with
899	deep learning baseline models. $arXiv \ e-prints, \ 2205(00865)$. Retrieved from
900	https://arxiv.org/abs/2205.00865
901	Ghazvinian, M., Zhang, Y., Seo, D., He, M., & Fernando, N. (2021). A novel hybrid
902	artificial neural network-Parametric scheme for postprocessing medium-range
903	precipitation forecasts. Advances in Water Resources, 151, 103907. Retrieved
904	$\mathrm{from}\ \mathtt{https://doi.org/10.1016/j.advwatres.2021.103907}$
905	Gil, Y., Pierce, S., Babaie, H., Banerjee, A., Borne, K., Bust, G., Shekhar, S.
906	(2019). Intelligent systems for geosciences: An essential research agenda.
907	Communications of the Association for Computing Machinery, $62(1)$, 76-84.
908	Retrieved from https://dl.acm.org/doi/10.1145/3192335
909	Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press. Re-
910	${ m trieved} \ { m from} \ { m https://www.deeplearningbook.org}$
911	Hamill, T. (2001). Interpretation of rank histograms for verifying ensemble forecasts.
912	Monthly Weather Review, 129(3), 550-560. Retrieved from https://doi.org/
913	10.1175/1520-0493(2001)129%3C0550:IORHFV%3E2.0.CO;2
914	Haynes, K., Lagerquist, R., McGraw, M., Musgrave, K., & Ebert-Uphoff, I.
915	(2023). Creating and evaluating uncertainty estimates with neural net-
916	works for environmental-science applications. Artificial Intelligence for the
917	Earth Systems, 2(2), 1-29. Retrieved from https://doi.org/10.1175/
918	AIES-D-22-0061.1
919	Hertel, V., Chow, C., Wani, O., Wieland, M., & Martinis, S. (2023). Probabilistic

920	SAR-based water segmentation with adapted Bayesian convolutional neural
921	network. Remote Sensing of Environment, 285, 113388. Retrieved from
922	https://doi.org/10.1016/j.rse.2022.113388
923	Hoffman, M., Blei, D., Wang, C., & Paisley, J. (2013). Stochastic variational infer-
924	ence. Journal of Machine Learning Research, 14(1), 1303-1347. Retrieved from
925	https://www.jmlr.org/papers/volume14/hoffman13a/hoffman13a.pdf
926	Hsu, W. & Murphy, A. (1986). The attributes diagram: A geometrical frame-
920	work for assessing the quality of probability forecasts <i>International Journal of</i>
028	Forecasting 2(3) 285-293 Retrieved from https://doi.org/10.1016/0169
020	-2070(86)90048-8
020	Hüllermeier E & Waegeman W (2021) Aleatoric and epistemic uncertainty in
930	machine learning: An introduction to concepts and methods Machine Learn-
931	ina 110(3) 457-506 Retrieved from https://doi.org/10.1007/s10994-021
932	-05946-3
955	Jacono M. Dolamoro I. Mlawor F. Shonhard M. Clough S. & Collins W.
934	(2008) Radiative forcing by long lived grouphouse gases: Calculations with the
935	(2006). Radiative forcing by long-inved greenhouse gases. Calculations with the
936	112(D12) Detrieved from https://doi.org/10.1020/2008 D000044
937	Lenin L. Leni H. Devenid F. Dunting W. & Devenium M. (2022). Herde
938	Jospin, L., Laga, H., Boussaid, F., Buntine, W., & Bennamoun, M. (2022). Hands-
939	on Bayesian neural networks – A tutorial for deep learning users. <i>IEEE Com</i> -
940	putational Intelligence Magazine, 17(2), 29-48. Retrieved from https://dol
941	.org/10.1109/MC1.2022.3155327
942	Kim, P., & Song, H. (2022). Usefulness of automatic hyperparameter optimiza-
943	tion in developing radiation emulator in a numerical weather prediction $f(x) = f(x) = f(x)$
944	model. Atmosphere, 13(5), 721. Retrieved from https://doi.org/10.3390/
945	atmos13050721
946	Kingma, D., & Welling, M. (2013). Auto-encoding variational Bayes. arXiv e-prints,
947	1312(6114). Retrieved from https://arxiv.org/abs/1312.6114
948	Klotz, D., Kratzert, F., Gauch, M., Sampson, A., Brandstetter, J., Klambauer, G.,
949	Nearing, G. (2022). Uncertainty estimation with deep learning for rain-
950	fall-runoff modeling. Hydrology and Earth System Sciences, 26(6), 1673-1693.
951	Retrieved from https://doi.org/10.5194/hess-26-1673-2022
952	Krasnopolsky, V., Belochitski, A., Hou, Y., Lord, S., & Yang, F. (2012). Accu-
953	rate and fast neural network emulations of long and short wave radiation for
954	the NCEP Global Forecast System model (Vol. Office Note 471; Tech. Rep.).
955	Retrieved from https://repository.library.noaa.gov/view/noaa/6951
956	Lagerquist, R., Turner, D., Ebert-Uphoff, I., & Stewart, J. (2023). Estimat-
957	ing full longwave and shortwave radiative transfer with neural networks
958	of varying complexity. Journal of Atmospheric and Oceanic Technol-
959	ogy, conditionally accepted. Retrieved from https://doi.org/10.22541/
960	essoar.168319865.58439449/v1
961	Lagerquist, R., Turner, D., Ebert-Uphoff, I., Stewart, J., & Hagerty, V. (2021). Us-
962	ing deep learning to emulate and accelerate a radiative transfer model. Journal
963	of Atmospheric and Oceanic Technology, 38(10), 1673-1696. Retrieved from
964	https://doi.org/10.1175/JTECH-D-21-0007.1
965	Orescanin, M., Petković, V., Powell, S., Marsh, B., & Heslin, S. (2021). Bayesian
966	deep learning for passive microwave precipitation type detection. <i>IEEE</i>
967	Geoscience and Remote Sensing Letters, 19, 1-5. Retrieved from https://
968	doi.org/10.1109/LGRS.2021.3090743
969	Ortiz, P., Orescanin, M., Petković, V., Powell, S., & Marsh, B. (2022). Decomposing
970	satellite-based classification uncertainties in large earth science datasets. $\ensuremath{\mathit{IEEE}}$
971	Transactions on Geoscience and Remote Sensing, 60, 1-11. Retrieved from
972	https://doi.org/10.1109/TGRS.2022.3152516
973	Ranganath, R., Gerrish, S., & Blei, D. (2014). Black box variational inference. Pro-
	and in a of Machine Learning Research 22 814 822 Detriound from http://

975	proceedings.mlr.press/v33/ranganath14
976	Rasp, S., Pritchard, M., & Gentine, P. (2018). Deep learning to represent sub-
977	grid processes in climate models. Proceedings of the National Academy of Sci-
978	ences, 115(39), 9684-9689. Retrieved from https://doi.org/10.1073/pnas
979	.1810286115
980	Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais,
981	N., & Prabhat. (2019). Deep learning and process understanding for
982	data-driven Earth system science. <i>Nature</i> , 566, 195-204. Retrieved from
983	https://doi.org/10.1038/s41586-019-0912-1
984	Rezende, D., Mohamed, S., & Wierstra, D. (2014). Stochastic backpropagation and
985	variational inference in deep latent Gaussian models. International Conference
986	on Machine Learning, 2, 2. Retrieved from http://web2.cs.columbia.edu/
987	\sim blei/fogm/2018F/materials/RezendeMohamedWierstra2014.pdf
988	Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for
989	biomedical image segmentation. In International conference on medical image
990	computing and computer-assisted intervention. Munich, Germany. Retrieved
991	from https://doi.org/10.1007/978-3-319-24574-4_28
992	Scher, S., & Messori, G. (2021). Ensemble methods for neural network-based
993	weather forecasts. Journal of Advances in Modeling Earth Systems, $13(2)$.
994	Retrieved from https://doi.org/10.1029/2020MS002331
995	Scheuerer, M., Switanek, M., Worsnop, R., & Hamill, T. (2020). Using artificial
996	neural networks for generating probabilistic subseasonal precipitation forecasts
997	over California. Monthly Weather Review, $148(8)$, $3489-3506$. Retrieved from
998	https://doi.org/10.1175/MWR-D-20-0096.1
999	Schulz, B., & Lerch, S. (2022). Machine learning methods for postprocessing en-
1000	semble forecasts of wind gusts: A systematic comparison. Monthly Weather
1001	<i>Review</i> , 150(1), 235-257. Retrieved from https://doi.org/10.1175/
1002	MWR-D-21-0150.1
1003	Slovnik, W. (1992). International Meteorological Vocabulary (Tech. Rep.). World
1004	Meteorological Organization.
1005	Song, H., & Roh, S. (2021). Improved weather forecasting using neural network em-
1006	ulation for radiation parameterization. Journal of Advances in Modeling Earth
1007	Systems, 13(10). Retrieved from https://doi.org/10.1029/2021MS002609
1008	Van, T., Nguyen, T., Tran, N., Nguyen, H., Doan, L., Dao, H., & Minh, T.
1009	(2020). Interpreting the latent space of generative adversarial networks
1010	using supervised learning. International Conference on Advanced Comput-
1011	ing and Applications, 49-54. Retrieved from https://doi.org/10.1109/
1012	ACUMP50827.2020.00015
1013	veidkamp, S., Whan, K., Dirksen, S., & Schmeits, S. (2021). Statistical postpro-
1014	Cessing of which speed forecasts using convolutional neural networks. Montaly Weather Devices $1/0(4)$ 1141 1152 Detrieved from https://doi.org/
1015	Weather Review, 149(4), 1141-1152. Retrieved from https://doi.org/
1016	Wellage I & Hebba D (2006) Atmospheric Science: An Introductory Survey
1017	(Vol. 2) Elsevier
1018	(VOI. 2). Elsevier. Won V. Vicel D. De J. Tran D. & Crosse P. (2018). Elipout: Efficient recorde
1019	independent weight perturbations on mini batches International Conference
1020	an Learning Representations — Betrioved from https://orviv.org/pdf/1802
1021	0/1386 pdf
1022	Wimmers A Velden C & Cossuth I (2010) Using doop learning to estimate
1023	tropical cyclone intensity from satellite passive microwave imagery Monthly
1024	Weather Review $1/7(6)$ 2261-2282 Retrieved from https://doi.org/
1025	$10 \ 1175 / MWR - D - 18 - 0.391 \ 1$
1027	Zhou Z Siddiquee M Taibakhsh N & Liang I (2010) Unet $\pm\pm\cdot$ Redesigning
1028	skip connections to exploit multiscale features in image segmentation IEEE
1029	Transactions on Medical Imagina, 39(6), 1856-1867 Retrieved from https://

doi.org/10.1109/TMI.2019.2959609

Machine-learned uncertainty quantification is not magic: Lessons learned from emulating radiative transfer with ML Supplemental material

Ryan Lagerquist^{$1,2^*$}, Imme Ebert-Uphoff^{1,3}, David D. Turner², and Jebb Q. Stewart²

7	¹ Cooperative Institute for Research in the Atmosphere (CIRA), Colorado State University, Fort Collins,
8	Colorado ² National Oceanic and Atmospheric Administration (NOAA) Global Systems Laboratory (GSL) Boulder
9 10	Colorado
11	³ Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, Colorado

12 1 Methods for data perturbation

This section explains how each atmospheric property is perturbed, with the exception of near-surface temperature (see Section 2d of main text).

15 **1.1** Near-surface humidity

1

2

3

4

5

Our motivation is to mimic the lower-tropospheric moistening expected with climate change. The procedure has three parameters: maximum depth of the moist layer (D_{max}) , minimum surface relative humidity ($\mathrm{RH}_{\mathrm{sfc}}^{\min}$), and maximum ($\mathrm{RH}_{\mathrm{sfc}}^{\max}$). Parameter settings are shown in Table S1; the procedure is shown schematically in Figure S1. After the numbered procedure below, we recompute the two moisture variables used as predictors (relative and specific humidity), based on the new mixing ratio and untouched temperature/pressure.

23	1.	Sample to determine the depth of the moist layer: $D \in \mathcal{U}[0, D_{\max}]$.
24	2.	Sample to determine the surface RH: $RH_{sfc} \in \mathcal{U}\left[RH_{sfc}^{min}, RH_{sfc}^{max}\right]$.
25	3.	Compare the new and original (unperturbed) surface-RH values. If $RH_{sfc} \leq RH_{sfc}^{orig}$,
26		do nothing and end the procedure.
27	4.	Convert surface RH to surface mixing ratio, $w_{\rm sfc}$.
28	5.	Calculate the increase in surface mixing ratio: $\Delta w_{\rm sfc} = w_{\rm sfc} - w_{\rm sfc}^{\rm orig}$.
29	6.	At each height in the moist layer, scale the mixing-ratio increase linearly from $\Delta w_{\rm sfc}$
30		at the surface to 0 at height D above the surface. See Figures S1a-c.
31	7.	If step 6 led to any height with dewpoint > temperature, reduce dewpoint to tem-
32		perature. See Figure S1d.
33	8.	If step 6 led to any mixing ratio above 40 g kg ⁻¹ , reduce to 40 g kg ⁻¹ . See Figure
34		S1e.

^{*325} Broadway, R/GSL6, Boulder, CO 80305

Corresponding author: Ryan Lagerquist, ralager@colostate.edu



Figure S1: Procedure for perturbing near-surface humidity. In this example, the moist-layer depth D is 3 km and the new surface relative humidity RH_{sfc} is 100%. In panel c, the new mixing ratio is obtained by adding panels a and b; the new dewpoint is then computed from the new mixing ratio. In panel d, the new dewpoint is obtained by taking the minimum of dewpoint and temperature at each height; the new mixing ratio is obtained by reducing to 40 g kg⁻¹ at each height if necessary; the new dewpoint is then computed from the new mixing ratio.

Parameter	Setting for	Setting for	Setting for						
	lightly perturbed	validation data	testing data						
	training data								
Near-surface temperature									
D_{\max}	1.25 km	2.5 km	5 km						
$\Delta T_{ m sfc}^{ m max}$	2 K	4 K	8 K						
Near-surface humidity									
D_{\max}	1.25 km	2.5 km	5 km						
$\mathrm{RH}_{\mathrm{sfc}}^{\mathrm{min}}$	50%	50%	50%						
$\mathrm{RH}_{\mathrm{sfc}}^{\mathrm{max}}$	62.5%	75%	100%						
Liquid cloud									
N _{max}	2	3	5						
D_{\max}	5 km	5 km	5 km						
LWC_{center}^{max}	2 g m^{-3}	2.5 g m^{-3}	$5 \mathrm{g} \mathrm{m}^{-3}$						
$\sigma_{ m LWC}$	$0.25~\mathrm{g~m}^{\text{-3}}$	$0.5~{\rm g~m}^{-3}$	1 g m ⁻³						
Ice cloud									
$N_{ m max}$	2	3	5						
D_{\max}	5 km	5 km	5 km						
IWC_{center}^{max}	2 g m^{-3}	2.5 g m^{-3}	5 g m^{-3}						
$\sigma_{ m IWC}$	0.25 g m^{-3}	0.5 g m^{-3}	1 g m ⁻³						
Ozone									
D_{\min} and D_{\max}	40 and 60 km	$25~{\rm and}~60~{\rm km}$	0.1 and 60 km						
z_{center}^{\min} and z_{center}^{\max}	20 and 50 km AGL	$20~{\rm and}~50~{\rm km}~{\rm AGL}$	20 and 50 km AGL						
w_{center}^{max}	20 mg kg^{-1}	25 mg kg^{-1}	50 mg kg^{-1}						
$\sigma_{ m w}$	0.25 mg kg^{-1}	0.5 mg kg^{-1}	1 mg kg ⁻¹						

Table S1: Parameter settings for perturbation of predictor variables.

1.2 Liquid cloud

35

42

43

Our motivation is to create more complex cloud profiles, as well as denser and deeper clouds, than seen in the real atmosphere. The procedure has four parameters: maximum number of cloud layers (N_{max}) , maximum layer depth (D_{max}) , maximum liquid-water content at layer center (LWC^{max}_{center}), and noise level for LWC (σ_{LWC}). Parameter settings are shown in Table S1; the procedure is shown schematically in Figure S2.

1. Sample to determine the number of cloud layers: $N \in \mathcal{U}[0, N_{\max}]$.

2. For each cloud layer i from 1...N:

- (a) Sample to determine the depth of the i^{th} cloud: $D \in \mathcal{U}[0, D_{\text{max}}]$.
- (b) Sample to determine the height of the cloud top above the surface:

$z_{\text{top}} \in \mathcal{U}[0, z_{\text{tropopause}} + 2 \text{ km}]^{2,3}$

45

46

47

51

52

53

54

55

56

57

58

59

(c) Calculate the height of the cloud bottom: $z_{\text{bottom}} = z_{\text{top}} - D$. If this leads to $z_{\text{bottom}} < 0 \text{ km}$ AGL, increase to 0 km AGL. See Figure S2b.

(d) Find all grid points in the cloud. These are grid points with a height in $[z_{\text{bottom}}, z_{\text{top}}]$ and temperature ≥ -40 °C that have not already been assigned to another liquidcloud layer.⁴ See Figure S2c.

(e) Sample to determine the LWC at the center of the cloud: $LWC_{center} \in \mathcal{U}[0, LWC_{center}^{max}]$.

(f) At each height in the cloud, scale the LWC linearly from LWC_{center} at the center to 0 g m⁻³ at both the top and bottom. See Figure S2d.

(g) Add Gaussian noise to the LWC profile for this cloud. Specifically, at each height in the cloud, add an offset δ , sampled from a normal distribution with mean = 0 g m⁻³ and standard deviation = σ_{LWC} . Symbolically, $\delta \in \mathcal{N}(0, \sigma_{LWC})$. See Figure S2e.

(h) If the previous step led to any LWC < 0 g m⁻³, increase to 0 g m⁻³. If the previous step led to any LWC > LWC_{center}, reduce to LWC_{center}. See Figure S2f.

² We use the World Meteorological Organization (Slovnik, 1992) definition of the first trop opause: the lowest height at which lapse rate decreases to < 2 K km⁻¹ (let this be z'), provided that the mean lapse rate between z' and z' + 2 km does not exceed 2 K km⁻¹.

³ A maximum height of $z_{\text{tropopause}} + 2$ km allows clouds to reach 2 km into the stratosphere, as in the overshooting tops of thunderstorms. Figure 12 of Cooney et al. (2018) shows that few overshooting tops reach further than 2 km into the stratosphere.

 $^{^4}$ Supercooled liquid droplets can exist at temperatures down to ~-40 °C; see Figure 6.29 of Wallace and Hobbs (2006).



Figure S2: Procedure for creating a new liquid cloud layer. [a] Original profiles of LWC and temperature. Temperature is not perturbed in this procedure, but it is shown as a reference variable, because liquid droplets cannot exist at temperatures below -40 °C. [b] Same as panel a, but the extent of the proposed new cloud is shaded in grey. In this example, the proposed new cloud has depth D = 5 km. [c] Same as panel b, but corrected to exclude temperatures < -40°C and overlap with other liquid cloud. [d] LWC profile after adding new cloud. In this example, the LWC at the center of the cloud is LWC_{center} = 5 g m⁻³. [e] Same as panel d, but after adding Gaussian noise for new cloud. In this example, the noise parameter is $\sigma_{LWC} = 1$ g m⁻³. [f] Same as panel e, but after removing unwanted values created by Gaussian noise.

60 1.3 Ice cloud

⁶¹ This procedure has the same parameters as for liquid cloud, but replacing liquid-⁶² water content with ice-water content (IWC). In other words, replace LWC^{max}_{center} with IWC^{max}_{center} ⁶³ and σ_{LWC} with σ_{IWC} . Parameter settings are shown in Table S1.

The procedure itself – shown schematically in Figure S3 – is the same as for liquid cloud, except in step 2d. The criterion "temperature \geq -40 °C" is replaced with "temperature < 0 °C".



Figure S3: Procedure for creating a new ice cloud layer. [a] Original profiles of IWC and temperature. [b] Same as panel a, but the extent of the proposed new cloud is shaded in grey. In this example, the proposed new cloud has depth D = 5 km. [c] Same as panel b, but corrected to exclude temperatures $\geq 0^{\circ}$ C and overlap with other ice cloud. [d] IWC profile after adding new cloud. In this example, the IWC at the center of the cloud is

IWC_{center} = 5 g m⁻³. [e] Same as panel d, but after adding Gaussian noise for new cloud. In this example, the noise parameter is $\sigma_{IWC} = 1$ g m⁻³. [f] Same as panel e, but after removing unwanted values created by Gaussian noise.

1.4 Ozone

67

75

76

77

78

79

80

81

82

⁶⁸ Our motivation is to create more complex ozone layers – over a wider range of lo-⁶⁹ cations, depths, and mixing ratios – than seen in the real atmosphere. The procedure ⁷⁰ has six parameters: minimum and maximum depth (D_{\min} and D_{\max}), minimum and max-⁷¹ imum height of layer center (z_{center}^{\min} and z_{center}^{\max}), maximum ozone mixing ratio at layer ⁷² center (w_{center}^{\max}), and noise level for mixing ratio (σ_w). Parameter settings are shown in ⁷³ Table S1; the procedure is shown schematically in Figure S4.

- 1. Sample to determine the ozone-layer depth: $D \in \mathcal{U}[D_{\min}, D_{\max}]$.
 - 2. Sample to determine the height of the layer center: $z_{\text{center}} \in \mathcal{U}\left[z_{\text{center}}^{\min}, z_{\text{center}}^{\max}\right]$.
 - 3. Find all grid points in the ozone layer. These are grid points with a height in
 - $\left[z_{\text{center}} \frac{1}{2}D, z_{\text{center}} + \frac{1}{2}D\right]$ that are above the tropopause.
 - 4. Sample to determine the ozone mixing ratio at the layer center: $w_{\text{center}} \in \mathcal{U}[0, w_{\text{center}}^{\max}]$.
 - 5. At each height in the ozone layer, scale the mixing ratio linearly from w_{center} at the center to 0 mg kg⁻¹ at both the top and bottom.
 - 6. Add Gaussian noise to the mixing-ratio profile. Specifically, at each height in the ozone layer, add an offset δ sampled from $\mathcal{N}(0, \sigma_w)$.
- 7. If the previous step led to any $w < 0 \text{ mg kg}^{-1}$, increase to 0 mg kg⁻¹. If the previous step led to any $w > w_{\text{center}}$, reduce to w_{center} .



Figure S4: Procedure for creating a new ozone layer. [a] Original temperature profile, with proposed extent of ozone layer shaded in grey. In this example, the proposed ozone layer has depth D = 20 km and center $z_{center} = 25$ km AGL. The tropopause is at 16.5 km AGL, where temperature begins to increase with height. [b] Same as panel a, but corrected to exclude heights below the tropopause. [c] New ozone profile. In this example, $w_{center} = 10 \text{ mg kg}^{-1}$. [d] Same as panel c, but after adding Gaussian noise. In this example, the noise parameter is $\sigma_w = 1 \text{ mg kg}^{-1}$. [e] Same as panel d, but after removing unwanted values created by Gaussian noise.

⁸⁵ 2 Effects of data perturbation

Figures S5-S9 show the effects of different levels of data perturbation: light (for one set of training data), moderate (for the validation data), and heavy (for the testing data). Specific perturbation methods are discussed in Section 3d of the main text and Supplemental Section 1. A key property shown in these figures is that the perturbations to ozone are more drastic than those to liquid and ice water, which in turn are much more drastic than the perturbations to temperature and humidity.

Note that, as in Table S1, temperature is perturbed only at heights below {1.25, 2.5, 5} km AGL in the {lightly perturbed training, validation, testing} data. Thus, above 5 km, the temperature distribution is nearly identical across the four datasets (Figure S5). All differences above 5 km are caused by differences among the *clean* datasets, *i.e.*, before perturbation. The same is true for other quantities not affected by perturbation: specific humidity above 5 km, liquid-water content in the stratosphere, ice-water content in the stratosphere, and ozone mixing ratio in the troposphere.



Figure S5: Distribution of temperature in [a] the clean training data, [b] the lightly perturbed training data, [c] the validation data, and [d] the testing data.



Figure S6: Same as Figure S5 but for specific humidity.



Figure S7: Same as Figure S5 but for liquid-water content (LWC).



Figure S8: Same as Figure S5 but for ice-water content (IWC).



Figure S9: Same as Figure S5 but for ozone mixing ratio.

⁹⁹ 3 Hyperparameter experiment to determine the best BNN architec ¹⁰⁰ ture

We conduct four hyperparameter experiments, each to determine the best BNN (Bayesian neural network) architecture in a given context. The contexts are:

• for the BNN-only UQ method, trained with clean data;

103

104

105

106

- for the BNN/CRPS method, trained with clean data;
 - for the BNN-only method, trained with lightly perturbed data;
 - for the BNN/CRPS method, trained with lightly perturbed data.

For each context we optimize four hyperparameters: the number of Bayesian fully connected layers (N_b^{fully}) , number of Bayesian upsampling connections $(N_b^{\text{upsampling}})$, number of Bayesian skip connections (N_b^{skip}) , and training method for Bayesian layers (reparameterization or flipout). Bayesian fully connected layers allow the network to do UQ for scalar outputs (fluxes), while Bayesian upsampling and skip connections allow the network to do UQ for vector outputs (heating rates). The attempted values for each hyperparameter are listed in Table S2. Table S2: Hyperparameters optimized for BNNs.

HyperparameterValues attemptedNumber of Bayesian fully connected layers2, 3, 4Number of Bayesian upsampling connections1, 2Number of Bayesian skip connections1, 2Training method for Bayesian layersReparameterization, flipoutSpectral complexity64, 128

We experiment with the first four hyperparameters in particular – defined in Sec-114 tion 3a of the main text – because they are the main choices involved in changing a U-115 net from deterministic to Bayesian. When making fully connected layers Bayesian, we 116 start at the layer nearest to the scalar outputs and work backwards. For example, if N_b^{fully} 117 = 3, we make the two output layers (labeled "A" in Figure S10) and the preceding layer 118 ("B" in Figure S10) Bayesian. Similarly, when making upsampling and skip connections 119 Bayesian, we start at the layer nearest to the vector outputs and work backwards. For example, if $N_b^{\text{upsampling}} = 2$ and $N_b^{\text{skip}} = 1$, the two upsampling connections with high-120 121 est spatial resolution ("C" in Figure S10) and one skip connection with highest spatial 122 resolution ("D" in Figure S10) are made Bayesian. The vector output layer ("E" in Fig-123 ure S10) is always Bayesian if the network is Bayesian. Within a network, the training 124 method (reparameterization or flipout) is the same for all Bayesian layers. 125



Figure S10: The optimal U-net++ setup for one context: the BNN-only UQ method with lightly perturbed training data. Since BNN-only uses the deterministic loss function rather than the CRPS, N (the ensemble size) = 1. This figure, which shows one specific U-net++ setup, is analogous to Figure 4 in the main text, which shows the generic U-net++ setup for UQ. In Figure 4 the double arrows indicate a component that *might be* Bayesian, while in this figure the double arrows indicate a component that *is* Bayesian.

For all hyperparameters not related to UQ, we use the optimal value determined in Lagerquist et al. (2023, henceforth L23), with one exception. The exception is spectral complexity, defined as the number of filters in the first convolutional layer (128 in Figure S10, the optimal value determined by L23). With the BNN/CRPS UQ method, this spectral complexity makes the network so large that training leads to out-of-memory errors. Thus, we attempt values of 64 and 128, as L23 found that a spectral complexity of 64 is nearly optimal.

The hyperparameter experiment is a grid search (Section 11.4.3 of Goodfellow et al., 2016), meaning that we try all 48 possible combinations of the hyperparameters. We use the validation data to find the best model, *i.e.*, best combination of hyperparame-

ter values. Because we wish to optimize both point predictions (ensemble means) and 136 uncertainty estimates for two variables (heating rates [HR] and fluxes), choosing the "best" 137 model is non-trivial, as model performance can be described by a wide variety of error 138 metrics. We examine the 12 metrics shown in Figure 5 of the main text: {MAE, REL, 139 SSREL, SSRAT, PITD, and MF} for {HR, flux}. MAE is mean absolute error; REL is 140 reliability; SSREL is spread-skill reliability; SSRAT is spread-skill ratio; PITD is devi-141 ation from the perfect probability integral transform (PIT) histogram; and MF is mono-142 tonicity fraction measured from the discard test. MAE and REL concern point predic-143 tions only, while the other metrics concern the entire predicted distribution, including 144 uncertainty. MAE, REL, SSREL, and PITD are negatively oriented with a perfect value 145 of 0.0; MF is positively oriented with a perfect value of 1.0; and SSRAT has a perfect 146 value of 1.0, with values of $[0,\infty)$ possible. We choose the best model by subjectively 147 combining results across the 12 metrics. 148

For brevity – and because results are similar across the four contexts – here we 149 show the results for only one context, namely the BNN-only method trained with lightly 150 perturbed data. These results are shown in Figures S11-S22. In each figure, the circle 151 represents the selected model (based on all 12 metrics) and the star represents the model 152 with the best value of the given metric. Note that the best values for HR MF (Figure 153 S16) and flux MF (Figure S22) are a multi-way tie, as many models achieve a perfect 154 MF of 1.0. In this case, the tie is broken arbitrarily and the star marks one of the many 155 models with perfect MF. The selected model has $N_b^{\text{fully}} = 3$, $N_b^{\text{upsampling}} = 2$, $N_b^{\text{skip}} =$ 156 1, spectral complexity of 128, and uses the reparameterization method instead of flipout. 157 The first four of these hyperparameter values are represented in Figure S10. 158



Figure S11: MAE for heating rate, computed on validation data for each set of hyperparameters. "Dense" here is a synonym for "fully connected". Each panel shows one spectral complexity and one Bayesian training method (reparameterization or flipout); within each panel the other three hyperparameters vary. The white circle marks the selected model, and the white star (hidden behind the white circle) marks the model with the lowest value for this error metric.



Figure S12: REL for heating rate, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S13: SSREL for heating rate, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S14: SSRAT for heating rate, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S15: PITD for heating rate, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S16: MF for heating rate, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S17: MAE for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S18: REL for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S19: SSREL for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S20: SSRAT for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S21: PITD for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.



Figure S22: MF for flux variables, computed on validation data for each set of hyperparameters. Formatting is explained in the caption of Figure S11.

¹⁵⁹ 4 Further results of Experiment 1 (clean training data)

Figure S23 shows error metrics, based on the *heavily perturbed* testing data, for all five UQ methods. This is analogous to Figure 5 in the main text, which shows results on the *moderately perturbed* validation data. The main purpose of Figure S23 is to show that the overall ranking of UQ methods is similar between the validation and testing data. Most importantly, BNN/CRPS appears to be the best method for both datasets.



Figure S23: Comparison of UQ methods on testing data, for models trained with clean (unperturbed) data. In panel g, higher is better; in panel e, closer to 1.0 is better; in all other panels, lower is better. "CEF" in panel h is catastrophic-error frequency.

Figure S24 shows detailed results for the BNN/CRPS model on the subset of the testing data with perturbed near-surface temperature. This is analogous to Figure 7 in the main text, which shows results on the entire testing set. Figures S25-S28 are analogous to S24 but for different perturbation types: near-surface moisture, liquid cloud, ice cloud, and ozone. Some broad conclusions from these figures are highlighted in the main text.



Figure S24: Detailed results of the BNN/CRPS method, for a model trained with clean data, on the subset of the testing data with perturbed near-surface temperature. Formatting is explained in the caption of Figure 6 in the main text.



Figure S25: Same as Figure S24 but for the subset of testing data with perturbed near-surface moisture.



Figure S26: Same as Figure S24 but for the subset of testing data with perturbed liquid cloud.



Figure S27: Same as Figure S24 but for the subset of testing data with perturbed ice cloud.



Figure S28: Same as Figure S24 but for the subset of testing data with perturbed ozone.

5 Further results of Experiment 2 (lightly perturbed training data)

Figure S29 shows error metrics, based on the *heavily perturbed* testing data, for all five UQ methods. This is analogous to Figure 11 in the main text, which shows results on the *moderately perturbed* validation data. The main purpose of this new figure is to show that the BNN/CRPS method appears to be best for both datasets.



Figure S29: Comparison of UQ methods on testing data, for models trained with lightly perturbed data. In panel g, higher is better; in panel e, closer to 1.0 is better; in all other panels, lower is better. "CEF" in panel h is catastrophic-error frequency.

Figure S30 shows detailed results for the BNN/CRPS model on the validation data. This is analogous to Figure 12 in the main text, which shows results on the testing set.

- ¹⁷⁸ The main purpose of this new figure is to show that concerning properties of the test-
- ¹⁷⁹ ing results are much less concerning in the validation results. Specifically, the positive
- bias for large HR (when ensemble mean $\gtrsim 38$ K day⁻¹) decreases from ~ 5 K day⁻¹ to ~ 1
- $_{181}$ K day⁻¹; struggles with perturbed ozone improve (catastrophic-error frequency in the up-
- $_{182}$ per stratosphere decreases from 14% to 5%); and overall underestimation of HR uncer-
- tainty improves, with SSRAT increasing from 0.520 to 0.884.


Figure S30: Detailed results of the BNN/CRPS method, for a model trained with lightly perturbed data, on the validation data. Formatting is explained in the caption of Figure 6 in the main text.

Figure S31 shows detailed results for the BNN/CRPS model on the subset of the testing data with perturbed near-surface temperature. This is analogous to Figure 12 in the main text, which shows results on the entire testing set. Figures S32-S35 are analogous to S31 but for different perturbation types: near-surface moisture, liquid cloud, ice cloud, and ozone. Some broad conclusions from these figures are highlighted in the main text.



Figure S31: Detailed results of the BNN/CRPS method, for a model trained with lightly perturbed data, on the subset of the testing data with perturbed near-surface temperature. Formatting is explained in the caption of Figure 6 in the main text.



Figure S32: Same as Figure S31 but for the subset of testing data with perturbed near-surface moisture.



Figure S33: Same as Figure S31 but for the subset of testing data with perturbed liquid cloud.



Figure S34: Same as Figure S31 but for the subset of testing data with perturbed ice cloud.



Figure S35: Same as Figure S31 but for the subset of testing data with perturbed ozone.

190 References

- 191Cooney, J., Bowman, K., Homeyer, C., & Fenske, T.(2018).Ten year analy-192sis of tropopause-overshooting convection using GridRad data.Journal193of Geophysical Research: Atmospheres, 123(1), 329-343.Retrieved from194https://doi.org/10.1002/2017JD027718
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press. Re trieved from https://www.deeplearningbook.org
- Lagerquist, R., Turner, D., Ebert-Uphoff, I., & Stewart, J. (2023). Estimat ing full longwave and shortwave radiative transfer with neural networks
- 199of varying complexity.Journal of Atmospheric and Oceanic Technol-200ogy, conditionally accepted.Retrieved from https://doi.org/10.22541/201essoar.168319865.58439449/v1
- Slovnik, W. (1992). International Meteorological Vocabulary (Tech. Rep.). World
 Meteorological Organization.
- Wallace, J., & Hobbs, P. (2006). Atmospheric Science: An Introductory Survey
 (Vol. 2). Elsevier.