# Uncertainty Quantification of a Machine Learning Subgrid-Scale Parameterization for Atmospheric Gravity Waves

Laura A Mansfield<sup>1</sup> and Aditi Sheshadri<sup>1</sup>

<sup>1</sup>Stanford University

February 28, 2024

#### Abstract

Subgrid-scale processes, such as atmospheric gravity waves, play a pivotal role in shaping the Earth's climate but cannot be explicitly resolved in climate models due to limitations on resolution. Instead, subgrid-scale parameterizations are used to capture their effects. Recently, machine learning has emerged as a promising approach to learn parameterizations. In this study, we explore uncertainties associated with a machine learning parameterization for atmospheric gravity waves. Focusing on the uncertainties in the training process (parametric uncertainty), we use an ensemble of neural networks to emulate an existing gravity wave parameterization. We estimate both offline uncertainties in raw neural network output and online uncertainties in climate model output, after the neural networks are coupled. We find that online parametric uncertainty contributes a significant source of uncertainty in climate model output that must be considered when introducing neural network parameterizations. This uncertainty quantification provides valuable insights into the reliability and robustness of machine learning-based gravity wave parameterizations, thus advancing our understanding of their potential applications in climate modeling.

#### Hosted file

AGU\_SuppInfo\_UncertaintyQuantification.docx available at https://authorea.com/users/ 709595/articles/717814-uncertainty-quantification-of-a-machine-learning-subgrid-scaleparameterization-for-atmospheric-gravity-waves

#### Hosted file

mov02.gif available at https://authorea.com/users/709595/articles/717814-uncertaintyquantification-of-a-machine-learning-subgrid-scale-parameterization-for-atmosphericgravity-waves

#### Hosted file

mov03.gif available at https://authorea.com/users/709595/articles/717814-uncertaintyquantification-of-a-machine-learning-subgrid-scale-parameterization-for-atmosphericgravity-waves

#### Hosted file

mov01.gif available at https://authorea.com/users/709595/articles/717814-uncertaintyquantification-of-a-machine-learning-subgrid-scale-parameterization-for-atmosphericgravity-waves



# AGU Word Manuscript Template

1 2 3 4	Uncertainty Quantification of a Machine Learning Subgrid-Scale Parameterization for Atmospheric Gravity Waves
5	L. A. Mansfield <sup>1</sup> , and A. Sheshadri <sup>1</sup>
6 7	<sup>1</sup> Earth System Science, Doerr School of Sustainability, Stanford University, California
8	Corresponding author: Laura A. Mansfield (lauraman@stanford.edu)
9	Key Points:
10 11	• Using ensembles of neural networks, we learn parametric uncertainties associated with an emulator of a gravity wave parameterization.
12 13	• When coupled to the climate model, the ensemble of neural networks reveals increased climate variability.
14 15 16	• Parametric uncertainty dominates the Quasi-Biennial Oscillation statistics, although polar vortex properties remain robust to parameters.

### 17 Abstract

- 18 Subgrid-scale processes, such as atmospheric gravity waves, play a pivotal role in shaping the
- 19 Earth's climate but cannot be explicitly resolved in climate models due to limitations on
- 20 resolution. Instead, subgrid-scale parameterizations are used to capture their effects. Recently,
- 21 machine learning has emerged as a promising approach to learn parameterizations. In this study,
- 22 we explore uncertainties associated with a machine learning parameterization for atmospheric
- 23 gravity waves. Focusing on the uncertainties in the training process (parametric uncertainty), we
- 24 use an ensemble of neural networks to emulate an existing gravity wave parameterization. We
- estimate both offline uncertainties in raw neural network output and online uncertainties in
- climate model output, after the neural networks are coupled. We find that online parametric
   uncertainty contributes a significant source of uncertainty in climate model output that must be
- considered when introducing neural network parameterizations. This uncertainty quantification
- provides valuable insights into the reliability and robustness of machine learning-based gravity
- 30 wave parameterizations, thus advancing our understanding of their potential applications in
- 31 climate modeling.
- 32

### 33 Plain Language Summary

- 34 Climate models are unable to resolve processes that vary on length and time scales smaller than
- the model resolution and timestep. For example, atmospheric gravity waves, which are waves
- 36 created when winds encounter disturbances to the flow, such as mountains, convection and
- 37 fronts, can have wavelengths smaller than the spacing between grid cells. Climate models use
- <sup>38</sup> "parameterizations" to capture the effect of these processes. Machine learning based
- 39 parameterizations are becoming popular because they can learn relationships purely from data.
- However, we do not have a good understanding of the uncertainties introduced through machine
   learning parameterizations. This study estimates uncertainties associated with training a neural
- learning parameterizations. This study estimates uncertainties associated with training a neural
   network gravity wave parameterization. We explore uncertainties in the neural network output,
- network gravity wave parameterization. We explore uncertainties in the neural network output,
   as well as the uncertainties in the climate model output, when the neural network is used for the
- as well as the uncertainties in the climate model output, when the neural network is used for thgravity wave parameterization.
- 45 **1 Introduction**
- 46

## 1.1. Subgrid-scale parameterizations

Global climate models (GCMs) simulate the entire Earth system by coupling a dynamical 47 core, which numerically solves the primitive equations for atmospheric flow, with other physical 48 components called "subgrid-scale parameterizations". The latter includes dynamical processes 49 occurring on scales smaller than the grid-scale (generally O(100 km) for a typical GCM; Chen 50 et al., 2021), such as convection and short wavelength gravity waves, and non-dynamical 51 processes, such as radiation, atmospheric chemistry, and cloud and aerosol microphysics. 52 Subgrid-scale parameterizations make up a large portion of the computational cost associated 53 with GCM simulations and sometimes make drastic assumptions for the sake of computational 54 cost, which can introduce additional sources of model uncertainty. This has motivated the 55 demand for faster and/or higher accuracy schemes that use machine learning (ML)/artificial 56 intelligence (AI), which hold out the potential for training on large volumes of training data and 57

58 performing fast inferences when invoked.

### 59

ML-based subgrid-scale parameterizations have demonstrated skill across a wide range of 60 atmospheric processes including convection, clouds, aerosols, radiation and gravity waves (e.g., 61 Brenowitz et al., 2020; Brenowitz & Bretherton, 2019; Chantry et al., 2021; Chevallier et al., 62 2000; Espinosa et al., 2022; Gentine et al., 2018; Harder et al., 2022; Krasnopolsky & Fox-63 Rabinovitz, 2006; O'Gorman & Dwyer, 2018; Perkins et al., 2023; Rasp et al., 2018; Ukkonen, 64 2022; Yu et al., 2023; Yuval et al., 2021; Yuval & O'Gorman, 2020). However, few studies have 65 explored the uncertainties associated with these. Stochastic subgrid-scale parameterizations have 66 been developed by sampling from parametric distributions, learned through neural networks 67 (Guillaumin & Zanna, 2021) and generative adversarial networks (GANs) (Gagne II et al., 2020; 68 Nadiga et al., 2022; Perezhogin et al., 2023). These studies focus on stochastic representations to 69 improve model accuracy since they may better represent scaling properties (Palmer, 2019). 70 Including uncertainty estimates can also be beneficial in assessing the trustworthiness of model 71 predictions (Haynes et al., 2023; McGovern et al., 2022), and has gained some attention in 72 weather and climate prediction studies (e.g., Delaunay & Christensen, 2022; Gagne et al., 2014, 73 2017; Gordon & Barnes, 2022; Weyn et al., 2021). Here, we explore uncertainty quantification 74 in a machine learning subgrid-scale parameterization (a type of *model uncertainty*; Hawkins & 75 Sutton, 2009; Palmer, 2019), focusing on gravity wave parameterizations. 76

77 78 79

### **1.2 Atmospheric Gravity Waves**

Atmospheric gravity waves (GWs) are important drivers of middle atmosphere 80 circulation as they transport momentum upwards and away from their sources in the lower 81 troposphere (Fritts & Alexander, 2003). They are forced by perturbations to a stable stratified 82 flow, for instance, orography, convection, and frontogenesis. They propagate primarily in the 83 84 vertical and, due to the decreasing density in the upper atmosphere, grow in amplitude until reaching a critical level, at which point they break and deposit momentum. This provides a 85 forcing on the mean flow in the middle and upper atmosphere and has a substantial impact on 86 atmospheric circulation, including in driving the Quasi-Biennial Oscillation (QBO) in the 87 equatorial stratosphere (Baldwin et al., 2001) and affecting the occurrence of Sudden 88 Stratospheric Warmings in the polar vortex during winter (Wang & Alexander, 2009), described 89 90 further in Section 1.3.

91

GW wavelengths can range from O(1 km) to O(1000 km), which presents a challenge 92 for accurate representation in global climate models (GCMs). While the primitive equations do 93 capture GW dynamics, typical GCM resolutions are O(100 km), resulting in a large portion of 94 the GW spectrum being un- or under-resolved. Parameterizations must be employed to model the 95 impacts of subgrid-scale GWs on the mean flow and are critical for obtaining realistic 96 97 circulation, for example, to induce a spontaneous QBO (Bushell et al., 2020). Some studies find GW parameterizations to be necessary even in kilometer-scale resolution simulations (Achatz et 98 al., 2023; Polichtchouk et al., 2023), suggesting that the need for accurate parameterizations will 99 persist even as modeling centers move towards high resolution GCMs (or "digital twins"; e.g., 100 Bauer et al., 2021). 101 102

- 103 **1.2.2 Gravity wave parameterizations**
- 104

GCMs usually make use of both an orographic and a non-orographic GW parameterization to capture their effects. Machine learning alternatives to GW parameterizations have recently gained attention in several forms. Chantry et al. (2021), Espinosa et al. (2022) and Hardiman et al. (2023) present machine learning emulators of existing non-orographic gravity wave schemes, while Dong et al. (2023) and Sun et al. (2023) use machine learning to learn gravity wave momentum fluxes from high resolution simulations.

111

This study can be viewed as a continuation of the work by Espinosa et al. (2022), which 112 develops an emulator of a non-orographic GW parameterization designed primarily for 113 convectively forced GWs (Alexander & Dunkerton, 1999). Note that this machine learning 114 parameterization is, at best, as accurate as the scheme it aims to emulate and is not significantly 115 faster than the original physics-based scheme, which could be due to coupling of the neural 116 network within a Fortran-based GCM (Cambridge-ICCS, 2023). Rather, this neural network 117 emulator is used as a first step towards probing uncertainties introduced when replacing a gravity 118 wave parameterization with an emulator, when we have a "ground truth" parameterization for 119 reference. 120

121 122

1.3 Gravity wave effects

123 124

## 1.3.1 Quasi-Biennial Oscillation

Gravity waves strongly influence the stratospheric circulation. In the tropical stratosphere, the dominant mode of variability is the Quasi-Biennial Oscillation (QBO), in which the equatorial stratospheric zonal winds alternate between easterly and westerly and descend downwards with time (Gray, 2010). The change in direction is driven by breaking waves across a range of scales (Baldwin et al., 2001; Lindzen & Holton, 1968), with modeling studies suggesting that non-orographic gravity wave parameterizations contribute around half of the forcing required for a simulated QBO (Holt et al., 2020).

133

In this study, we measure the performance of gravity wave parameterizations through the simulated QBO period and amplitudes at 10 hPa, where the QBO amplitude is generally a maximum (Bushell et al., 2020; Richter et al., 2020). We consider the QBO winds to be defined by the zonal mean zonal winds between  $5^{\circ}S$  and  $5^{\circ}N$ . Following Schenzinger et al., (2017), we estimate the period of a QBO cycle by the length between transition times from westward and eastward flow, after applying a 5-month binomial filter to remove high frequency variability. The amplitude is estimated as the absolute maximum of the QBO winds during each cycle.

141 142

## 1.3.2 Stratospheric Polar Vortex

As well as driving the equatorial stratospheric circulation, gravity waves are also influential at high latitudes. Gravity waves affect the stratospheric polar vortex in both hemispheres, as they contribute to the breakdown of the polar vortices, influencing the frequency and properties of Sudden Stratospheric Warmings (SSWs) (Siskind et al., 2007, 2010; Wang & Alexander, 2009; Whiteway et al., 1997; Wright et al., 2010) and the timing of the Spring final warming (Gupta et al., 2021). SSWs are defined as a reversal of the zonal mean zonal winds at 60°N at 10 hPa (Butler et al., 2015) which is followed by large and rapid temperature increases

(>30-40 K) in the polar stratosphere. They occur around 6 times per decade in the Northern 151

152 hemisphere, but are not common in the Southern hemisphere. In this study, we consider gravity

wave parameterization effects on the number of Northern hemisphere SSWs per decade and the 153

timing of the final warming of the Southern hemisphere polar vortex. 154

155 156

158

#### 2. Uncertainty Quantification 157

Uncertainties can be categorized into two types: *aleatoric uncertainty* and *epistemic* 159 uncertainty (Hüllermeier & Waegeman, 2021). Aleatoric uncertainty is used to describe the 160 variability in a system that is due to inherently random effects (Haynes et al., 2023; Hüllermeier 161 & Waegeman, 2021). It represents the statistical or stochastic nature of a system, such as flipping 162 a coin or rolling a dice and in ML literature, refers to uncertainty in the data. It includes *internal* 163 variability of the system and observational uncertainties in the data. In contrast, epistemic 164 uncertainty is caused by a lack of knowledge about the best model for a system and refers to 165 uncertainty in the model. It includes structural uncertainties from the choice of ML architecture, 166 parametric uncertainties in estimating of model parameters, and out-of-sample uncertainties 167 which arise when predicting outside of the range of the training data. 168

In this study, we aim to quantify parametric uncertainty, a type of epistemic uncertainty, 169 in an ML-based parameterization for gravity waves. We expect this to also capture out-of-sample 170 uncertainties, i.e., increased uncertainty when generalizing to a situation that lies outside of the 171 training data distribution. For simplicity, we do not estimate aleatoric uncertainty in the training 172 data, and we also do not consider structural uncertainty. Future studies may wish to account for 173 these additional types of uncertainty for a more complete picture. There are several methods that 174 could be used to estimate parametric uncertainty (Abdar et al., 2021). Here, we use an ensemble 175 of deep neural networks or "deep ensembles", which involves training multiple identical neural 176 networks, each with a different initialization (Lakshminarayanan et al., 2017). Each neural 177 network converges upon slightly different parameters which are then used to predict an 178 ensemble, from which statistics can be obtained. This is a relatively simple approach to 179 implement, although can be costly as it requires repetition during training and evaluation. Deep 180 ensembles have been used in climate model applications for prediction (Weyn et al., 2021), but 181 have not been used for subgrid-scale parameterizations. In this context, deep ensembles could be 182 viewed as a machine learning complement to "perturbed parameter ensembles" (PPE), which 183 involve perturbing physics-based parameters for uncertainty quantification (e.g., Murphy et al., 184 2007; Sengupta et al., 2021; Sexton et al., 2021). 185

#### 186 187 **3** Methods

- 188
- 189
- 190

### **3.1 Gravity Wave Parameterization Setup**

Alexander & Dunkerton (1999; hereafter AD99) present a simple non-orographic, gravity 191 wave parameterization that has been used in various GCMs, including GFDL's Atmospheric 192 Model 3 (Donner et al., 2011), Isca (Vallis et al., 2018), and MiMA (Jucker & Gerber, 2017). 193 AD99 estimates gravity wave drag (GWD) in both the zonal and meridional directions for each 194

195 level in a column, at each grid-cell and timestep. When coupled into a climate model, gravity

wave drag or forcing acts to accelerate or decelerate winds (i.e., it is a wind tendency). As a 196

spectral parameterization, AD99 defines a spectrum of gravity waves at a source level with momentum flux distributed by phase speeds, assumed to follow a Gaussian distribution centered at 0 m/s with half-width 35 m/s. This spectrum of gravity waves propagates upwards until the waves reach the critical level (when the wind speed equals the phase speed of the waves), when breaking occurs and drag is deposited.

202 203

204

## 3.2 Atmospheric Model Setup

We use an intermediate complexity GCM, a Model of an idealized Moist Atmosphere 205 (MiMA) (Jucker & Gerber, 2017). It is run at spectral resolution T42, corresponding to 64 206 latitudes by 128 longitudes (approximately 2.8 degrees or 300 km grid spacing at the equator), 207 with 40 model levels. The level top is 0.18 hPa, with a strong dissipating sponge layer in the 208 upper three levels (0.85-0.18 hPa). AD99 is coupled into MiMA with the parameters described 209 above and with a fixed source level defined to be 315 hPa in the tropics and decreasing in height 210 with latitude, roughly in line with the tropopause. The model is run with an advection timestep of 211 10 minutes and a physics timestep, which includes calling the gravity wave parameterization, of 212 213 3 hours.

214 215

### 3.2 Machine Learning Setup

216 We use the neural network (NN) gravity wave parameterization developed by Espinosa et 217 al. (2022). This is trained on MiMA simulations using the AD99 gravity wave parameterization, 218 described above (Alexander & Dunkerton, 1999). Espinosa et al. (2022) show that the NN 219 emulator, trained on one year of data, achieves an accurate representation of the AD99 scheme 220 both offline and online. For the online tests, Espinosa et al. (2022) replace the original AD99 221 scheme in MiMA with the NN emulator within MiMA and show that these coupled NN 222 simulations produce a Quasi-Biennial Oscillation consistent with original AD99 simulation. 223 Furthermore, when tested on an out-of-sample climate under 4xCO2 forcing, the NN simulations 224 remained stable and reproduced similar changes to the QBO as the AD99 simulations. 225 226

Espinosa et al. (2022) emulate the zonal and meridional GW drag with two independently 227 trained but almost identical fully connected NNs. The inputs to the zonal GW drag network are 228 229 zonal winds at all levels, u, temperature at all levels, T, surface pressure,  $p_s$ , and latitude,  $\lambda$ , and similarly for the meridional GW drag the inputs are meridional winds at all levels,  $v, T, p_s$ , and 230  $\lambda$ . MiMA uses 40 pressure levels, giving a total of 82 inputs into the NN. The architecture 231 consists of four shared hidden layers followed by another four pressure level specific layers (see 232 233 Supporting Information of Espinosa et al., 2022). The network outputs the zonal/meridional GW drag for all 40 pressure levels. Note that the pressure levels closest the surface always predict 234 zero, where there is no GW drag below the source of the GWs. Although these layers are 235 redundant, we include them because the AD99 gravity wave source level changes with latitude to 236 237 follow the approximate level of the tropopause. Following Espinosa et al. (2022), we normalize the input and output data to have a zero mean and standard deviation of 1. For the pressure levels 238 below the source level, where all GW drag values are exactly zero and standard deviation is 239 undefined, we fix the outputs to zero. Although we follow the same architecture as Espinosa et 240 al. (2022), there are some software differences in our implementation. Firstly, we opt for 241 PyTorch (Paszke et al., 2019) rather than Keras and TensorFlow (Abadi et al., 2015; Chollet & 242

others, 2015) for the machine learning library. Secondly, Espinosa et al. (2022) use the forpy

software (Rabel, 2019) to call python code in the fortran-based climate model. This resulted in a

slow-down of roughly 2.5x when replacing AD99 with the NN emulator. Instead, we use FTorch

- (Cambridge-ICCS, 2023), a software package that directly calls the existing Torch C++ interface
   from Fortran resulting in faster inference. We find a 20% slow-down in the NN simulations
- from Fortran resulting in faster inference. We find a 20% slow-down in the NN simulations relative to the AD99 simulations, although we have not explored if this could be optimized
- 248 relative to the AD99 simulations, although we have not explored if this could 249 further.
- 250

In this study, we capture parametric uncertainty of the NN emulator presented in 251 Espinosa et al. (2022) using deep ensembles (Lakshminarayanan et al., 2017). We repeatedly 252 train an ensemble of size 30 independent NNs, each with the same architecture and trained on the 253 same data but with different random seed initializations. The random seed affects the 254 initialization of the NN parameters and the shuffling order of data during training, leading to 255 slightly different parameters when converged. Following Espinosa et al. (2022), we train the 256 NNs with one year of data, selected so that it contains a typical QBO cycle with a period and 257 amplitude similar to the long-term mean period and amplitude. We use the following one year of 258 data for the validation dataset, and the following 20 years are used for the test dataset, requiring 259 22 years of simulation data in total. Figure 1 shows (a) the QBO zonal winds and (b) the QBO 260 261



262

Figure 1 The QBO (a) zonal winds and (b) zonal gravity wave drag for the training, validation,
and test dataset.

265 266 267 **4 Results** 268

269 **4.1 Offline predictions** 

270 271 Figure 2 shows an example of gravity wave drag (GWD) profiles for a single grid cell close to the equator for a) the zonal component and b) the meridional component, with the black 272 273 line indicating the ground truth from the AD99 parameterization and the red line indicating the mean prediction across all NN ensemble members. The orange shading represents 1 standard 274 deviation across all ensemble members. Animations showing the evolution of this GWD profile 275 can be found in the Supporting Materials. The NNs agree well on the gravity wave profiles and 276 the ground truth falls within the 1 standard deviation range for across most model levels for the 277 zonal component. The meridional component generally captures the patterns within the profile 278 279 but is found to be less accurate, even when considering the uncertainty estimates.

- 280
- 281





283 284

285

Figure 2 Example profiles of a) zonal and b) meridional gravity wave drag at one grid-cell and
one timestep in the tropics where the black line indicates the ground truth from the AD99
parameterization, the red line indicates the mean prediction across all neural network ensembles
and the orange shading indicates 1 standard deviation across these ensembles.

290 To measure the errors, we calculate the continuous ranked probability score (CRPS), a generalization of mean absolute error that allows for comparison of probability distributions. The 291 use of CRPS to measure error between a predicted probability distribution and a single ground 292 293 truth has long been used for verification of ensemble weather forecasts (Hersbach, 2000), and 294 has recently been adopted for probabilistic machine learning (Gneiting & Raftery, 2007). Figure 3 shows CRPS for a) zonal and b) meridional gravity wave drag predictions over a range of 295 latitudes. Note the scale of the axis is reduced by 10x relative to the gravity wave drag 296 297 magnitudes in Figure 2. We find lower errors in the lower and mid-stratosphere that increase with height, where gravity wave drag magnitudes also increase. We see good performance across 298 all latitudes. 299



300 Zonal GWD (ms<sup>-2</sup>) Le<sup>-6</sup> Meridional GWD (ms<sup>-2</sup>) Le<sup>-6</sup> 301 Figure 3 Continuous Ranked Probability Score for a) zonal and b) meridional gravity wave drag 302 for different latitudes over the test dataset.

303 304

305

### 4.2 Offline uncertainty estimates

306 One common problem in uncertainty quantification of deep learning algorithms is in ensuring that uncertainty estimates are reasonable, often known as calibration of uncertainty 307 308 (Lakshminarayanan et al., 2017). A well-calibrated machine learning model should predict low 309 uncertainties when errors are small and high uncertainties when errors are large (for instance, 310 when the data is out-of-sample). Figure 4 shows the 1 standard deviation uncertainty estimates against the ensemble mean absolute errors estimated for the test dataset, with the colors 311 312 representing the density of points. Ideally, these should be correlated and lie approximately along the y = x line shown in the dashed line. Points above the y = x line are underconfident and 313 points below are overconfident. Although the errors and predicted uncertainties are correlated, 314 we see that the NNs suffer from overconfidence and frequently underestimate the uncertainty 315 relative to the error. This is typical behavior for machine learning uncertainty estimates, 316 including those based on deep ensembles (Abdar et al., 2021), and may be not be surprising 317 given we only consider one type of uncertainty (parametric uncertainty) and do not consider 318 structural uncertainty or data uncertainty in these estimates. This overconfidence is systematic 319 across all levels of the stratosphere and occurs for both zonal and meridional NNs, but especially 320 for the meridional predictions. 321 322

- 323
- 324

#### manuscript submitted to JAMES



Confidence of neural networks at latitudes 5°S-5°N at 10.9 hPa

Figure 4 Ensemble uncertainty (measured as 1 standard deviation amongst the ensemble predictions) against ensemble error (measured as the mean absolute error across all ensemble predictions) for a) zonal and b) meridional gravity wave drag for test dataset between 5°S-5°N at 10 hPa. Each individual point represents a single prediction at one timestep and grid-cell and they are shaded according to density. The black dashed line shows the y = x line.

331

325

332 333

341

### 4.3 Offline and Online Probability Distributions

Once coupled online into MiMA, the ensembles begin to diverge from each other even though they are initialized from the same state. This is partly due to the chaotic nature of the atmosphere where minute differences in one atmospheric variable can lead to very different atmospheric states after some time. Even introducing relatively minor differences in the GWD profiles, such as those in Figure 2, can lead to very different atmospheric states. Here, we aim to quantify how uncertainties in Figure 2 propagate into the GCM. We examine long-term statistics in order to separate out the NN parametric uncertainty from the internal variability.

342 We consider GWD in the tropics, due to its influence on the QBO. Figure 5 shows distributions of gravity wave drag in the upper stratosphere at 10 hPa for (a) zonal and (b) 343 meridional components, where the black line indicates ground truth from the AD99 MiMA 344 simulations, the blue line indicates the offline NN predicted gravity wave drag and the red line 345 indicates the online NN predicted GWD. Both offline and online distributions are centered over 346 347 the same location as AD99, indicating that the NN does not introduce a bias. In the lower stratosphere, the distributions are virtually indistinguishable (not shown). However, in the upper 348 stratosphere at 10 hPa, the NN distributions take a different shape than AD99. This is 349 particularly notable around the low negative zonal gravity wave drag values, where AD99 350 predicts an asymmetric gravity wave drag distribution with a positive skew. The NN 351 distributions are more symmetric between positive and negative values. This may because 352 353 machine learning optimizes for RMSE which may overly smooth gravity wave drag profiles, reducing asymmetry between positive and negative drag. The online NN distributions are 354 slightly smoother than the offline NN distributions. We suggest that this must be caused by the 355 interaction between the predicted gravity wave drag and the winds when coupled online. This is 356

verified by Figure 6a, which shows distributions of zonal winds near the equator at 10 hPa,
 where online distributions tend to be smoother and weaker than the AD99 distributions.

360 Figure 5 b shows that the online and offline meridional distributions are highly similar, 361 even though they are smoothed out at low magnitudes. Even though the meridional NN is 362 generally less accurate (e.g., Figure 2b), the meridional component of gravity wave drag does not 363 appear to diverge when coupled online. Similarly, Figure 6 b shows the distribution of the 364 meridional winds to be unchanged when the NNs are coupled. This indicates that the meridional 365 circulation is not highly sensitive to the effects of subgrid-scale gravity wave drag, possibly due 366 to lower magnitude of the meridional winds.

- 367
- 368

Distributions of Gravity Wave Drag for Equator at 10.9 hPa



374 5°N.

#### a) Meridional Zonal b) 0.08 AD99 AD99 Online Online 0.07 0.06 0.05 0.04 0.03 0.02 0.01 0.00 -40 40 -60 -40 -20 0 20 40 60 -60 -20 0 20 60 Zonal wind (ms<sup>-1</sup>) Meridional wind (ms<sup>-1</sup>)

Distributions of Wind for Equator at 10.9 hPa

375

Figure 6 a) zonal and b) meridional wind distributions for AD99 (black) and online NN

377 simulations (red) at 10 hPa between 5°S-5°N.

### **4.4 QBO uncertainties**

380 Ultimately, we are interested in how the NN estimations for GWD influence the 381 climatology and its variability when coupled into a GCM. We examine statistics of the OBO in 382 MiMA by calculating the QBO period and amplitude at 10 hPa for each QBO cycle within 400 383 years of AD99 simulations and the 600 years of NN simulations (from 30 simulations each of 20 384 385 year simulations), shown in Figure 7. While the mean period of the QBO across all simulation years are similar, the NN ensembles show increased variability that can be attributed to the 386 387 parametric uncertainty. The NNs also appear to introduce a bias that reduces the QBO amplitude, consistent with the reduction in QBO zonal winds (Figure 6). These increases in QBO variability 388 originate from differences between NN ensemble members (and therefore from the learned NN 389 parameters), each of which tend to maintain fairly consistent OBO periods and amplitudes within 390 the 20 year simulation. 391



378



393

394

Figure 7 Violin plots showing distributions of QBO a) period and b) amplitude for the AD99
simulations in grey and for NN simulations in orange. The boxplots also show the median, upper
and lower quartiles and each point represents a single OBO cycle.

398



Equation 1

401 ensemble of NNs, the additional variability from the uncertainty in parameters,  $\sigma_{param}$ , can be calculated as 402  $\sigma_{param}^2 = \sigma_{AD99}^2 + \sigma_{NNs}^2$ 

403

404

where  $\sigma_{AD99}^2$  is the variance in the AD99 simulations and  $\sigma_{NNs}^2$  is the total variance across 405 all NN ensemble members. These results are shown in Table 1. Notably, the parametric 406 uncertainty is significantly larger than the internal variability in the AD99 simulations, for both 407 the QBO period and amplitude. It is possible that these uncertainties are underestimates of the 408 409 true parametric uncertainty, given the overconfidence noted in offline tests (Figure 4). Still, the 410 uncertainties in NN parameters are much greater than uncertainties in the parameters in the physics-based scheme AD99, estimated to be 1.53 months and 2.14 m/s for the period and 411 amplitude respectively, in Mansfield & Sheshadri (2022) under the same model set-up. This 412 413 highlights the importance of uncertainty quantification, regardless of whether the parameterization is physics-based or machine learning based. 414

415

Table 1 Mean and variability of QBO calculated across MiMA simulations using AD99 vs. the 416

ensemble of NNs. Means are estimated across all QBO cycles in a 400 year long MiMA 417

simulation using AD99 and in 600 years of simulations from the 30-member, 20 year long 418

419 simulations from the ensemble of NNs. Variability is measured as 1 standard deviation between

all OBO cycles. Parametric uncertainty is calculated assuming OBO cycles are normally 420

distributed (Equation 1). 421

	Mean		Variability (mea	asured as 1 stand	ard deviation)
	AD99	Ensemble of NNs	Internal variability in AD99 simulations	Total variability in ensemble of NNs	Parametric uncertainty
Period (months)	25.32	26.78	2.03	3.82	3.25
Amplitude (m/s)	28.29	25.91	2.17	3.86	3.18

422

423 424

425 426

### 4.5 Polar vortex uncertainties

427 The QBO is just one phenomenon that is strongly influenced by gravity wave dynamics. 428 The stratospheric polar vortices in both hemispheres also depend upon gravity wave activity. In 429 particular, the breakdown of the polar vortices during sudden stratospheric warmings (SSWs) 430 and in the springtime final warming is driven by both planetary-scale and subgrid-scale gravity 431 waves, and the variability of these events could also be impacted by changes to the gravity wave 432

### manuscript submitted to JAMES

- 433 parameterization. For the northern hemisphere polar vortex, we consider the frequency of SSWs
- and for the southern hemisphere, we consider polar vortex lifetime. Figure 8 shows there is no
- obvious distinction between the variability of these properties between the AD99 and NN
   simulations, thus making the attribution of extratropical changes (and therefore, the calibration of
- extratropical parameters in AD99 and other schemes; Mansfield & Sheshadri, 2022) rather
- 438 challenging. This may be because the breakdown of the polar vortices is driven by both
- 439 planetary-scale waves and subgrid-scale gravity waves, thereby reducing the impact of any
- 440 changes to the parameterization. Furthermore, some studies find there may be a compensation
- 441 effect between resolved Rossby waves and unresolved gravity waves during SSW events (e.g.,
- 442 Cohen & Gerber, 2013), while some studies suggest that small scale gravity waves influence
- 443 polar vortex recovery after a SSW more strongly than the breakdown itself (Wicker et al., 2023).
- 444



445

446 Figure 8 Histograms showing a) the Northern hemisphere number of SSWs per decade and b)

the Southern hemisphere polar vortex lifetime for AD99 simulations in grey and the NN

448 *simulations in orange.* 

449

450

### 451 **5** Conclusions

452

This study uses deep neural network ensembles to quantify parametric uncertainties in a machine learning parameterization of gravity wave drag. We use the neural network architecture of Espinosa et al. (2022) trained on one year of data simulated by the intermediate complexity GCM, MiMA, which uses AD99 gravity wave parameterization (Alexander & Dunkerton, 1999; Jucker & Gerber, 2017). An ensemble of 30 identical neural networks are trained, each initialized with a different random seed. This ensemble allows us to estimate parametric uncertainties in neural network weights and biases. First, we assessed uncertainties in raw GWD output, which we refer to as *offline uncertainties*. We find fairly consistent results across all neural networks. Then, we used the FTorch library to couple the neural network into MiMA, allowing for GCM simulations that use the machine learning parameterization in place of the traditional physics-based scheme (Cambridge-ICCS, 2023). We assess uncertainties in GCM

- 464 output for gravity wave drag and wind, refering to these as *online uncertainties*. We find 465 increased online uncertainty, particularly for zonal winds.
- 466

487

Comparing long-term statistics of the climate within MiMA using the physics-based 467 scheme AD99 and the ensemble of neural networks, showed that the use of NN emulators can 468 alter the circulation significantly. We found that the NNs from the ensemble produce a bias in 469 the QBO towards reduced amplitudes and dramatically increase the variability of the QBO, with 470 uncertainty from NN parameters increasing the variability between OBO cycles by over 50%. 471 Uncertainty quantification of parameterizations should therefore not be overlooked when 472 developing ML-based schemes for future climate models. Our findings reiterate results from 473 previous studies that find that, even when offline tests indicate "good" NN performance with 474 relatively low uncertainties, the coupling of machine learning schemes into climate models can 475 still introduce a significant source of uncertainty (Brenowitz et al., 2020; Lin et al., 2023). 476 Learning distributions on the model parameters could provide a basis for further parameter 477 refinement, for example, acting as a Bayesian prior distribution that could be constrained through 478 online calibration, such as derivative-free optimization Ensemble Kalman methods (Pahlavan et 479 al., 2023). As with traditional parameterization calibration, this could lead to improved QBO 480 statistics and reduced parametric uncertainty. Interestingly, we find that the behavior and 481 breakdown of the polar vortex is not strongly dependent on the parameterization, which may be 482 partially due to influences from planetary-scale waves. This suggests that it may not be possible 483 to further calibrate neural network parameters to polar vortex properties, and is comparable to 484 the difficulties in calibration of extratropical parameters of AD99 (Mansfield & Sheshadri, 485 2022). 486

We only scratch the surface of uncertainty quantification for machine learning 488 parameterizations. Firstly, we describe only one type of uncertainty: parametric uncertainty, a 489 type of epistemic (model) uncertainty. There exist a wide range of machine learning approaches 490 that could be used for this task, including Bayesian Neural Networks, Monte Carlo dropout 491 generative models and deep ensembles (Abdar et al., 2021). We used deep ensemble methods for 492 this task (Lakshminarayanan et al., 2017), due to their simplicity to implement. However, this 493 approach is computationally costly during both training and evaluation, requiring the use of 494 ensembles which is not feasible for long climate model integrations. A more complete picture 495 would be given by also assessing aleatoric (data) uncertainties. We note that our parametric 496 uncertainty estimates would change given a different training dataset, which makes detangling 497 the effects of epistemic and aleatoric uncertainty a challenge (Haynes et al., 2023; Hüllermeier & 498 Waegeman, 2021). Still, learning the relative contributions between model and data uncertainties 499 would be insightful when designing machine learning parameterizations. Aleatoric uncertainties 500 501 could be estimated through the use of Bayesian neural networks or Monte Carlo dropout (Abdar et al., 2021), by parameterizing gravity wave outputs as a distribution (Guillaumin & Zanna, 502

2021; Haynes et al., 2023), or through generative models such as GANs (Gagne II et al., 2020;
Nadiga et al., 2022; Perezhogin et al., 2023).

505

Secondly, the machine learning parameterization used here is an emulator of an existing 506 scheme, allowing us to compare against a ground truth simulation. Future studies may wish to 507 extend this to train ML models on gravity-wave resolving simulations e.g., with kilometer-scale 508 resolution models such as IFS (Anantharaj et al., 2022), WRF (Sun et al., 2023) or ICON 509 (Hohenegger et al., 2023). When using novel training datasets from high resolution simulations, 510 we do not have online "true" distributions to compare against, which could present challenges 511 when disentangling the various sources of variability. Furthermore, it also raises the issue of 512 understanding the role of aleatoric uncertainty, e.g., in the choice of training data and method for 513 estimating gravity wave drag (Sun et al., 2023). 514

516 Thirdly, MiMA is an intermediate complexity atmospheric circulation model. One may expect that coupling this atmospheric model to other Earth system components, such as the 517 ocean, land, and sea-ice, would introduce further uncertainties. Therefore, we might consider the 518 results presented here as a lower bound on the uncertainties we could expect to see in fully 519 operational Earth system models that employ ML parameterizations. Extending this study to 520 higher complexity Earth system models would be significantly more costly, however, this could 521 be worthwhile towards better informing the design of ML parameterizations, which ultimately 522 could lead to efficient but accurate hybrid GCMs that combine traditional dynamical solvers with 523 novel machine learning parameterizations. 524

525 526

515

### 527 Acknowledgments

- 528 This research was made possible by Schmidt Sciences, a philanthropic initiative founded by Eric
- and Wendy Schmidt, as part of the Virtual Earth System Research Institute (VESRI). AS
- acknowledges support from the National Science Foundation through grant OAC-2004492.
- 531 We would also like to thank our Datawave colleagues, in particular L. Minah Yang and Dave
- 532 Connelly for their work on the PyTorch implementation of the machine learning model, and
- 533 Simon Clifford, Jack Atkinson, Dominic Orchard and others at ICCS, for their help in setting up
- the FTorch coupler with the Fortran-based climate model. We also appreciate the Stanford high
- 535 performance computing resources that made this work possible.
- 536
- 537 **Open Research**

### manuscript submitted to JAMES

538	The code to run simulations, train neural networks and replicate plots presented in this paper is
539	available at https://github.com/lm2612/WaveNet_UQ. The data generated will be made available
540	on the Stanford Digital Repository on publication. The FTorch library for coupling PyTorch to
541	Fortrain is maintained by ICCS and can be found https://github.com/Cambridge-ICCS/FTorch.
542	The Model of an idealized Moist Atmosphere (MiMA) is maintained by Martin Jucker and is
543	available at https://github.com/mjucker/MiMA. The version of MiMA that uses FTorch for
544	coupling to the PyTorch emulator used in this study can be found at
545	https://github.com/lm2612/MiMA/tree/ML-laura.
546	
547	
548	
549	References
549 550	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,
549 550 551	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,
549 550 551 552	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> .
549 550 551 552 553	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> . https://www.tensorflow.org/
549 550 551 552 553 554	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> </ul>
549 550 551 552 553 554 555	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> </ul>
549 550 551 552 553 554 555 556	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> <li>deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297.</li> </ul>
549 550 551 552 553 554 555 556 557	References         Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,         Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,         M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> .         https://www.tensorflow.org/         Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,         A., Acharya, U. R., Makarenkov, V., & Nahavandi, S. (2021). A review of uncertainty quantification in         deep learning: Techniques, applications and challenges. <i>Information Fusion</i> , 76, 243–297.         https://doi.org/10.1016/j.inffus.2021.05.008
549 550 551 552 553 554 555 556 557 558	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> <li>deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297.</li> <li>https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I.,</li> </ul>
549 550 551 552 553 554 555 556 557 558 559	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>. https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297. https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I., Sato, K., Sheshadri, A., Stephan, C. C., Niekerk, A. van, &amp; Wright, C. J. (2023). Atmospheric Gravity</li> </ul>
549 550 551 552 553 554 555 556 557 558 559 560	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>. https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297. https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I., Sato, K., Sheshadri, A., Stephan, C. C., Niekerk, A. van, &amp; Wright, C. J. (2023). Atmospheric Gravity Waves: Processes and Parameterization. <i>Journal of the Atmospheric Sciences</i>, <i>1</i>(aop).</li> </ul>

- 562 Alexander, M. J., & Dunkerton, T. J. (1999). A Spectral Parameterization of Mean-Flow Forcing due to Breaking
- 563 Gravity Waves. Journal of the Atmospheric Sciences, 56(24), 4167–4182. https://doi.org/10.1175/1520-

564 0469(1999)056<4167:ASPOMF>2.0.CO;2

- Anantharaj, V., Hatfield, S., Polichtchouk, I., Wedi, N., O'Neill, M. E., Papatheodore, T., & Dueben, P. (2022). An
   open science exploration of global 1-km simulations of the earth's atmosphere. 2022 IEEE 18th
- 567 International Conference on E-Science (e-Science), 427–428.
- 568 https://doi.org/10.1109/eScience55777.2022.00071
- 569 Baldwin, M. P., Gray, L. J., Dunkerton, T. J., Hamilton, K., Haynes, P. H., Randel, W. J., Holton, J. R., Alexander,
- 570 M. J., Hirota, I., Horinouchi, T., Jones, D. B. A., Kinnersley, J. S., Marquardt, C., Sato, K., & Takahashi,
- 571 M. (2001). The quasi-biennial oscillation. *Reviews of Geophysics*, 39(2), 179–229.
- 572 https://doi.org/10.1029/1999RG000073
- Bauer, P., Stevens, B., & Hazeleger, W. (2021). A digital twin of Earth for the green transition. *Nature Climate Change*, *11*(2), Article 2. https://doi.org/10.1038/s41558-021-00986-y
- Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and Stabilizing Machine Learning Parametrizations of Convection. *Journal of the Atmospheric Sciences*, 77(12), 4357–4375.

577 https://doi.org/10.1175/JAS-D-20-0082.1

Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially Extended Tests of a Neural Network Parametrization
 Trained by Coarse-Graining. *Journal of Advances in Modeling Earth Systems*, *11*(8), 2728–2744.

580 https://doi.org/10.1029/2019MS001711

- 581 Bushell, A. C., Anstey, J. A., Butchart, N., Kawatani, Y., Osprey, S. M., Richter, J. H., Serva, F., Braesicke, P.,
- 582 Cagnazzo, C., Chen, C.-C., Chun, H.-Y., Garcia, R. R., Gray, L. J., Hamilton, K., Kerzenmacher, T., Kim,
- 583 Y.-H., Lott, F., McLandress, C., Naoe, H., ... Yukimoto, S. (2020). Evaluation of the Quasi-Biennial
- 584 Oscillation in global climate models for the SPARC QBO-initiative. *Quarterly Journal of the Royal*
- 585 *Meteorological Society*, *n/a*(n/a). https://doi.org/10.1002/qj.3765
- Butler, A. H., Seidel, D. J., Hardiman, S. C., Butchart, N., Birner, T., & Match, A. (2015). Defining Sudden
   Stratospheric Warmings. *Bulletin of the American Meteorological Society*, *96*(11), 1913–1928.
- 588 https://doi.org/10.1175/BAMS-D-13-00173.1

- 589 Cambridge-ICCS. (2023). FTorch: A library for coupling (Py)Torch machine learning models to Fortran
- 590 [Computer software]. https://github.com/Cambridge-ICCS/FTorch
- Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Machine Learning Emulation of 591
- 592 Gravity Wave Drag in Numerical Weather Forecasting. Journal of Advances in Modeling Earth Systems, 593 13(7), e2021MS002477. https://doi.org/10.1029/2021MS002477
- 594 Chen, D., Rojas, M., Samset, B. H., Cobb, K., Diongue Niang, A., Edwards, P., Emori, S., Faria, S. H., Hawkins, E.,
- 595 Hope, P., Huybrechts, P., Meinshausen, M., Mustafa, S. K., Plattner, G.-K., & Tréguier, A.-M. (2021).
- 596 Framing, Context, and Methods. In V. Masson-Delmotte, P. Zhai, A. Pirani, S. L. Connors, C. Péan, S.
- 597 Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M. Huang, K. Leitzell, E. Lonnoy, J. B. R. Matthews,
- 598 T. K. Mavcock, T. Waterfield, O. Yelekci, R. Yu, & B. Zhou (Eds.), Climate Change 2021: The Physical
- 599 Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental
- 600 Panel on Climate Change (pp. 147–286). Cambridge University Press.
- 601 https://doi.org/10.1017/9781009157896.003
- 602 Chevallier, F., Morcrette, J.-J., Chéruy, F., & Scott, N. A. (2000). Use of a neural-network-based long-wave 603 radiative-transfer scheme in the ECMWF atmospheric model. Ouarterly Journal of the Royal
- Meteorological Society, 126(563), 761-776. https://doi.org/10.1002/qj.49712656318 604
- Chollet, F. & others. (2015). Keras. GitHub. https://github.com/fchollet/keras 605
- Cohen, N. Y., & Edwin P. Gerber, and O. B. (2013). Compensation between Resolved and Unresolved Wave 606
- 607 Driving in the Stratosphere: Implications for Downward Control. Journal of the Atmospheric Sciences,
- 608 70(12), 3780-3798. https://doi.org/10.1175/JAS-D-12-0346.1
- 609 Delaunay, A., & Christensen, H. M. (2022). Interpretable Deep Learning for Probabilistic MJO Prediction. Geophysical Research Letters, 49(16), e2022GL098566. https://doi.org/10.1029/2022GL098566
- 610
- 611 Dong, W., Fritts, D. C., Liu, A. Z., Lund, T. S., Liu, H.-L., & Snively, J. (2023). Accelerating Atmospheric Gravity
- Wave Simulations Using Machine Learning: Kelvin-Helmholtz Instability and Mountain Wave Sources 612
- Driving Gravity Wave Breaking and Secondary Gravity Wave Generation. Geophysical Research Letters, 613
- 614 50(15), e2023GL104668. https://doi.org/10.1029/2023GL104668
- Donner, L. J., Wyman, B. L., Hemler, R. S., Horowitz, L. W., Ming, Y., Zhao, M., Golaz, J.-C., Ginoux, P., Lin, S.-615
- 616 J., Schwarzkopf, M. D., Austin, J., Alaka, G., Cooke, W. F., Delworth, T. L., Freidenreich, S. M., Gordon,

- 617 C. T., Griffies, S. M., Held, I. M., Hurlin, W. J., ... Zeng, F. (2011). The Dynamical Core, Physical
- 618 Parameterizations, and Basic Simulation Characteristics of the Atmospheric Component AM3 of the GFDL
- 619 Global Coupled Model CM3. *Journal of Climate*, *24*(13), 3484–3519.
- 620 https://doi.org/10.1175/2011JCLI3955.1
- Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J. (2022). Machine Learning Gravity
   Wave Parameterization Generalizes to Capture the QBO and Response to Increased CO2. *Geophysical Research Letters*, 49(8), e2022GL098174. https://doi.org/10.1029/2022GL098174
- Fritts, D. C., & Alexander, M. J. (2003). Gravity wave dynamics and effects in the middle atmosphere. *Reviews of Geophysics*, *41*(1). https://doi.org/10.1029/2001RG000106
- 626 Gagne, D. J., McGovern, A., Haupt, S. E., Sobash, R. A., Williams, J. K., & Xue, M. (2017). Storm-Based
- Probabilistic Hail Forecasting with Machine Learning Applied to Convection-Allowing Ensembles.
   *Weather and Forecasting*, *32*(5), 1819–1840. https://doi.org/10.1175/WAF-D-17-0010.1
- Gagne, D. J., McGovern, A., & Xue, M. (2014). Machine Learning Enhancement of Storm-Scale Ensemble
   Probabilistic Quantitative Precipitation Forecasts. *Weather and Forecasting*, *29*(4), 1024–1043.
   https://doi.org/10.1175/WAF-D-13-00108.1
- 632 Gagne II, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine Learning for
- Stochastic Parameterization: Generative Adversarial Networks in the Lorenz '96 Model. *Journal of Advances in Modeling Earth Systems*, *12*(3), e2019MS001896. https://doi.org/10.1029/2019MS001896
- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could Machine Learning Break the
   Convection Parameterization Deadlock? *Geophysical Research Letters*, 45(11), 5742–5751.
- 637 https://doi.org/10.1029/2018GL078202
- Gneiting, T., & Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, 102(477), 359–378. https://doi.org/10.1198/016214506000001437
- 640 Gordon, E. M., & Barnes, E. A. (2022). Incorporating Uncertainty Into a Regression Neural Network Enables
- Identification of Decadal State-Dependent Predictability in CESM2. *Geophysical Research Letters*, 49(15),
   e2022GL098635. https://doi.org/10.1029/2022GL098635
- Gray, L. J. (2010). Stratospheric Equatorial Dynamics. In *The Stratosphere: Dynamics, Transport, and Chemistry* (pp. 93–107). American Geophysical Union (AGU). https://doi.org/10.1002/97811186666630.ch5

Guillaumin, A. P., & Zanna, L. (2021). Stochastic-Deep Learning Parameterization of Ocean Momentum Forcing.
 *Journal of Advances in Modeling Earth Systems*, 13(9), e2021MS002534.

647 https://doi.org/10.1029/2021MS002534

- 648 Gupta, A., Birner, T., Dörnbrack, A., & Polichtchouk, I. (2021). Importance of Gravity Wave Forcing for
- 649 Springtime Southern Polar Vortex Breakdown as Revealed by ERA5. *Geophysical Research Letters*,

650 *48*(10), e2021GL092762. https://doi.org/10.1029/2021GL092762

- Harder, P., Watson-Parris, D., Stier, P., Strassel, D., Gauger, N. R., & Keuper, J. (2022). Physics-informed learning
  of aerosol microphysics. *Environmental Data Science*, *1*, e20. https://doi.org/10.1017/eds.2022.22
- Hardiman, S. C., Scaife, A. A., Niekerk, A. van, Prudden, R., Owen, A., Adams, S. V., Dunstan, T., Dunstone, N. J.,

654 & Madge, S. (2023). Machine learning for non-orographic gravity waves in a climate model. *Artificial*655 *Intelligence for the Earth Systems*, *1*(aop). https://doi.org/10.1175/AIES-D-22-0081.1

- Hawkins, E., & Sutton, R. (2009). The Potential to Narrow Uncertainty in Regional Climate Predictions. *Bulletin of the American Meteorological Society*, *90*(8), 1095–1108. https://doi.org/10.1175/2009BAMS2607.1
- Haynes, K., Lagerquist, R., McGraw, M., Musgrave, K., & Ebert-Uphoff, I. (2023). Creating and Evaluating
- 659 Uncertainty Estimates with Neural Networks for Environmental-Science Applications. *Artificial*

660 Intelligence for the Earth Systems, 2(2). https://doi.org/10.1175/AIES-D-22-0061.1

- Hersbach, H. (2000). Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction Systems.
   *Weather and Forecasting*, 15(5), 559–570. https://doi.org/10.1175/1520-
- 663 0434(2000)015<0559:DOTCRP>2.0.CO;2
- Hohenegger, C., Korn, P., Linardakis, L., Redler, R., Schnur, R., Adamidis, P., Bao, J., Bastin, S., Behravesh, M.,
- 665 Bergemann, M., Biercamp, J., Bockelmann, H., Brokopf, R., Brüggemann, N., Casaroli, L., Chegini, F.,
- Datseris, G., Esch, M., George, G., ... Stevens, B. (2023). ICON-Sapphire: Simulating the components of
- the Earth system and their interactions at kilometer and subkilometer scales. *Geoscientific Model*
- 668 Development, 16(2), 779–811. https://doi.org/10.5194/gmd-16-779-2023
- 669 Holt, L. A., Lott, F., Garcia, R. R., Kiladis, G. N., Cheng, Y.-M., Anstey, J. A., Braesicke, P., Bushell, A. C.,
- 670 Butchart, N., Cagnazzo, C., Chen, C.-C., Chun, H.-Y., Kawatani, Y., Kerzenmacher, T., Kim, Y.-H.,
- 671 McLandress, C., Naoe, H., Osprey, S., Richter, J. H., ... Yukimoto, S. (2020). An evaluation of tropical

- waves and wave forcing of the QBO in the QBOi models. *Quarterly Journal of the Royal Meteorological Society*, *n/a*(n/a). https://doi.org/10.1002/qj.3827
- Hüllermeier, E., & Waegeman, W. (2021). Aleatoric and epistemic uncertainty in machine learning: An introduction
  to concepts and methods. *Machine Learning*, *110*(3), 457–506. https://doi.org/10.1007/s10994-021-059463
- Jucker, M., & Gerber, E. P. (2017). Untangling the Annual Cycle of the Tropical Tropopause Layer with an
- 678 Idealized Moist Model. Journal of Climate, 30(18), 7339–7358. https://doi.org/10.1175/JCLI-D-17-0127.1
- 679 Krasnopolsky, V. M., & Fox-Rabinovitz, M. S. (2006). Complex hybrid models combining deterministic and
- 680 machine learning components for numerical climate modeling and weather prediction. *Neural Networks*,
- 681 19(2), 122–134. https://doi.org/10.1016/j.neunet.2006.01.002
- Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). *Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles* (arXiv:1612.01474). arXiv. http://arxiv.org/abs/1612.01474
- Lin, J., Yu, S., Beucler, T., Gentine, P., Walling, D., & Pritchard, M. (2023). Systematic Sampling and Validation of
   Machine Learning-Parameterizations in Climate Models (arXiv:2309.16177). arXiv.
- 686 http://arxiv.org/abs/2309.16177
- Lindzen, R. S., & Holton, J. R. (1968). A Theory of the Quasi-Biennial Oscillation. *Journal of the Atmospheric Sciences*, 25(6), 1095–1107. https://doi.org/10.1175/1520-0469(1968)025<1095:ATOTQB>2.0.CO;2
- 689 Mansfield, L. A., & Sheshadri, A. (2022). Calibration and Uncertainty Quantification of a Gravity Wave
- 690 Parameterization: A Case Study of the Quasi-Biennial Oscillation in an Intermediate Complexity Climate
- 691 Model. Journal of Advances in Modeling Earth Systems, 14(11), e2022MS003245.
- 692 https://doi.org/10.1029/2022MS003245
- 693 McGovern, A., Bostrom, A., Davis, P., Demuth, J. L., Ebert-Uphoff, I., He, R., Hickey, J., Ii, D. J. G., Snook, N.,
- 694 Stewart, J. Q., Thorncroft, C., Tissot, P., & Williams, J. K. (2022). NSF AI Institute for Research on
- Trustworthy AI in Weather, Climate, and Coastal Oceanography (AI2ES). *Bulletin of the American Meteorological Society*, *103*(7), E1658–E1668. https://doi.org/10.1175/BAMS-D-21-0020.1
- Murphy, J. M., Booth, B. B. B., Collins, M., Harris, G. R., Sexton, D. M. H., & Webb, M. J. (2007). A methodology
   for probabilistic predictions of regional climate change from perturbed physics ensembles. *Philosophical*

- Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 365(1857), 1993–
  2028. https://doi.org/10.1098/rsta.2007.2077
- Nadiga, B. T., Sun, X., & Nash, C. (2022). Stochastic parameterization of column physics using generative
   adversarial networks. *Environmental Data Science*, *1*, e22. https://doi.org/10.1017/eds.2022.32
- O'Gorman, P. A., & Dwyer, J. G. (2018). Using Machine Learning to Parameterize Moist Convection: Potential for
   Modeling of Climate, Climate Change, and Extreme Events. *Journal of Advances in Modeling Earth*

705 Systems, 10(10), 2548–2563. https://doi.org/10.1029/2018MS001351

Pahlavan, H. A., Hassanzadeh, P., & Alexander, M. J. (2023). *Explainable Offline-Online Training of Neural Networks for Parameterizations: A 1D Gravity Wave-OBO Testbed in the Small-data Regime*

708 (arXiv:2309.09024). arXiv. https://doi.org/10.48550/arXiv.2309.09024

Palmer, T. N. (2019). Stochastic weather and climate models. *Nature Reviews Physics*, 1(7), Article 7.

710 https://doi.org/10.1038/s42254-019-0062-2

711 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L.,

712 Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang,

L., ... Chintala, S. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. In

714 *Advances in Neural Information Processing Systems 32* (pp. 8024–8035). Curran Associates, Inc.

- 715 http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-
- 716 library.pdf
- Perezhogin, P., Zanna, L., & Fernandez-Granda, C. (2023). *Generative data-driven approaches for stochastic subgrid parameterizations in an idealized ocean model* (arXiv:2302.07984). arXiv.

719 http://arxiv.org/abs/2302.07984

Perkins, W. A., Brenowitz, N. D., Bretherton, C. S., & Nugent, J. M. (2023). *Emulation of cloud microphysics in a climate model* [Preprint]. Preprints. https://doi.org/10.22541/essoar.168614667.71811888/v1

722 Polichtchouk, I., Niekerk, A. van, & Wedi, N. (2023). Resolved Gravity Waves in the Extratropical Stratosphere:

Effect of Horizontal Resolution Increase from O(10) to O(1) km. Journal of the Atmospheric Sciences,

724 80(2), 473–486. https://doi.org/10.1175/JAS-D-22-0138.1

725 Rabel, E. (2019). forpy: A library for Fortran-Python interoperability [Computer software].

726 https://github.com/ylikx/forpy

- 727 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models.
- 728 Proceedings of the National Academy of Sciences, 115(39), 9684–9689.

729 https://doi.org/10.1073/pnas.1810286115

- 730 Richter, J. H., Anstey, J. A., Butchart, N., Kawatani, Y., Meehl, G. A., Osprey, S., & Simpson, I. R. (2020).
- Progress in Simulating the Quasi-Biennial Oscillation in CMIP Models. *Journal of Geophysical Research: Atmospheres*, *125*(8), e2019JD032362. https://doi.org/10.1029/2019JD032362
- Schenzinger, V., Osprey, S., Gray, L., & Butchart, N. (2017). Defining metrics of the Quasi-Biennial Oscillation in
   global climate models. *Geoscientific Model Development*, *10*(6), 2157–2168. https://doi.org/10.5194/gmd 10-2157-2017
- 736 Sengupta, K., Pringle, K., Johnson, J. S., Reddington, C., Browse, J., Scott, C. E., & Carslaw, K. (2021). A global
- model perturbed parameter ensemble study of secondary organic aerosol formation. *Atmospheric Chemistry and Physics*, 21(4), 2693–2723. https://doi.org/10.5194/acp-21-2693-2021
- 739 Sexton, D. M. H., McSweeney, C. F., Rostron, J. W., Yamazaki, K., Booth, B. B. B., Murphy, J. M., Regayre, L.,
- Johnson, J. S., & Karmalkar, A. V. (2021). A perturbed parameter ensemble of HadGEM3-GC3.05 coupled
- 741 model projections: Part 1: selecting the parameter combinations. *Climate Dynamics*, *56*(11), 3395–3436.
- 742 https://doi.org/10.1007/s00382-021-05709-9
- 743 Siskind, D., Eckermann, S. D., Coy, L., McCormack, J. P., & Randall, C. E. (2007). On recent interannual
- 744 variability of the Arctic winter mesosphere: Implications for tracer descent: MESOSPHERIC
- 745 INTERANNUAL VARIABILITY. Geophysical Research Letters, 34(9).
- 746 https://doi.org/10.1029/2007GL029293
- 747 Siskind, D., Eckermann, S., McCormack, J., Coy, L., Hoppel, K., & Baker, N. (2010). Case studies of the
- mesospheric response to recent minor, major, and extended stratospheric warmings. J. Geophys. Res, 115,
- 749 0–3. https://doi.org/10.1029/2010JD014114
- Sun, Y. Q., Hassanzadeh, P., Alexander, M. J., & Kruse, C. G. (2023). Quantifying 3D Gravity Wave Drag in a
- 751 Library of Tropical Convection-Permitting Simulations for Data-Driven Parameterizations. *Journal of*
- 752 Advances in Modeling Earth Systems, 15(5), e2022MS003585. https://doi.org/10.1029/2022MS003585

- 753 Ukkonen, P. (2022). Exploring Pathways to More Accurate Machine Learning Emulation of Atmospheric Radiative
- Transfer. Journal of Advances in Modeling Earth Systems, 14(4), e2021MS002875.

755 https://doi.org/10.1029/2021MS002875

- 756 Vallis, G. K., Colyer, G., Geen, R., Gerber, E., Jucker, M., Maher, P., Paterson, A., Pietschnig, M., Penn, J., &
- 757 Thomson, S. I. (2018). Isca, v1.0: A framework for the global modelling of the atmospheres of Earth and
- other planets at varying levels of complexity. *Geoscientific Model Development*, 11(3), 843–859.
- 759 https://doi.org/10.5194/gmd-11-843-2018
- Wang, L., & Alexander, M. J. (2009). Gravity wave activity during stratospheric sudden warmings in the 2007–2008
   Northern Hemisphere winter. *Journal of Geophysical Research: Atmospheres*, *114*(D18).
- 762 https://doi.org/10.1029/2009JD011867
- Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-Seasonal Forecasting With a Large
   Ensemble of Deep-Learning Weather Prediction Models. *Journal of Advances in Modeling Earth Systems*,
   *13*(7), e2021MS002502. https://doi.org/10.1029/2021MS002502
- Whiteway, J. A., Duck, T. J., Donovan, D. P., Bird, J. C., Pal, S. R., & Carswell, A. I. (1997). Measurements of
   gravity wave activity within and around the Arctic stratospheric vortex. *Geophysical Research Letters*,
   24(11), 1387–1390. https://doi.org/10.1029/97GL01322
- 769 Wicker, W., Polichtchouk, I., & Domeisen, D. I. V. (2023). Increased vertical resolution in the stratosphere reveals
- role of gravity waves after sudden stratospheric warmings. *Weather and Climate Dynamics*, 4(1), 81–93.
  https://doi.org/10.5194/wcd-4-81-2023
- Wright, C. J., Osprey, S. M., Barnett, J. J., Gray, L. J., & Gille, J. C. (2010). High Resolution Dynamics Limb
   Sounder measurements of gravity wave activity in the 2006 Arctic stratosphere. *Journal of Geophysical Research: Atmospheres*, *115*(D2). https://doi.org/10.1029/2009JD011858
- Yu, S., Hannah, W. M., Peng, L., Lin, J., Bhouri, M. A., Gupta, R., Lütjens, B., Will, J. C., Behrens, G., Busecke, J.
- J. M., Loose, N., Stern, C., Beucler, T., Harrop, B. E., Hilman, B. R., Jenney, A. M., Ferretti, S. L., Liu, N.,
- 777 Anandkumar, A., ... Pritchard, M. S. (2023). *ClimSim: An open large-scale dataset for training high-*
- 778 resolution physics emulators in hybrid multi-scale climate simulators (arXiv:2306.08754). arXiv.
- 779 https://doi.org/10.48550/arXiv.2306.08754

- 780 Yuval, J., & O'Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate
- 781 modeling at a range of resolutions. *Nature Communications*, *11*(1), Article 1.

782 https://doi.org/10.1038/s41467-020-17142-3

- 783 Yuval, J., O'Gorman, P. A., & Hill, C. N. (2021). Use of Neural Networks for Stable, Accurate and Physically
- 784 Consistent Parameterization of Subgrid Atmospheric Processes With Good Performance at Reduced
- 785 Precision. *Geophysical Research Letters*, *48*(6), e2020GL091363. https://doi.org/10.1029/2020GL091363

786



# AGU Word Manuscript Template

1 2 3 4	Uncertainty Quantification of a Machine Learning Subgrid-Scale Parameterization for Atmospheric Gravity Waves
5	L. A. Mansfield <sup>1</sup> , and A. Sheshadri <sup>1</sup>
6 7	<sup>1</sup> Earth System Science, Doerr School of Sustainability, Stanford University, California
8	Corresponding author: Laura A. Mansfield (lauraman@stanford.edu)
9	Key Points:
10 11	• Using ensembles of neural networks, we learn parametric uncertainties associated with an emulator of a gravity wave parameterization.
12 13	• When coupled to the climate model, the ensemble of neural networks reveals increased climate variability.
14 15 16	• Parametric uncertainty dominates the Quasi-Biennial Oscillation statistics, although polar vortex properties remain robust to parameters.

### 17 Abstract

- 18 Subgrid-scale processes, such as atmospheric gravity waves, play a pivotal role in shaping the
- 19 Earth's climate but cannot be explicitly resolved in climate models due to limitations on
- 20 resolution. Instead, subgrid-scale parameterizations are used to capture their effects. Recently,
- 21 machine learning has emerged as a promising approach to learn parameterizations. In this study,
- 22 we explore uncertainties associated with a machine learning parameterization for atmospheric
- 23 gravity waves. Focusing on the uncertainties in the training process (parametric uncertainty), we
- 24 use an ensemble of neural networks to emulate an existing gravity wave parameterization. We
- estimate both offline uncertainties in raw neural network output and online uncertainties in
- climate model output, after the neural networks are coupled. We find that online parametric
   uncertainty contributes a significant source of uncertainty in climate model output that must be
- considered when introducing neural network parameterizations. This uncertainty quantification
- provides valuable insights into the reliability and robustness of machine learning-based gravity
- 30 wave parameterizations, thus advancing our understanding of their potential applications in
- 31 climate modeling.
- 32

### 33 Plain Language Summary

- 34 Climate models are unable to resolve processes that vary on length and time scales smaller than
- the model resolution and timestep. For example, atmospheric gravity waves, which are waves
- 36 created when winds encounter disturbances to the flow, such as mountains, convection and
- 37 fronts, can have wavelengths smaller than the spacing between grid cells. Climate models use
- <sup>38</sup> "parameterizations" to capture the effect of these processes. Machine learning based
- 39 parameterizations are becoming popular because they can learn relationships purely from data.
- However, we do not have a good understanding of the uncertainties introduced through machine
   learning parameterizations. This study estimates uncertainties associated with training a neural
- learning parameterizations. This study estimates uncertainties associated with training a neural
   network gravity wave parameterization. We explore uncertainties in the neural network output,
- network gravity wave parameterization. We explore uncertainties in the neural network output,
   as well as the uncertainties in the climate model output, when the neural network is used for the
- as well as the uncertainties in the climate model output, when the neural network is used for thgravity wave parameterization.
- 45 **1 Introduction**
- 46

## 1.1. Subgrid-scale parameterizations

Global climate models (GCMs) simulate the entire Earth system by coupling a dynamical 47 core, which numerically solves the primitive equations for atmospheric flow, with other physical 48 components called "subgrid-scale parameterizations". The latter includes dynamical processes 49 occurring on scales smaller than the grid-scale (generally O(100 km) for a typical GCM; Chen 50 et al., 2021), such as convection and short wavelength gravity waves, and non-dynamical 51 processes, such as radiation, atmospheric chemistry, and cloud and aerosol microphysics. 52 Subgrid-scale parameterizations make up a large portion of the computational cost associated 53 with GCM simulations and sometimes make drastic assumptions for the sake of computational 54 cost, which can introduce additional sources of model uncertainty. This has motivated the 55 demand for faster and/or higher accuracy schemes that use machine learning (ML)/artificial 56 intelligence (AI), which hold out the potential for training on large volumes of training data and 57

58 performing fast inferences when invoked.

### 59

ML-based subgrid-scale parameterizations have demonstrated skill across a wide range of 60 atmospheric processes including convection, clouds, aerosols, radiation and gravity waves (e.g., 61 Brenowitz et al., 2020; Brenowitz & Bretherton, 2019; Chantry et al., 2021; Chevallier et al., 62 2000; Espinosa et al., 2022; Gentine et al., 2018; Harder et al., 2022; Krasnopolsky & Fox-63 Rabinovitz, 2006; O'Gorman & Dwyer, 2018; Perkins et al., 2023; Rasp et al., 2018; Ukkonen, 64 2022; Yu et al., 2023; Yuval et al., 2021; Yuval & O'Gorman, 2020). However, few studies have 65 explored the uncertainties associated with these. Stochastic subgrid-scale parameterizations have 66 been developed by sampling from parametric distributions, learned through neural networks 67 (Guillaumin & Zanna, 2021) and generative adversarial networks (GANs) (Gagne II et al., 2020; 68 Nadiga et al., 2022; Perezhogin et al., 2023). These studies focus on stochastic representations to 69 improve model accuracy since they may better represent scaling properties (Palmer, 2019). 70 Including uncertainty estimates can also be beneficial in assessing the trustworthiness of model 71 predictions (Haynes et al., 2023; McGovern et al., 2022), and has gained some attention in 72 weather and climate prediction studies (e.g., Delaunay & Christensen, 2022; Gagne et al., 2014, 73 2017; Gordon & Barnes, 2022; Weyn et al., 2021). Here, we explore uncertainty quantification 74 in a machine learning subgrid-scale parameterization (a type of *model uncertainty*; Hawkins & 75 Sutton, 2009; Palmer, 2019), focusing on gravity wave parameterizations. 76

77 78 79

### **1.2 Atmospheric Gravity Waves**

Atmospheric gravity waves (GWs) are important drivers of middle atmosphere 80 circulation as they transport momentum upwards and away from their sources in the lower 81 troposphere (Fritts & Alexander, 2003). They are forced by perturbations to a stable stratified 82 flow, for instance, orography, convection, and frontogenesis. They propagate primarily in the 83 84 vertical and, due to the decreasing density in the upper atmosphere, grow in amplitude until reaching a critical level, at which point they break and deposit momentum. This provides a 85 forcing on the mean flow in the middle and upper atmosphere and has a substantial impact on 86 atmospheric circulation, including in driving the Quasi-Biennial Oscillation (QBO) in the 87 equatorial stratosphere (Baldwin et al., 2001) and affecting the occurrence of Sudden 88 Stratospheric Warmings in the polar vortex during winter (Wang & Alexander, 2009), described 89 90 further in Section 1.3.

91

GW wavelengths can range from O(1 km) to O(1000 km), which presents a challenge 92 for accurate representation in global climate models (GCMs). While the primitive equations do 93 capture GW dynamics, typical GCM resolutions are O(100 km), resulting in a large portion of 94 the GW spectrum being un- or under-resolved. Parameterizations must be employed to model the 95 impacts of subgrid-scale GWs on the mean flow and are critical for obtaining realistic 96 97 circulation, for example, to induce a spontaneous QBO (Bushell et al., 2020). Some studies find GW parameterizations to be necessary even in kilometer-scale resolution simulations (Achatz et 98 al., 2023; Polichtchouk et al., 2023), suggesting that the need for accurate parameterizations will 99 persist even as modeling centers move towards high resolution GCMs (or "digital twins"; e.g., 100 Bauer et al., 2021). 101 102

- 103 **1.2.2 Gravity wave parameterizations**
- 104

GCMs usually make use of both an orographic and a non-orographic GW parameterization to capture their effects. Machine learning alternatives to GW parameterizations have recently gained attention in several forms. Chantry et al. (2021), Espinosa et al. (2022) and Hardiman et al. (2023) present machine learning emulators of existing non-orographic gravity wave schemes, while Dong et al. (2023) and Sun et al. (2023) use machine learning to learn gravity wave momentum fluxes from high resolution simulations.

111

This study can be viewed as a continuation of the work by Espinosa et al. (2022), which 112 develops an emulator of a non-orographic GW parameterization designed primarily for 113 convectively forced GWs (Alexander & Dunkerton, 1999). Note that this machine learning 114 parameterization is, at best, as accurate as the scheme it aims to emulate and is not significantly 115 faster than the original physics-based scheme, which could be due to coupling of the neural 116 network within a Fortran-based GCM (Cambridge-ICCS, 2023). Rather, this neural network 117 emulator is used as a first step towards probing uncertainties introduced when replacing a gravity 118 wave parameterization with an emulator, when we have a "ground truth" parameterization for 119 reference. 120

121 122

1.3 Gravity wave effects

123 124

## 1.3.1 Quasi-Biennial Oscillation

Gravity waves strongly influence the stratospheric circulation. In the tropical stratosphere, the dominant mode of variability is the Quasi-Biennial Oscillation (QBO), in which the equatorial stratospheric zonal winds alternate between easterly and westerly and descend downwards with time (Gray, 2010). The change in direction is driven by breaking waves across a range of scales (Baldwin et al., 2001; Lindzen & Holton, 1968), with modeling studies suggesting that non-orographic gravity wave parameterizations contribute around half of the forcing required for a simulated QBO (Holt et al., 2020).

133

In this study, we measure the performance of gravity wave parameterizations through the simulated QBO period and amplitudes at 10 hPa, where the QBO amplitude is generally a maximum (Bushell et al., 2020; Richter et al., 2020). We consider the QBO winds to be defined by the zonal mean zonal winds between  $5^{\circ}S$  and  $5^{\circ}N$ . Following Schenzinger et al., (2017), we estimate the period of a QBO cycle by the length between transition times from westward and eastward flow, after applying a 5-month binomial filter to remove high frequency variability. The amplitude is estimated as the absolute maximum of the QBO winds during each cycle.

141 142

## 1.3.2 Stratospheric Polar Vortex

As well as driving the equatorial stratospheric circulation, gravity waves are also influential at high latitudes. Gravity waves affect the stratospheric polar vortex in both hemispheres, as they contribute to the breakdown of the polar vortices, influencing the frequency and properties of Sudden Stratospheric Warmings (SSWs) (Siskind et al., 2007, 2010; Wang & Alexander, 2009; Whiteway et al., 1997; Wright et al., 2010) and the timing of the Spring final warming (Gupta et al., 2021). SSWs are defined as a reversal of the zonal mean zonal winds at 60°N at 10 hPa (Butler et al., 2015) which is followed by large and rapid temperature increases

(>30-40 K) in the polar stratosphere. They occur around 6 times per decade in the Northern 151

152 hemisphere, but are not common in the Southern hemisphere. In this study, we consider gravity

wave parameterization effects on the number of Northern hemisphere SSWs per decade and the 153

timing of the final warming of the Southern hemisphere polar vortex. 154

155 156

158

#### 2. Uncertainty Quantification 157

Uncertainties can be categorized into two types: *aleatoric uncertainty* and *epistemic* 159 uncertainty (Hüllermeier & Waegeman, 2021). Aleatoric uncertainty is used to describe the 160 variability in a system that is due to inherently random effects (Haynes et al., 2023; Hüllermeier 161 & Waegeman, 2021). It represents the statistical or stochastic nature of a system, such as flipping 162 a coin or rolling a dice and in ML literature, refers to uncertainty in the data. It includes *internal* 163 variability of the system and observational uncertainties in the data. In contrast, epistemic 164 uncertainty is caused by a lack of knowledge about the best model for a system and refers to 165 uncertainty in the model. It includes structural uncertainties from the choice of ML architecture, 166 parametric uncertainties in estimating of model parameters, and out-of-sample uncertainties 167 which arise when predicting outside of the range of the training data. 168

In this study, we aim to quantify parametric uncertainty, a type of epistemic uncertainty, 169 in an ML-based parameterization for gravity waves. We expect this to also capture out-of-sample 170 uncertainties, i.e., increased uncertainty when generalizing to a situation that lies outside of the 171 training data distribution. For simplicity, we do not estimate aleatoric uncertainty in the training 172 data, and we also do not consider structural uncertainty. Future studies may wish to account for 173 these additional types of uncertainty for a more complete picture. There are several methods that 174 could be used to estimate parametric uncertainty (Abdar et al., 2021). Here, we use an ensemble 175 of deep neural networks or "deep ensembles", which involves training multiple identical neural 176 networks, each with a different initialization (Lakshminarayanan et al., 2017). Each neural 177 network converges upon slightly different parameters which are then used to predict an 178 ensemble, from which statistics can be obtained. This is a relatively simple approach to 179 implement, although can be costly as it requires repetition during training and evaluation. Deep 180 ensembles have been used in climate model applications for prediction (Weyn et al., 2021), but 181 have not been used for subgrid-scale parameterizations. In this context, deep ensembles could be 182 viewed as a machine learning complement to "perturbed parameter ensembles" (PPE), which 183 involve perturbing physics-based parameters for uncertainty quantification (e.g., Murphy et al., 184 2007; Sengupta et al., 2021; Sexton et al., 2021). 185

#### 186 187 **3** Methods

- 188
- 189
- 190

### **3.1 Gravity Wave Parameterization Setup**

Alexander & Dunkerton (1999; hereafter AD99) present a simple non-orographic, gravity 191 wave parameterization that has been used in various GCMs, including GFDL's Atmospheric 192 Model 3 (Donner et al., 2011), Isca (Vallis et al., 2018), and MiMA (Jucker & Gerber, 2017). 193 AD99 estimates gravity wave drag (GWD) in both the zonal and meridional directions for each 194

195 level in a column, at each grid-cell and timestep. When coupled into a climate model, gravity

wave drag or forcing acts to accelerate or decelerate winds (i.e., it is a wind tendency). As a 196

spectral parameterization, AD99 defines a spectrum of gravity waves at a source level with momentum flux distributed by phase speeds, assumed to follow a Gaussian distribution centered at 0 m/s with half-width 35 m/s. This spectrum of gravity waves propagates upwards until the waves reach the critical level (when the wind speed equals the phase speed of the waves), when breaking occurs and drag is deposited.

202 203

204

## 3.2 Atmospheric Model Setup

We use an intermediate complexity GCM, a Model of an idealized Moist Atmosphere 205 (MiMA) (Jucker & Gerber, 2017). It is run at spectral resolution T42, corresponding to 64 206 latitudes by 128 longitudes (approximately 2.8 degrees or 300 km grid spacing at the equator), 207 with 40 model levels. The level top is 0.18 hPa, with a strong dissipating sponge layer in the 208 upper three levels (0.85-0.18 hPa). AD99 is coupled into MiMA with the parameters described 209 above and with a fixed source level defined to be 315 hPa in the tropics and decreasing in height 210 with latitude, roughly in line with the tropopause. The model is run with an advection timestep of 211 10 minutes and a physics timestep, which includes calling the gravity wave parameterization, of 212 213 3 hours.

214 215

### 3.2 Machine Learning Setup

216 We use the neural network (NN) gravity wave parameterization developed by Espinosa et 217 al. (2022). This is trained on MiMA simulations using the AD99 gravity wave parameterization, 218 described above (Alexander & Dunkerton, 1999). Espinosa et al. (2022) show that the NN 219 emulator, trained on one year of data, achieves an accurate representation of the AD99 scheme 220 both offline and online. For the online tests, Espinosa et al. (2022) replace the original AD99 221 scheme in MiMA with the NN emulator within MiMA and show that these coupled NN 222 simulations produce a Quasi-Biennial Oscillation consistent with original AD99 simulation. 223 Furthermore, when tested on an out-of-sample climate under 4xCO2 forcing, the NN simulations 224 remained stable and reproduced similar changes to the QBO as the AD99 simulations. 225 226

Espinosa et al. (2022) emulate the zonal and meridional GW drag with two independently 227 trained but almost identical fully connected NNs. The inputs to the zonal GW drag network are 228 229 zonal winds at all levels, u, temperature at all levels, T, surface pressure,  $p_s$ , and latitude,  $\lambda$ , and similarly for the meridional GW drag the inputs are meridional winds at all levels,  $v, T, p_s$ , and 230  $\lambda$ . MiMA uses 40 pressure levels, giving a total of 82 inputs into the NN. The architecture 231 consists of four shared hidden layers followed by another four pressure level specific layers (see 232 233 Supporting Information of Espinosa et al., 2022). The network outputs the zonal/meridional GW drag for all 40 pressure levels. Note that the pressure levels closest the surface always predict 234 zero, where there is no GW drag below the source of the GWs. Although these layers are 235 redundant, we include them because the AD99 gravity wave source level changes with latitude to 236 237 follow the approximate level of the tropopause. Following Espinosa et al. (2022), we normalize the input and output data to have a zero mean and standard deviation of 1. For the pressure levels 238 below the source level, where all GW drag values are exactly zero and standard deviation is 239 undefined, we fix the outputs to zero. Although we follow the same architecture as Espinosa et 240 al. (2022), there are some software differences in our implementation. Firstly, we opt for 241 PyTorch (Paszke et al., 2019) rather than Keras and TensorFlow (Abadi et al., 2015; Chollet & 242

others, 2015) for the machine learning library. Secondly, Espinosa et al. (2022) use the forpy

software (Rabel, 2019) to call python code in the fortran-based climate model. This resulted in a

slow-down of roughly 2.5x when replacing AD99 with the NN emulator. Instead, we use FTorch

- (Cambridge-ICCS, 2023), a software package that directly calls the existing Torch C++ interface
   from Fortran resulting in faster inference. We find a 20% slow-down in the NN simulations
- from Fortran resulting in faster inference. We find a 20% slow-down in the NN simulations relative to the AD99 simulations, although we have not explored if this could be optimized
- 248 relative to the AD99 simulations, although we have not explored if this could 249 further.
- 250

In this study, we capture parametric uncertainty of the NN emulator presented in 251 Espinosa et al. (2022) using deep ensembles (Lakshminarayanan et al., 2017). We repeatedly 252 train an ensemble of size 30 independent NNs, each with the same architecture and trained on the 253 same data but with different random seed initializations. The random seed affects the 254 initialization of the NN parameters and the shuffling order of data during training, leading to 255 slightly different parameters when converged. Following Espinosa et al. (2022), we train the 256 NNs with one year of data, selected so that it contains a typical QBO cycle with a period and 257 amplitude similar to the long-term mean period and amplitude. We use the following one year of 258 data for the validation dataset, and the following 20 years are used for the test dataset, requiring 259 22 years of simulation data in total. Figure 1 shows (a) the QBO zonal winds and (b) the QBO 260 261



262

Figure 1 The QBO (a) zonal winds and (b) zonal gravity wave drag for the training, validation,
and test dataset.

265 266 267 **4 Results** 268

269 **4.1 Offline predictions** 

270 271 Figure 2 shows an example of gravity wave drag (GWD) profiles for a single grid cell close to the equator for a) the zonal component and b) the meridional component, with the black 272 273 line indicating the ground truth from the AD99 parameterization and the red line indicating the mean prediction across all NN ensemble members. The orange shading represents 1 standard 274 deviation across all ensemble members. Animations showing the evolution of this GWD profile 275 can be found in the Supporting Materials. The NNs agree well on the gravity wave profiles and 276 the ground truth falls within the 1 standard deviation range for across most model levels for the 277 zonal component. The meridional component generally captures the patterns within the profile 278 279 but is found to be less accurate, even when considering the uncertainty estimates.

- 280
- 281





283 284

285

Figure 2 Example profiles of a) zonal and b) meridional gravity wave drag at one grid-cell and
one timestep in the tropics where the black line indicates the ground truth from the AD99
parameterization, the red line indicates the mean prediction across all neural network ensembles
and the orange shading indicates 1 standard deviation across these ensembles.

290 To measure the errors, we calculate the continuous ranked probability score (CRPS), a generalization of mean absolute error that allows for comparison of probability distributions. The 291 use of CRPS to measure error between a predicted probability distribution and a single ground 292 293 truth has long been used for verification of ensemble weather forecasts (Hersbach, 2000), and 294 has recently been adopted for probabilistic machine learning (Gneiting & Raftery, 2007). Figure 3 shows CRPS for a) zonal and b) meridional gravity wave drag predictions over a range of 295 latitudes. Note the scale of the axis is reduced by 10x relative to the gravity wave drag 296 297 magnitudes in Figure 2. We find lower errors in the lower and mid-stratosphere that increase with height, where gravity wave drag magnitudes also increase. We see good performance across 298 all latitudes. 299



300 Zonal GWD (ms<sup>-2</sup>) Le<sup>-6</sup> Meridional GWD (ms<sup>-2</sup>) Le<sup>-6</sup> 301 Figure 3 Continuous Ranked Probability Score for a) zonal and b) meridional gravity wave drag 302 for different latitudes over the test dataset.

303 304

305

### 4.2 Offline uncertainty estimates

306 One common problem in uncertainty quantification of deep learning algorithms is in ensuring that uncertainty estimates are reasonable, often known as calibration of uncertainty 307 308 (Lakshminarayanan et al., 2017). A well-calibrated machine learning model should predict low 309 uncertainties when errors are small and high uncertainties when errors are large (for instance, 310 when the data is out-of-sample). Figure 4 shows the 1 standard deviation uncertainty estimates against the ensemble mean absolute errors estimated for the test dataset, with the colors 311 312 representing the density of points. Ideally, these should be correlated and lie approximately along the y = x line shown in the dashed line. Points above the y = x line are underconfident and 313 points below are overconfident. Although the errors and predicted uncertainties are correlated, 314 we see that the NNs suffer from overconfidence and frequently underestimate the uncertainty 315 relative to the error. This is typical behavior for machine learning uncertainty estimates, 316 including those based on deep ensembles (Abdar et al., 2021), and may be not be surprising 317 given we only consider one type of uncertainty (parametric uncertainty) and do not consider 318 structural uncertainty or data uncertainty in these estimates. This overconfidence is systematic 319 across all levels of the stratosphere and occurs for both zonal and meridional NNs, but especially 320 for the meridional predictions. 321 322

- 323
- 324

#### manuscript submitted to JAMES



Confidence of neural networks at latitudes 5°S-5°N at 10.9 hPa

Figure 4 Ensemble uncertainty (measured as 1 standard deviation amongst the ensemble predictions) against ensemble error (measured as the mean absolute error across all ensemble predictions) for a) zonal and b) meridional gravity wave drag for test dataset between 5°S-5°N at 10 hPa. Each individual point represents a single prediction at one timestep and grid-cell and they are shaded according to density. The black dashed line shows the y = x line.

331

325

332 333

341

### 4.3 Offline and Online Probability Distributions

Once coupled online into MiMA, the ensembles begin to diverge from each other even though they are initialized from the same state. This is partly due to the chaotic nature of the atmosphere where minute differences in one atmospheric variable can lead to very different atmospheric states after some time. Even introducing relatively minor differences in the GWD profiles, such as those in Figure 2, can lead to very different atmospheric states. Here, we aim to quantify how uncertainties in Figure 2 propagate into the GCM. We examine long-term statistics in order to separate out the NN parametric uncertainty from the internal variability.

342 We consider GWD in the tropics, due to its influence on the QBO. Figure 5 shows distributions of gravity wave drag in the upper stratosphere at 10 hPa for (a) zonal and (b) 343 meridional components, where the black line indicates ground truth from the AD99 MiMA 344 simulations, the blue line indicates the offline NN predicted gravity wave drag and the red line 345 indicates the online NN predicted GWD. Both offline and online distributions are centered over 346 347 the same location as AD99, indicating that the NN does not introduce a bias. In the lower stratosphere, the distributions are virtually indistinguishable (not shown). However, in the upper 348 stratosphere at 10 hPa, the NN distributions take a different shape than AD99. This is 349 particularly notable around the low negative zonal gravity wave drag values, where AD99 350 predicts an asymmetric gravity wave drag distribution with a positive skew. The NN 351 distributions are more symmetric between positive and negative values. This may because 352 353 machine learning optimizes for RMSE which may overly smooth gravity wave drag profiles, reducing asymmetry between positive and negative drag. The online NN distributions are 354 slightly smoother than the offline NN distributions. We suggest that this must be caused by the 355 interaction between the predicted gravity wave drag and the winds when coupled online. This is 356

verified by Figure 6a, which shows distributions of zonal winds near the equator at 10 hPa,
 where online distributions tend to be smoother and weaker than the AD99 distributions.

360 Figure 5 b shows that the online and offline meridional distributions are highly similar, 361 even though they are smoothed out at low magnitudes. Even though the meridional NN is 362 generally less accurate (e.g., Figure 2b), the meridional component of gravity wave drag does not 363 appear to diverge when coupled online. Similarly, Figure 6 b shows the distribution of the 364 meridional winds to be unchanged when the NNs are coupled. This indicates that the meridional 365 circulation is not highly sensitive to the effects of subgrid-scale gravity wave drag, possibly due 366 to lower magnitude of the meridional winds.

- 367
- 368

Distributions of Gravity Wave Drag for Equator at 10.9 hPa



374 5°N.

#### a) Meridional Zonal b) 0.08 AD99 AD99 Online Online 0.07 0.06 0.05 0.04 0.03 0.02 0.01 0.00 -40 40 -60 -40 -20 0 20 40 60 -60 -20 0 20 60 Zonal wind (ms<sup>-1</sup>) Meridional wind (ms<sup>-1</sup>)

Distributions of Wind for Equator at 10.9 hPa

375

Figure 6 a) zonal and b) meridional wind distributions for AD99 (black) and online NN

377 simulations (red) at 10 hPa between 5°S-5°N.

### **4.4 QBO uncertainties**

380 Ultimately, we are interested in how the NN estimations for GWD influence the 381 climatology and its variability when coupled into a GCM. We examine statistics of the OBO in 382 MiMA by calculating the QBO period and amplitude at 10 hPa for each QBO cycle within 400 383 years of AD99 simulations and the 600 years of NN simulations (from 30 simulations each of 20 384 385 year simulations), shown in Figure 7. While the mean period of the QBO across all simulation years are similar, the NN ensembles show increased variability that can be attributed to the 386 387 parametric uncertainty. The NNs also appear to introduce a bias that reduces the QBO amplitude, consistent with the reduction in QBO zonal winds (Figure 6). These increases in QBO variability 388 originate from differences between NN ensemble members (and therefore from the learned NN 389 parameters), each of which tend to maintain fairly consistent OBO periods and amplitudes within 390 the 20 year simulation. 391



378



393

394

Figure 7 Violin plots showing distributions of QBO a) period and b) amplitude for the AD99
simulations in grey and for NN simulations in orange. The boxplots also show the median, upper
and lower quartiles and each point represents a single OBO cycle.

398



Equation 1

401 ensemble of NNs, the additional variability from the uncertainty in parameters,  $\sigma_{param}$ , can be calculated as 402  $\sigma_{param}^2 = \sigma_{AD99}^2 + \sigma_{NNs}^2$ 

403

404

where  $\sigma_{AD99}^2$  is the variance in the AD99 simulations and  $\sigma_{NNs}^2$  is the total variance across 405 all NN ensemble members. These results are shown in Table 1. Notably, the parametric 406 uncertainty is significantly larger than the internal variability in the AD99 simulations, for both 407 the QBO period and amplitude. It is possible that these uncertainties are underestimates of the 408 409 true parametric uncertainty, given the overconfidence noted in offline tests (Figure 4). Still, the 410 uncertainties in NN parameters are much greater than uncertainties in the parameters in the physics-based scheme AD99, estimated to be 1.53 months and 2.14 m/s for the period and 411 amplitude respectively, in Mansfield & Sheshadri (2022) under the same model set-up. This 412 413 highlights the importance of uncertainty quantification, regardless of whether the parameterization is physics-based or machine learning based. 414

415

Table 1 Mean and variability of QBO calculated across MiMA simulations using AD99 vs. the 416

ensemble of NNs. Means are estimated across all QBO cycles in a 400 year long MiMA 417

simulation using AD99 and in 600 years of simulations from the 30-member, 20 year long 418

419 simulations from the ensemble of NNs. Variability is measured as 1 standard deviation between

all OBO cycles. Parametric uncertainty is calculated assuming OBO cycles are normally 420

distributed (Equation 1). 421

	Mean		Variability (mea	asured as 1 stand	ard deviation)
	AD99	Ensemble of NNs	Internal variability in AD99 simulations	Total variability in ensemble of NNs	Parametric uncertainty
Period (months)	25.32	26.78	2.03	3.82	3.25
Amplitude (m/s)	28.29	25.91	2.17	3.86	3.18

422

423 424

425 426

### 4.5 Polar vortex uncertainties

427 The QBO is just one phenomenon that is strongly influenced by gravity wave dynamics. 428 The stratospheric polar vortices in both hemispheres also depend upon gravity wave activity. In 429 particular, the breakdown of the polar vortices during sudden stratospheric warmings (SSWs) 430 and in the springtime final warming is driven by both planetary-scale and subgrid-scale gravity 431 waves, and the variability of these events could also be impacted by changes to the gravity wave 432

### manuscript submitted to JAMES

- 433 parameterization. For the northern hemisphere polar vortex, we consider the frequency of SSWs
- and for the southern hemisphere, we consider polar vortex lifetime. Figure 8 shows there is no
- obvious distinction between the variability of these properties between the AD99 and NN
   simulations, thus making the attribution of extratropical changes (and therefore, the calibration of
- extratropical parameters in AD99 and other schemes; Mansfield & Sheshadri, 2022) rather
- 438 challenging. This may be because the breakdown of the polar vortices is driven by both
- 439 planetary-scale waves and subgrid-scale gravity waves, thereby reducing the impact of any
- 440 changes to the parameterization. Furthermore, some studies find there may be a compensation
- 441 effect between resolved Rossby waves and unresolved gravity waves during SSW events (e.g.,
- 442 Cohen & Gerber, 2013), while some studies suggest that small scale gravity waves influence
- 443 polar vortex recovery after a SSW more strongly than the breakdown itself (Wicker et al., 2023).
- 444



445

446 Figure 8 Histograms showing a) the Northern hemisphere number of SSWs per decade and b)

the Southern hemisphere polar vortex lifetime for AD99 simulations in grey and the NN

448 *simulations in orange.* 

449

450

### 451 **5** Conclusions

452

This study uses deep neural network ensembles to quantify parametric uncertainties in a machine learning parameterization of gravity wave drag. We use the neural network architecture of Espinosa et al. (2022) trained on one year of data simulated by the intermediate complexity GCM, MiMA, which uses AD99 gravity wave parameterization (Alexander & Dunkerton, 1999; Jucker & Gerber, 2017). An ensemble of 30 identical neural networks are trained, each initialized with a different random seed. This ensemble allows us to estimate parametric uncertainties in neural network weights and biases. First, we assessed uncertainties in raw GWD output, which we refer to as *offline uncertainties*. We find fairly consistent results across all neural networks. Then, we used the FTorch library to couple the neural network into MiMA, allowing for GCM simulations that use the machine learning parameterization in place of the traditional physics-based scheme (Cambridge-ICCS, 2023). We assess uncertainties in GCM

- 464 output for gravity wave drag and wind, refering to these as *online uncertainties*. We find 465 increased online uncertainty, particularly for zonal winds.
- 466

487

Comparing long-term statistics of the climate within MiMA using the physics-based 467 scheme AD99 and the ensemble of neural networks, showed that the use of NN emulators can 468 alter the circulation significantly. We found that the NNs from the ensemble produce a bias in 469 the QBO towards reduced amplitudes and dramatically increase the variability of the QBO, with 470 uncertainty from NN parameters increasing the variability between OBO cycles by over 50%. 471 Uncertainty quantification of parameterizations should therefore not be overlooked when 472 developing ML-based schemes for future climate models. Our findings reiterate results from 473 previous studies that find that, even when offline tests indicate "good" NN performance with 474 relatively low uncertainties, the coupling of machine learning schemes into climate models can 475 still introduce a significant source of uncertainty (Brenowitz et al., 2020; Lin et al., 2023). 476 Learning distributions on the model parameters could provide a basis for further parameter 477 refinement, for example, acting as a Bayesian prior distribution that could be constrained through 478 online calibration, such as derivative-free optimization Ensemble Kalman methods (Pahlavan et 479 al., 2023). As with traditional parameterization calibration, this could lead to improved QBO 480 statistics and reduced parametric uncertainty. Interestingly, we find that the behavior and 481 breakdown of the polar vortex is not strongly dependent on the parameterization, which may be 482 partially due to influences from planetary-scale waves. This suggests that it may not be possible 483 to further calibrate neural network parameters to polar vortex properties, and is comparable to 484 the difficulties in calibration of extratropical parameters of AD99 (Mansfield & Sheshadri, 485 2022). 486

We only scratch the surface of uncertainty quantification for machine learning 488 parameterizations. Firstly, we describe only one type of uncertainty: parametric uncertainty, a 489 type of epistemic (model) uncertainty. There exist a wide range of machine learning approaches 490 that could be used for this task, including Bayesian Neural Networks, Monte Carlo dropout 491 generative models and deep ensembles (Abdar et al., 2021). We used deep ensemble methods for 492 this task (Lakshminarayanan et al., 2017), due to their simplicity to implement. However, this 493 approach is computationally costly during both training and evaluation, requiring the use of 494 ensembles which is not feasible for long climate model integrations. A more complete picture 495 would be given by also assessing aleatoric (data) uncertainties. We note that our parametric 496 uncertainty estimates would change given a different training dataset, which makes detangling 497 the effects of epistemic and aleatoric uncertainty a challenge (Haynes et al., 2023; Hüllermeier & 498 Waegeman, 2021). Still, learning the relative contributions between model and data uncertainties 499 would be insightful when designing machine learning parameterizations. Aleatoric uncertainties 500 501 could be estimated through the use of Bayesian neural networks or Monte Carlo dropout (Abdar et al., 2021), by parameterizing gravity wave outputs as a distribution (Guillaumin & Zanna, 502

2021; Haynes et al., 2023), or through generative models such as GANs (Gagne II et al., 2020;
Nadiga et al., 2022; Perezhogin et al., 2023).

505

Secondly, the machine learning parameterization used here is an emulator of an existing 506 scheme, allowing us to compare against a ground truth simulation. Future studies may wish to 507 extend this to train ML models on gravity-wave resolving simulations e.g., with kilometer-scale 508 resolution models such as IFS (Anantharaj et al., 2022), WRF (Sun et al., 2023) or ICON 509 (Hohenegger et al., 2023). When using novel training datasets from high resolution simulations, 510 we do not have online "true" distributions to compare against, which could present challenges 511 when disentangling the various sources of variability. Furthermore, it also raises the issue of 512 understanding the role of aleatoric uncertainty, e.g., in the choice of training data and method for 513 estimating gravity wave drag (Sun et al., 2023). 514

516 Thirdly, MiMA is an intermediate complexity atmospheric circulation model. One may expect that coupling this atmospheric model to other Earth system components, such as the 517 ocean, land, and sea-ice, would introduce further uncertainties. Therefore, we might consider the 518 results presented here as a lower bound on the uncertainties we could expect to see in fully 519 operational Earth system models that employ ML parameterizations. Extending this study to 520 higher complexity Earth system models would be significantly more costly, however, this could 521 be worthwhile towards better informing the design of ML parameterizations, which ultimately 522 could lead to efficient but accurate hybrid GCMs that combine traditional dynamical solvers with 523 novel machine learning parameterizations. 524

525 526

515

### 527 Acknowledgments

- 528 This research was made possible by Schmidt Sciences, a philanthropic initiative founded by Eric
- and Wendy Schmidt, as part of the Virtual Earth System Research Institute (VESRI). AS
- acknowledges support from the National Science Foundation through grant OAC-2004492.
- 531 We would also like to thank our Datawave colleagues, in particular L. Minah Yang and Dave
- 532 Connelly for their work on the PyTorch implementation of the machine learning model, and
- 533 Simon Clifford, Jack Atkinson, Dominic Orchard and others at ICCS, for their help in setting up
- the FTorch coupler with the Fortran-based climate model. We also appreciate the Stanford high
- 535 performance computing resources that made this work possible.
- 536
- 537 **Open Research**

### manuscript submitted to JAMES

538	The code to run simulations, train neural networks and replicate plots presented in this paper is
539	available at https://github.com/lm2612/WaveNet_UQ. The data generated will be made available
540	on the Stanford Digital Repository on publication. The FTorch library for coupling PyTorch to
541	Fortrain is maintained by ICCS and can be found https://github.com/Cambridge-ICCS/FTorch.
542	The Model of an idealized Moist Atmosphere (MiMA) is maintained by Martin Jucker and is
543	available at https://github.com/mjucker/MiMA. The version of MiMA that uses FTorch for
544	coupling to the PyTorch emulator used in this study can be found at
545	https://github.com/lm2612/MiMA/tree/ML-laura.
546	
547	
548	
549	References
549 550	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,
549 550 551	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,
549 550 551 552	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> .
549 550 551 552 553	References Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> . https://www.tensorflow.org/
549 550 551 552 553 554	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> </ul>
549 550 551 552 553 554 555	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> </ul>
549 550 551 552 553 554 555 556	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> <li>deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297.</li> </ul>
549 550 551 552 553 554 555 556 557	References         Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,         Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,         M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i> .         https://www.tensorflow.org/         Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,         A., Acharya, U. R., Makarenkov, V., & Nahavandi, S. (2021). A review of uncertainty quantification in         deep learning: Techniques, applications and challenges. <i>Information Fusion</i> , 76, 243–297.         https://doi.org/10.1016/j.inffus.2021.05.008
549 550 551 552 553 554 555 556 557 558	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M.,</li> <li>Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur,</li> <li>M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>.</li> <li>https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi,</li> <li>A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in</li> <li>deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297.</li> <li>https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I.,</li> </ul>
549 550 551 552 553 554 555 556 557 558 559	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>. https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297. https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I., Sato, K., Sheshadri, A., Stephan, C. C., Niekerk, A. van, &amp; Wright, C. J. (2023). Atmospheric Gravity</li> </ul>
549 550 551 552 553 554 555 556 557 558 559 560	<ul> <li>References</li> <li>Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Zheng, X. (2015). <i>TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems</i>. https://www.tensorflow.org/</li> <li>Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarenkov, V., &amp; Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. <i>Information Fusion</i>, <i>76</i>, 243–297. https://doi.org/10.1016/j.inffus.2021.05.008</li> <li>Achatz, U., Alexander, M. J., Becker, E., Chun, HY., Dörnbrack, A., Holt, L., Plougonven, R., Polichtchouk, I., Sato, K., Sheshadri, A., Stephan, C. C., Niekerk, A. van, &amp; Wright, C. J. (2023). Atmospheric Gravity Waves: Processes and Parameterization. <i>Journal of the Atmospheric Sciences</i>, <i>1</i>(aop).</li> </ul>

- 562 Alexander, M. J., & Dunkerton, T. J. (1999). A Spectral Parameterization of Mean-Flow Forcing due to Breaking
- 563 Gravity Waves. Journal of the Atmospheric Sciences, 56(24), 4167–4182. https://doi.org/10.1175/1520-

564 0469(1999)056<4167:ASPOMF>2.0.CO;2

- Anantharaj, V., Hatfield, S., Polichtchouk, I., Wedi, N., O'Neill, M. E., Papatheodore, T., & Dueben, P. (2022). An
   open science exploration of global 1-km simulations of the earth's atmosphere. 2022 IEEE 18th
- 567 International Conference on E-Science (e-Science), 427–428.
- 568 https://doi.org/10.1109/eScience55777.2022.00071
- 569 Baldwin, M. P., Gray, L. J., Dunkerton, T. J., Hamilton, K., Haynes, P. H., Randel, W. J., Holton, J. R., Alexander,
- 570 M. J., Hirota, I., Horinouchi, T., Jones, D. B. A., Kinnersley, J. S., Marquardt, C., Sato, K., & Takahashi,
- 571 M. (2001). The quasi-biennial oscillation. *Reviews of Geophysics*, 39(2), 179–229.
- 572 https://doi.org/10.1029/1999RG000073
- Bauer, P., Stevens, B., & Hazeleger, W. (2021). A digital twin of Earth for the green transition. *Nature Climate Change*, *11*(2), Article 2. https://doi.org/10.1038/s41558-021-00986-y
- Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and Stabilizing Machine Learning Parametrizations of Convection. *Journal of the Atmospheric Sciences*, 77(12), 4357–4375.

577 https://doi.org/10.1175/JAS-D-20-0082.1

Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially Extended Tests of a Neural Network Parametrization
 Trained by Coarse-Graining. *Journal of Advances in Modeling Earth Systems*, *11*(8), 2728–2744.

580 https://doi.org/10.1029/2019MS001711

- 581 Bushell, A. C., Anstey, J. A., Butchart, N., Kawatani, Y., Osprey, S. M., Richter, J. H., Serva, F., Braesicke, P.,
- 582 Cagnazzo, C., Chen, C.-C., Chun, H.-Y., Garcia, R. R., Gray, L. J., Hamilton, K., Kerzenmacher, T., Kim,
- 583 Y.-H., Lott, F., McLandress, C., Naoe, H., ... Yukimoto, S. (2020). Evaluation of the Quasi-Biennial
- 584 Oscillation in global climate models for the SPARC QBO-initiative. *Quarterly Journal of the Royal*
- 585 *Meteorological Society*, *n/a*(n/a). https://doi.org/10.1002/qj.3765
- Butler, A. H., Seidel, D. J., Hardiman, S. C., Butchart, N., Birner, T., & Match, A. (2015). Defining Sudden
   Stratospheric Warmings. *Bulletin of the American Meteorological Society*, *96*(11), 1913–1928.
- 588 https://doi.org/10.1175/BAMS-D-13-00173.1

- 589 Cambridge-ICCS. (2023). FTorch: A library for coupling (Py)Torch machine learning models to Fortran
- 590 [Computer software]. https://github.com/Cambridge-ICCS/FTorch
- Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Machine Learning Emulation of 591
- 592 Gravity Wave Drag in Numerical Weather Forecasting. Journal of Advances in Modeling Earth Systems, 593 13(7), e2021MS002477. https://doi.org/10.1029/2021MS002477
- 594 Chen, D., Rojas, M., Samset, B. H., Cobb, K., Diongue Niang, A., Edwards, P., Emori, S., Faria, S. H., Hawkins, E.,
- 595 Hope, P., Huybrechts, P., Meinshausen, M., Mustafa, S. K., Plattner, G.-K., & Tréguier, A.-M. (2021).
- 596 Framing, Context, and Methods. In V. Masson-Delmotte, P. Zhai, A. Pirani, S. L. Connors, C. Péan, S.
- 597 Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M. Huang, K. Leitzell, E. Lonnoy, J. B. R. Matthews,
- 598 T. K. Mavcock, T. Waterfield, O. Yelekci, R. Yu, & B. Zhou (Eds.), Climate Change 2021: The Physical
- 599 Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental
- 600 Panel on Climate Change (pp. 147–286). Cambridge University Press.
- 601 https://doi.org/10.1017/9781009157896.003
- 602 Chevallier, F., Morcrette, J.-J., Chéruy, F., & Scott, N. A. (2000). Use of a neural-network-based long-wave 603 radiative-transfer scheme in the ECMWF atmospheric model. Ouarterly Journal of the Royal
- Meteorological Society, 126(563), 761-776. https://doi.org/10.1002/qj.49712656318 604
- Chollet, F. & others. (2015). Keras. GitHub. https://github.com/fchollet/keras 605
- Cohen, N. Y., & Edwin P. Gerber, and O. B. (2013). Compensation between Resolved and Unresolved Wave 606
- 607 Driving in the Stratosphere: Implications for Downward Control. Journal of the Atmospheric Sciences,
- 608 70(12), 3780-3798. https://doi.org/10.1175/JAS-D-12-0346.1
- 609 Delaunay, A., & Christensen, H. M. (2022). Interpretable Deep Learning for Probabilistic MJO Prediction. Geophysical Research Letters, 49(16), e2022GL098566. https://doi.org/10.1029/2022GL098566
- 610
- 611 Dong, W., Fritts, D. C., Liu, A. Z., Lund, T. S., Liu, H.-L., & Snively, J. (2023). Accelerating Atmospheric Gravity
- Wave Simulations Using Machine Learning: Kelvin-Helmholtz Instability and Mountain Wave Sources 612
- Driving Gravity Wave Breaking and Secondary Gravity Wave Generation. Geophysical Research Letters, 613
- 614 50(15), e2023GL104668. https://doi.org/10.1029/2023GL104668
- Donner, L. J., Wyman, B. L., Hemler, R. S., Horowitz, L. W., Ming, Y., Zhao, M., Golaz, J.-C., Ginoux, P., Lin, S.-615
- 616 J., Schwarzkopf, M. D., Austin, J., Alaka, G., Cooke, W. F., Delworth, T. L., Freidenreich, S. M., Gordon,

- 617 C. T., Griffies, S. M., Held, I. M., Hurlin, W. J., ... Zeng, F. (2011). The Dynamical Core, Physical
- 618 Parameterizations, and Basic Simulation Characteristics of the Atmospheric Component AM3 of the GFDL
- 619 Global Coupled Model CM3. *Journal of Climate*, *24*(13), 3484–3519.
- 620 https://doi.org/10.1175/2011JCLI3955.1
- Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J. (2022). Machine Learning Gravity
   Wave Parameterization Generalizes to Capture the QBO and Response to Increased CO2. *Geophysical Research Letters*, 49(8), e2022GL098174. https://doi.org/10.1029/2022GL098174
- Fritts, D. C., & Alexander, M. J. (2003). Gravity wave dynamics and effects in the middle atmosphere. *Reviews of Geophysics*, *41*(1). https://doi.org/10.1029/2001RG000106
- 626 Gagne, D. J., McGovern, A., Haupt, S. E., Sobash, R. A., Williams, J. K., & Xue, M. (2017). Storm-Based
- Probabilistic Hail Forecasting with Machine Learning Applied to Convection-Allowing Ensembles.
   *Weather and Forecasting*, *32*(5), 1819–1840. https://doi.org/10.1175/WAF-D-17-0010.1
- Gagne, D. J., McGovern, A., & Xue, M. (2014). Machine Learning Enhancement of Storm-Scale Ensemble
   Probabilistic Quantitative Precipitation Forecasts. *Weather and Forecasting*, *29*(4), 1024–1043.
   https://doi.org/10.1175/WAF-D-13-00108.1
- 632 Gagne II, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine Learning for
- Stochastic Parameterization: Generative Adversarial Networks in the Lorenz '96 Model. *Journal of Advances in Modeling Earth Systems*, *12*(3), e2019MS001896. https://doi.org/10.1029/2019MS001896
- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could Machine Learning Break the
   Convection Parameterization Deadlock? *Geophysical Research Letters*, 45(11), 5742–5751.
- 637 https://doi.org/10.1029/2018GL078202
- Gneiting, T., & Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, 102(477), 359–378. https://doi.org/10.1198/016214506000001437
- 640 Gordon, E. M., & Barnes, E. A. (2022). Incorporating Uncertainty Into a Regression Neural Network Enables
- Identification of Decadal State-Dependent Predictability in CESM2. *Geophysical Research Letters*, 49(15),
   e2022GL098635. https://doi.org/10.1029/2022GL098635
- Gray, L. J. (2010). Stratospheric Equatorial Dynamics. In *The Stratosphere: Dynamics, Transport, and Chemistry* (pp. 93–107). American Geophysical Union (AGU). https://doi.org/10.1002/97811186666630.ch5

Guillaumin, A. P., & Zanna, L. (2021). Stochastic-Deep Learning Parameterization of Ocean Momentum Forcing.
 *Journal of Advances in Modeling Earth Systems*, 13(9), e2021MS002534.

647 https://doi.org/10.1029/2021MS002534

- 648 Gupta, A., Birner, T., Dörnbrack, A., & Polichtchouk, I. (2021). Importance of Gravity Wave Forcing for
- 649 Springtime Southern Polar Vortex Breakdown as Revealed by ERA5. *Geophysical Research Letters*,

650 *48*(10), e2021GL092762. https://doi.org/10.1029/2021GL092762

- Harder, P., Watson-Parris, D., Stier, P., Strassel, D., Gauger, N. R., & Keuper, J. (2022). Physics-informed learning
  of aerosol microphysics. *Environmental Data Science*, *1*, e20. https://doi.org/10.1017/eds.2022.22
- Hardiman, S. C., Scaife, A. A., Niekerk, A. van, Prudden, R., Owen, A., Adams, S. V., Dunstan, T., Dunstone, N. J.,

654 & Madge, S. (2023). Machine learning for non-orographic gravity waves in a climate model. *Artificial*655 *Intelligence for the Earth Systems*, *1*(aop). https://doi.org/10.1175/AIES-D-22-0081.1

- Hawkins, E., & Sutton, R. (2009). The Potential to Narrow Uncertainty in Regional Climate Predictions. *Bulletin of the American Meteorological Society*, *90*(8), 1095–1108. https://doi.org/10.1175/2009BAMS2607.1
- Haynes, K., Lagerquist, R., McGraw, M., Musgrave, K., & Ebert-Uphoff, I. (2023). Creating and Evaluating
- 659 Uncertainty Estimates with Neural Networks for Environmental-Science Applications. *Artificial*

660 Intelligence for the Earth Systems, 2(2). https://doi.org/10.1175/AIES-D-22-0061.1

- Hersbach, H. (2000). Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction Systems.
   *Weather and Forecasting*, 15(5), 559–570. https://doi.org/10.1175/1520-
- 663 0434(2000)015<0559:DOTCRP>2.0.CO;2
- Hohenegger, C., Korn, P., Linardakis, L., Redler, R., Schnur, R., Adamidis, P., Bao, J., Bastin, S., Behravesh, M.,
- 665 Bergemann, M., Biercamp, J., Bockelmann, H., Brokopf, R., Brüggemann, N., Casaroli, L., Chegini, F.,
- Datseris, G., Esch, M., George, G., ... Stevens, B. (2023). ICON-Sapphire: Simulating the components of
- the Earth system and their interactions at kilometer and subkilometer scales. *Geoscientific Model*
- 668 Development, 16(2), 779–811. https://doi.org/10.5194/gmd-16-779-2023
- 669 Holt, L. A., Lott, F., Garcia, R. R., Kiladis, G. N., Cheng, Y.-M., Anstey, J. A., Braesicke, P., Bushell, A. C.,
- 670 Butchart, N., Cagnazzo, C., Chen, C.-C., Chun, H.-Y., Kawatani, Y., Kerzenmacher, T., Kim, Y.-H.,
- 671 McLandress, C., Naoe, H., Osprey, S., Richter, J. H., ... Yukimoto, S. (2020). An evaluation of tropical

- waves and wave forcing of the QBO in the QBOi models. *Quarterly Journal of the Royal Meteorological Society*, *n/a*(n/a). https://doi.org/10.1002/qj.3827
- Hüllermeier, E., & Waegeman, W. (2021). Aleatoric and epistemic uncertainty in machine learning: An introduction
  to concepts and methods. *Machine Learning*, *110*(3), 457–506. https://doi.org/10.1007/s10994-021-059463
- Jucker, M., & Gerber, E. P. (2017). Untangling the Annual Cycle of the Tropical Tropopause Layer with an
- 678 Idealized Moist Model. Journal of Climate, 30(18), 7339–7358. https://doi.org/10.1175/JCLI-D-17-0127.1
- 679 Krasnopolsky, V. M., & Fox-Rabinovitz, M. S. (2006). Complex hybrid models combining deterministic and
- 680 machine learning components for numerical climate modeling and weather prediction. *Neural Networks*,
- 681 19(2), 122–134. https://doi.org/10.1016/j.neunet.2006.01.002
- Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). *Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles* (arXiv:1612.01474). arXiv. http://arxiv.org/abs/1612.01474
- Lin, J., Yu, S., Beucler, T., Gentine, P., Walling, D., & Pritchard, M. (2023). Systematic Sampling and Validation of
   Machine Learning-Parameterizations in Climate Models (arXiv:2309.16177). arXiv.
- 686 http://arxiv.org/abs/2309.16177
- Lindzen, R. S., & Holton, J. R. (1968). A Theory of the Quasi-Biennial Oscillation. *Journal of the Atmospheric Sciences*, 25(6), 1095–1107. https://doi.org/10.1175/1520-0469(1968)025<1095:ATOTQB>2.0.CO;2
- 689 Mansfield, L. A., & Sheshadri, A. (2022). Calibration and Uncertainty Quantification of a Gravity Wave
- 690 Parameterization: A Case Study of the Quasi-Biennial Oscillation in an Intermediate Complexity Climate
- 691 Model. Journal of Advances in Modeling Earth Systems, 14(11), e2022MS003245.
- 692 https://doi.org/10.1029/2022MS003245
- 693 McGovern, A., Bostrom, A., Davis, P., Demuth, J. L., Ebert-Uphoff, I., He, R., Hickey, J., Ii, D. J. G., Snook, N.,
- 694 Stewart, J. Q., Thorncroft, C., Tissot, P., & Williams, J. K. (2022). NSF AI Institute for Research on
- Trustworthy AI in Weather, Climate, and Coastal Oceanography (AI2ES). *Bulletin of the American Meteorological Society*, *103*(7), E1658–E1668. https://doi.org/10.1175/BAMS-D-21-0020.1
- Murphy, J. M., Booth, B. B. B., Collins, M., Harris, G. R., Sexton, D. M. H., & Webb, M. J. (2007). A methodology
   for probabilistic predictions of regional climate change from perturbed physics ensembles. *Philosophical*

- Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 365(1857), 1993–
  2028. https://doi.org/10.1098/rsta.2007.2077
- Nadiga, B. T., Sun, X., & Nash, C. (2022). Stochastic parameterization of column physics using generative
   adversarial networks. *Environmental Data Science*, *1*, e22. https://doi.org/10.1017/eds.2022.32
- O'Gorman, P. A., & Dwyer, J. G. (2018). Using Machine Learning to Parameterize Moist Convection: Potential for
   Modeling of Climate, Climate Change, and Extreme Events. *Journal of Advances in Modeling Earth*

705 Systems, 10(10), 2548–2563. https://doi.org/10.1029/2018MS001351

Pahlavan, H. A., Hassanzadeh, P., & Alexander, M. J. (2023). *Explainable Offline-Online Training of Neural Networks for Parameterizations: A 1D Gravity Wave-OBO Testbed in the Small-data Regime*

708 (arXiv:2309.09024). arXiv. https://doi.org/10.48550/arXiv.2309.09024

Palmer, T. N. (2019). Stochastic weather and climate models. *Nature Reviews Physics*, 1(7), Article 7.

710 https://doi.org/10.1038/s42254-019-0062-2

711 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L.,

712 Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang,

- L., ... Chintala, S. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. In
- 714 *Advances in Neural Information Processing Systems 32* (pp. 8024–8035). Curran Associates, Inc.
- 715 http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-
- 716 library.pdf
- Perezhogin, P., Zanna, L., & Fernandez-Granda, C. (2023). *Generative data-driven approaches for stochastic subgrid parameterizations in an idealized ocean model* (arXiv:2302.07984). arXiv.

719 http://arxiv.org/abs/2302.07984

Perkins, W. A., Brenowitz, N. D., Bretherton, C. S., & Nugent, J. M. (2023). *Emulation of cloud microphysics in a climate model* [Preprint]. Preprints. https://doi.org/10.22541/essoar.168614667.71811888/v1

722 Polichtchouk, I., Niekerk, A. van, & Wedi, N. (2023). Resolved Gravity Waves in the Extratropical Stratosphere:

Effect of Horizontal Resolution Increase from O(10) to O(1) km. Journal of the Atmospheric Sciences,

724 80(2), 473–486. https://doi.org/10.1175/JAS-D-22-0138.1

725 Rabel, E. (2019). forpy: A library for Fortran-Python interoperability [Computer software].

726 https://github.com/ylikx/forpy

- 727 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models.
- 728 Proceedings of the National Academy of Sciences, 115(39), 9684–9689.

729 https://doi.org/10.1073/pnas.1810286115

- 730 Richter, J. H., Anstey, J. A., Butchart, N., Kawatani, Y., Meehl, G. A., Osprey, S., & Simpson, I. R. (2020).
- Progress in Simulating the Quasi-Biennial Oscillation in CMIP Models. *Journal of Geophysical Research: Atmospheres*, *125*(8), e2019JD032362. https://doi.org/10.1029/2019JD032362
- Schenzinger, V., Osprey, S., Gray, L., & Butchart, N. (2017). Defining metrics of the Quasi-Biennial Oscillation in
   global climate models. *Geoscientific Model Development*, *10*(6), 2157–2168. https://doi.org/10.5194/gmd 10-2157-2017
- 736 Sengupta, K., Pringle, K., Johnson, J. S., Reddington, C., Browse, J., Scott, C. E., & Carslaw, K. (2021). A global
- model perturbed parameter ensemble study of secondary organic aerosol formation. *Atmospheric Chemistry and Physics*, 21(4), 2693–2723. https://doi.org/10.5194/acp-21-2693-2021
- 739 Sexton, D. M. H., McSweeney, C. F., Rostron, J. W., Yamazaki, K., Booth, B. B. B., Murphy, J. M., Regayre, L.,
- Johnson, J. S., & Karmalkar, A. V. (2021). A perturbed parameter ensemble of HadGEM3-GC3.05 coupled
- 741 model projections: Part 1: selecting the parameter combinations. *Climate Dynamics*, *56*(11), 3395–3436.
- 742 https://doi.org/10.1007/s00382-021-05709-9
- 743 Siskind, D., Eckermann, S. D., Coy, L., McCormack, J. P., & Randall, C. E. (2007). On recent interannual
- 744 variability of the Arctic winter mesosphere: Implications for tracer descent: MESOSPHERIC
- 745 INTERANNUAL VARIABILITY. Geophysical Research Letters, 34(9).
- 746 https://doi.org/10.1029/2007GL029293
- 747 Siskind, D., Eckermann, S., McCormack, J., Coy, L., Hoppel, K., & Baker, N. (2010). Case studies of the
- mesospheric response to recent minor, major, and extended stratospheric warmings. J. Geophys. Res, 115,
- 749 0–3. https://doi.org/10.1029/2010JD014114
- Sun, Y. Q., Hassanzadeh, P., Alexander, M. J., & Kruse, C. G. (2023). Quantifying 3D Gravity Wave Drag in a
- 751 Library of Tropical Convection-Permitting Simulations for Data-Driven Parameterizations. *Journal of*
- 752 Advances in Modeling Earth Systems, 15(5), e2022MS003585. https://doi.org/10.1029/2022MS003585

- 753 Ukkonen, P. (2022). Exploring Pathways to More Accurate Machine Learning Emulation of Atmospheric Radiative
- Transfer. Journal of Advances in Modeling Earth Systems, 14(4), e2021MS002875.

755 https://doi.org/10.1029/2021MS002875

- 756 Vallis, G. K., Colyer, G., Geen, R., Gerber, E., Jucker, M., Maher, P., Paterson, A., Pietschnig, M., Penn, J., &
- 757 Thomson, S. I. (2018). Isca, v1.0: A framework for the global modelling of the atmospheres of Earth and
- other planets at varying levels of complexity. *Geoscientific Model Development*, 11(3), 843–859.
- 759 https://doi.org/10.5194/gmd-11-843-2018
- Wang, L., & Alexander, M. J. (2009). Gravity wave activity during stratospheric sudden warmings in the 2007–2008
   Northern Hemisphere winter. *Journal of Geophysical Research: Atmospheres*, *114*(D18).
- 762 https://doi.org/10.1029/2009JD011867
- Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-Seasonal Forecasting With a Large
   Ensemble of Deep-Learning Weather Prediction Models. *Journal of Advances in Modeling Earth Systems*,
   *13*(7), e2021MS002502. https://doi.org/10.1029/2021MS002502
- Whiteway, J. A., Duck, T. J., Donovan, D. P., Bird, J. C., Pal, S. R., & Carswell, A. I. (1997). Measurements of
   gravity wave activity within and around the Arctic stratospheric vortex. *Geophysical Research Letters*,
   24(11), 1387–1390. https://doi.org/10.1029/97GL01322
- 769 Wicker, W., Polichtchouk, I., & Domeisen, D. I. V. (2023). Increased vertical resolution in the stratosphere reveals
- role of gravity waves after sudden stratospheric warmings. *Weather and Climate Dynamics*, 4(1), 81–93.
  https://doi.org/10.5194/wcd-4-81-2023
- Wright, C. J., Osprey, S. M., Barnett, J. J., Gray, L. J., & Gille, J. C. (2010). High Resolution Dynamics Limb
   Sounder measurements of gravity wave activity in the 2006 Arctic stratosphere. *Journal of Geophysical Research: Atmospheres*, *115*(D2). https://doi.org/10.1029/2009JD011858
- Yu, S., Hannah, W. M., Peng, L., Lin, J., Bhouri, M. A., Gupta, R., Lütjens, B., Will, J. C., Behrens, G., Busecke, J.
- J. M., Loose, N., Stern, C., Beucler, T., Harrop, B. E., Hilman, B. R., Jenney, A. M., Ferretti, S. L., Liu, N.,
- 777 Anandkumar, A., ... Pritchard, M. S. (2023). *ClimSim: An open large-scale dataset for training high-*
- 778 resolution physics emulators in hybrid multi-scale climate simulators (arXiv:2306.08754). arXiv.
- 779 https://doi.org/10.48550/arXiv.2306.08754

- 780 Yuval, J., & O'Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate
- 781 modeling at a range of resolutions. *Nature Communications*, *11*(1), Article 1.

782 https://doi.org/10.1038/s41467-020-17142-3

- 783 Yuval, J., O'Gorman, P. A., & Hill, C. N. (2021). Use of Neural Networks for Stable, Accurate and Physically
- 784 Consistent Parameterization of Subgrid Atmospheric Processes With Good Performance at Reduced
- 785 Precision. *Geophysical Research Letters*, *48*(6), e2020GL091363. https://doi.org/10.1029/2020GL091363

786