

Data Imbalance, Uncertainty Quantification, and Generalization via Transfer Learning in Data-driven Parameterizations: Lessons from the Emulation of Gravity Wave Momentum Transport in WACCM

Y. Qiang Sun¹, Hamid Pahlavan², Ashesh Chattopadhyay³, Pedram Hassanzadeh¹, Sandro W. Lubis⁴, M. Joan Alexander⁵, Edwin P Gerber⁶, Aditi Sheshadri⁷, and Yifei Guan¹

¹Rice University

²Unknown

³UC Santa Cruz

⁴Pacific Northwest National Laboratory (DOE)

⁵NorthWest Research Associates, CoRA Office

⁶New York University

⁷Stanford University

December 27, 2023

Abstract

Neural networks (NNs) are increasingly used for data-driven subgrid-scale parameterization in weather and climate models. While NNs are powerful tools for learning complex nonlinear relationships from data, there are several challenges in using them for parameterizations. Three of these challenges are 1) data imbalance related to learning rare (often large-amplitude) samples; 2) uncertainty quantification (UQ) of the predictions to provide an accuracy indicator; and 3) generalization to other climates, e.g., those with higher radiative forcing. Here, we examine the performance of methods for addressing these challenges using NN-based emulators of the Whole Atmosphere Community Climate Model (WACCM) physics-based gravity wave (GW) parameterizations as the test case. WACCM has complex, state-of-the-art parameterizations for orography-, convection- and frontal-driven GWs. Convection- and orography-driven GWs have significant data imbalance due to the absence of convection or orography in many grid points. We address data imbalance using resampling and/or weighted loss functions, enabling the successful emulation of parameterizations for all three sources. We demonstrate that three UQ methods (Bayesian NNs, variational auto-encoders, and dropouts) provide ensemble spreads that correspond to accuracy during testing, offering criteria on when a NN gives inaccurate predictions. Finally, we show that the accuracy of these NNs decreases for a warmer climate (4XCO₂). However, the generalization accuracy is significantly improved by applying transfer learning, e.g., re-training only one layer using ~1% new data from the warmer climate. The findings of this study offer insights for developing reliable and generalizable data-driven parameterizations for various processes, including (but not limited) to GWs.

Abstract

Neural networks (NNs) are increasingly used for data-driven subgrid-scale parameterization in weather and climate models. While NNs are powerful tools for learning complex nonlinear relationships from data, there are several challenges in using them for parameterizations. Three of these challenges are 1) data imbalance related to learning rare (often large-amplitude) samples; 2) uncertainty quantification (UQ) of the predictions to provide an accuracy indicator; and 3) generalization to other climates, e.g., those with higher radiative forcing. Here, we examine performance of methods for addressing these challenges using NN-based emulators of the Whole Atmosphere Community Climate Model (WACCM) physics-based gravity wave (GW) parameterizations as the test case. WACCM has complex, state-of-the-art parameterizations for orography-, convection- and frontal-driven GWs. Convection- and orography-driven GWs have significant data imbalance due to the absence of convection or orography in many grid points. We address data imbalance using resampling and/or weighted loss functions, enabling the successful emulation of parameterizations for all three sources. We demonstrate that three UQ methods (Bayesian NNs, variational auto-encoders, and dropouts) provide ensemble spreads that correspond to accuracy during testing, offering criteria on when a NN gives inaccurate predictions. Finally, we show that accuracy of these NNs decreases for a warmer climate ($4\times\text{CO}_2$). However, the generalization accuracy is significantly improved by applying transfer learning, e.g., re-training only one layer using $\sim 1\%$ new data from the warmer climate. The findings of this study offer insights for developing reliable and generalizable data-driven parameterizations for various processes, including (but not limited) to GWs.

Plain Language Summary

Scientists are increasingly using machine learning methods, especially neural networks (NNs), to improve weather and climate models. However, it can be challenging for a NN to learn rare, large-amplitude events, because they are infrequent in training data. Also, NNs need to express their confidence (certainty) about a prediction and work effectively across different climates, e.g., warmer climates due to increased CO_2 . Traditional NNs often struggle with these challenges. Here, we share insights gained from emulating the complex physics-based parameterization schemes for gravity waves in a state-of-the-art climate model. We propose specific strategies for addressing imbalanced data, uncertainty quantification (UQ), and making accurate predictions across various climates. For instance, to manage data balance, one such strategy involves amplifying the impact of infrequent events in the training data. We also demonstrate that several UQ methods could be useful in determining the accuracy of predictions. Furthermore, we show that NNs trained on simulations of the historical period do not perform as well in warmer climates. However, we improve the NNs' performance by employing transfer learning using limited data from warmer climates. This study provides lessons for developing robust and generalizable approaches for using NNs to improve models in the future.

1 Introduction

Small-scale processes such as moist convection, gravity waves, and turbulence are key players in the variability of the climate system and its response to increased greenhouse gases. However, as these processes cannot be resolved, entirely or partially, by the coarse-resolution general circulation models (GCMs), they need to be represented as functions of the resolved dynamics via subgrid-scale (SGS) parameterization schemes (e.g., Kim et al., 2003; Stensrud, 2007; Prein et al., 2015). Many of these parameterization schemes are based on heuristic approximations and simplifications, introducing large parametric and epistemic uncertainties in GCMs (Schneider et al., 2017; Hourdin et al., 2017; Palmer, 2019).

Recently, there has been a growing interest in developing data-driven SGS parameterizations for different complex processes in the Earth system using machine learning (ML) techniques, particularly deep neural networks (NNs). Promising results have been demonstrated in a wide range of idealized applications, including prototype systems (Maulik et al., 2019; Gagne et al., 2020; Rasp, 2020; Chattopadhyay, Subel, & Hassanzadeh, 2020; Frezat et al., 2022; Guan et al., 2022; Pahlavan et al., 2023), ocean turbulent processes (Bolton & Zanna, 2019; C. Zhang et al., 2023), moist convection in the atmosphere (O’Gorman & Dwyer, 2018; Brenowitz & Bretherton, 2019; Yuval & O’Gorman, 2020; Beucler et al., 2021; Iglesias-Suarez et al., 2023), radiation (Krasnopolsky et al., 2005; Belochitski & Krasnopolsky, 2021; Song & Roh, 2021), and microphysics (Seifert & Rasp, 2020; Gettelman et al., 2021). The ultimate promise of data-driven parameterizations, learned from observation-derived data and/or high-fidelity high-resolution simulations, is that they might have smaller parametric/structural errors, thus reducing the biases of GCMs and producing more reliable climate change projections (e.g., Schneider et al., 2017; Reichstein et al., 2019; Schneider et al., 2021).

However, there are major challenges in developing trustworthy, interpretable, stable, and generalizable data-driven parameterizations that can be used for such climate change projection efforts. Discussing and even listing all of these challenges is well beyond the scope of this paper. Well-known challenges such as interpretability and stability have been extensively discussed in a number of recent studies (e.g., McGovern et al., 2019; Beck et al., 2019; Brenowitz et al., 2020; Balaji, 2021; Clare et al., 2022; Mamalakis et al., 2022; Guan et al., 2022; Subel et al., 2023; Pahlavan et al., 2023). Here, we focus on three other key issues:

1. Data imbalance, and related to that, learning rare/extreme events,
2. Uncertainty quantification (UQ) of the NN-based SGS parameterization outputs,
3. Out-of-distribution (OOD) generalization (e.g., extrapolation to climates with higher radiative forcings).

Below we briefly discuss the importance of 1-3 and the current state-of-the-art methods in addressing them in the climate and ML literature. Data imbalance is a well-known problem in the ML literature, especially in the context of classification tasks (e.g., Japkowicz &

98 Stephen, 2002; G. Wu & Chang, 2003; Chawla et al., 2004; Sun et al., 2009; Huang et al.,
 99 2016; Ando & Huang, 2017; Buda et al., 2018; Johnson & Khoshgoftaar, 2019). The prob-
 100 lem becomes particularly significant when one aims to learn rare/extreme events (Maalouf
 101 & Trafalis, 2011; Maalouf & Siddiqi, 2014; Baldi et al., 2014; Liu et al., 2016; O’Gorman &
 102 Dwyer, 2018; Qi & Majda, 2020; Chattopadhyay, Nabizadeh, & Hassanzadeh, 2020; Milo-
 103 shevich et al., 2023; Finkel et al., 2023; Shamekh et al., 2023). For example, suppose we
 104 aim to learn the binary classification of the 99 percentile of temperature anomalies using a
 105 NN. In this case, label 0 (no extreme) will constitute 99% of the training (or testing) set
 106 while label 1 (extreme) will be just 1%. With many common loss functions such as mean
 107 squared error (MSE) or root-mean squared-error (RMSE), training a NN will result in one
 108 that predicts 0 for any sample (extreme or no extreme) while having a seemingly high ac-
 109 curacy of 99% (of course, other metrics such as precision/recall will show the shortcoming,
 110 see Chattopadhyay, Nabizadeh, & Hassanzadeh (2020)). The most common remedy to this
 111 problem for classification tasks is resampling. An example is down-sampling non-extreme
 112 cases by a factor of 100, which effectively balances the dataset.

113 In addition to *classification* tasks, Data imbalance also presents a significant challenge
 114 in *regression* tasks required for parameterization schemes in climate models. As highlighted
 115 by Chantry et al. (2021), such imbalances contributed to the unsuccessful emulation of
 116 their orographic gravity wave parameterization (GWP) scheme, largely because orography
 117 affects the gravity wave (GW) drag in only a fraction of the grid columns. This challenge also
 118 persists in emulating GWP for non-orographic GWs, especially when GWs are intricately
 119 linked to their sources. For instance, the presence of zero convective GW drag at numerous
 120 grid points due to the absence of convection creates a notably imbalanced dataset. This
 121 issue will be explored further in the results section. In regression tasks, data imbalance
 122 may also manifest in the form of difficulty in learning large-amplitude (extreme) outputs,
 123 which are rare and constitute only a small fraction of the training set. In the case of GWs,
 124 Observations have shown that gravity wave amplitudes are highly intermittent such that
 125 the largest 10% events alone can contribute more than 50% of the total momentum flux
 126 (Hertzog et al., 2012), so the extreme events will contribute an outsized fraction of the
 127 total drag. Nonetheless, poorly learning these large-amplitude outputs, like drag forces,
 128 can result in instabilities (e.g., Guan et al., 2022). Addressing data imbalance in climate
 129 applications has received relatively limited attention. In this study, we propose several
 130 remedies based on resampling techniques and weighted loss functions, demonstrating their
 131 advantages in enabling successful emulations of all GWP schemes and improving the learning
 132 of rare/extreme events.

133 Quantifying the uncertainties in outputs from NN-based parameterization schemes is
 134 essential when employing these schemes, particularly for high-stakes decision-making tasks
 135 such as climate change projections. Crucially, during testing when we are unable to di-
 136 rectly determine a prediction’s accuracy, we need a UQ method that can provide a credible
 137 *confidence level* for each prediction, serving as a reliable indicator of its accuracy. During
 138 inference, the output of an NN can be inaccurate for various reasons, including poor approx-
 139 imation (e.g., due to poor NN architecture), poor within-distribution generalization (e.g.,

140 for inputs that are rare events), or poor optimization (collectively referred to as *epistemic*
 141 *uncertainty*), as well as because of OOD generalization errors due to input samples from
 142 a distribution different from that of the training set (Abdar et al., 2021; Lu et al., 2021;
 143 Krueger et al., 2021; Miller et al., 2021; Shen et al., 2021; D. Wu et al., 2021; Ye et al., 2021;
 144 D. Zhang et al., 2021; Subel et al., 2023). Quantifying the level of uncertainty would then
 145 allow us to avoid using a data-driven parameterization scheme when it is inaccurate due to
 146 one of the aforementioned reasons (Maddox et al., 2019; Zhu et al., 2019; Li et al., 2022;
 147 Psaros et al., 2023). In the context of data-driven parameterization in climate modeling,
 148 the two most challenging sources of uncertainty are rare/extreme events and OOD gener-
 149 alization errors. The latter is a concern, particularly when the GCM is used for climate
 150 change studies (see below for more discussions).

151 Developing UQ methods for NNs is an active area of research in the ML community,
 152 and there is not a generally applicable rigorous method yet. For instance, techniques like
 153 Markov-Chain Monte Carlo can be prohibitively expensive, especially when dealing with
 154 high-dimensional systems (Oh et al., 2005; Ballnus et al., 2017; Chen & Majda, 2019). For a
 155 comprehensive review in the context of scientific ML, refer to Psaros et al. (2023). The topic
 156 has also started to increasingly gain attention in the climate literature (Guillaumin & Zanna,
 157 2021; Gordon & Barnes, 2022; Haynes et al., 2023; Barnes et al., 2023). In this study, we will
 158 assess the performance of three common UQ methods (Bayesian, dropout, and variational
 159 NNs) by analyzing the relationships between uncertainty and accuracy during inference
 160 testing. We will also consider scenarios involving OOD generalization errors resulting from
 161 global warming.

162 As already mentioned above, OOD generalization (extrapolation to a test data distri-
 163 bution different from that of the training set) is a major challenge for applications involving
 164 non-stationarity, like a changing climate. Studies have already shown that the lack of OOD
 165 generalization in data-driven parameterizations leads to inaccurate and unstable simula-
 166 tion (Rasp et al., 2018; O’Gorman & Dwyer, 2018; Chattopadhyay, Subel, & Hassanzadeh,
 167 2020; Guan et al., 2022; Nagarajan et al., 2020). A general and powerful method for im-
 168 proving the OOD generalization capability of NNs is transfer learning (TL), which involves
 169 re-training a few or all of the layers of a NN using a small amount of data from the new
 170 system (Yosinski et al., 2014). This approach has already shown remarkable success in
 171 enabling data-driven parameterization schemes to extrapolate across the parameter space
 172 (e.g., to $100\times$ higher Reynolds number) in canonical test cases (Chattopadhyay, Subel, &
 173 Hassanzadeh, 2020; Subel et al., 2021; Guan et al., 2023; Subel et al., 2023; C. Zhang et al.,
 174 2023). In particular, Subel et al. (2023) introduced SpArK (Spectral Analysis of Regression
 175 Kernels and Activations) showing that re-training even one layer can lead to successful OOD
 176 generalization, although this optimal layer, unlike the rule of thumb in the ML literature,
 177 may not be the deepest but the shallowest hidden layer. Here, we further leverage these
 178 studies and show how TL can enable OOD generalization of data-driven parameterization
 179 schemes in state-of-the-art GCMs.

180 The methods used in this study and the learned lessons apply to a broad range of
 181 processes and applications in climate modeling. However, the results are presented for a
 182 single test case, that is based on the emulation of complex physics-based GWP schemes in
 183 version 6 of the Whole Atmosphere Community Climate Model (WACCM), a state-of-the-
 184 art GCM (Gettelman et al., 2019). Here, we use the emulations of current physics-based
 185 parameterization schemes as a stepping stone towards learning data-driven parameteriza-
 186 tions from observations and high-fidelity simulations by testing ideas for addressing items
 187 1-3 listed earlier. Furthermore, developing better representations of un- and under-resolved
 188 GWs in GCMs is an important problem on its own (Kim et al., 2003; Alexander et al., 2010;
 189 Achatz, 2022). A number of recent studies have taken the first steps in learning data-driven
 190 GWP from observations and high-resolution simulations (Matsuoka et al., 2020; Amiramjadi
 191 et al., 2022; Sun et al., 2023; Dong et al., 2023), though careful and time-consuming steps
 192 are needed in producing, analyzing, and using such data. Furthermore, two recent stud-
 193 ies focused on emulators of simpler GWP schemes in a forecast model and idealized GCM
 194 have readily shown the usefulness of lessons learned from emulators (Chantry et al., 2021;
 195 Espinosa et al., 2022; Hardiman et al., 2023). This further motivates the focus on using
 196 emulators for testing ideas for addressing data imbalance, UQ, and OOD generalization.

197 This paper is structured as follows. Section 2 introduces the WACCM simulations and
 198 the NN architectures used in this study. The findings, detailed in Section 3, emphasize
 199 the insights gained in addressing data imbalance and UQ, alongside OOD generalization of
 200 the emulators under warmer climate conditions. Consistent with Chantry et al. (2021), we
 201 find that using an NN to emulate the parameterization of orographic GWs is significantly
 202 more challenging than non-orographic GWs. This necessitated additional steps to achieve
 203 reasonable offline performance, as detailed in Section 4. To the best of our knowledge, this
 204 stands as the first NN-based emulation of orographic GWs to address the challenges in
 205 Chantry et al. (2021). Finally, we provide a concluding summary in Section 5.

206 **2 Data and Methods**

207 **2.1 The Whole Atmosphere Community Climate Model (WACCM)**

208 The NCAR’s WACCM version 6 introduced in Gettelman et al. (2019) is used in this
 209 study. WACCM has state-of-the-art GWP schemes for GWs from three different sources:
 210 orography (OGWs), convection (CGWs), and fronts (FGWs). These complex sources make
 211 the emulation of the GWP schemes in WACCM a challenging task. This is, therefore, a
 212 suitable test case to investigate ideas for learning rare events, UQ, and OOD generalization to
 213 benefit the future efforts for the much more complex task, that is learning data-driven GWP
 214 schemes from observations and/or high-resolution GW-resolving simulations (Amiramjadi
 215 et al., 2022; Sun et al., 2023).

216 The configuration of the WACCM used in this study is identical to the public version in
 217 Gettelman et al. (2019), with a horizontal resolution of $0.95^\circ \times 1.25^\circ$ and 70 vertical levels.
 218 The two non-orographic GWP schemes in WACCM both follow Richter et al. (2010), yet

219 allow separate specifications of FGW and CGW sources. For OGWs, WACCM uses an up-
 220 dated planetary boundary layer form drag scheme from Beljaars et al. (2004), near-surface
 221 nonlinear drag processes following Scinocca & McFarlane (2000), and a ridge-finding algo-
 222 rithm to define orographic sources based on Bacmeister et al. (1994). A full documentation
 223 of WACCM OGWs can also be found in Kruse et al. (2022).

224 We conduct two sets of simulations: A 10-year pre-industrial “control” run, and a
 225 10-year pseudo-global-warming “future” run with $4\times\text{CO}_2$ and uniform +4 K sea-surface
 226 temperature increases. In each run, we save, on the native grid, all the inputs and outputs
 227 for each of the three GWPs every 3 hours to capture the diurnal cycle. A complete list of
 228 these inputs/outputs, which are used in the training of the NN-based emulators, is presented
 229 in Appendix A.

230 We train separate NNs for emulating the three GWP schemes that have different
 231 sources. We use the first 6 years of the control run for training and the last 4 years for
 232 validation (years 7 and 8) and testing (years 9 and 10). With a grid resolution of $\sim 1^\circ$,
 233 there are 55,296 columns for each time snapshot, resulting in approximately 960 million
 234 input/output columns during the 6-year training period. Given the strong temporal cor-
 235 relation between the 3-hourly samples, we perform sub-sampling on both the training and
 236 validation data to reduce the dataset size. To accomplish this, we begin by shuffling all
 237 the input/output column pairs in time at each latitude/longitude grid point. Then, we
 238 randomly select 2,000 input/output pairs at each location for training and 500 pairs for
 239 validation.

240 To give the readers a general idea of the parameterized GWs and large-scale circulation
 241 in WACCM, Figure 1 shows the zonal-mean climatology for zonal GW drag/forcing, here-
 242 inafter referred to as GWD, arises from the divergence of gravity wave momentum transport
 243 (fluxes), from all three sources, computed from the 6-year training period in the control run.
 244 The zonal-mean zonal wind climatology is also shown. Seasonal dependency for both the
 245 GWD and the circulation is observed in the simulations. At levels below 100 hPa, the ten-
 246 dencies of non-orographic GW are relatively small compared to those from OGWs; however,
 247 their amplitudes increase significantly at higher altitudes. While the parameterized effect of
 248 GWs is generally to decelerate the zonal flow, there are exceptions, notably in regions like
 249 the equatorward flanks of the stratospheric polar night jets, where FGWs can accelerate the
 250 flow. For more information on the GWP schemes and circulations in WACCM, see Garcia
 251 et al. (2017) and Gettelman et al. (2019).

252 **2.2 The NNs and UQ**

253 ***2.2.1 The Deterministic Fully Connected NN***

254 Here we briefly describe the general structure of the NN-based regression models trained
 255 as emulators for GWP schemes. For the deterministic artificial NN, denoted as ANN in this
 256 study, we use multilayer perceptrons (MLP). MLPs, which are feedforward fully connected
 257 NNs, take inputs through successive layers of linear transformation and non-linear activation

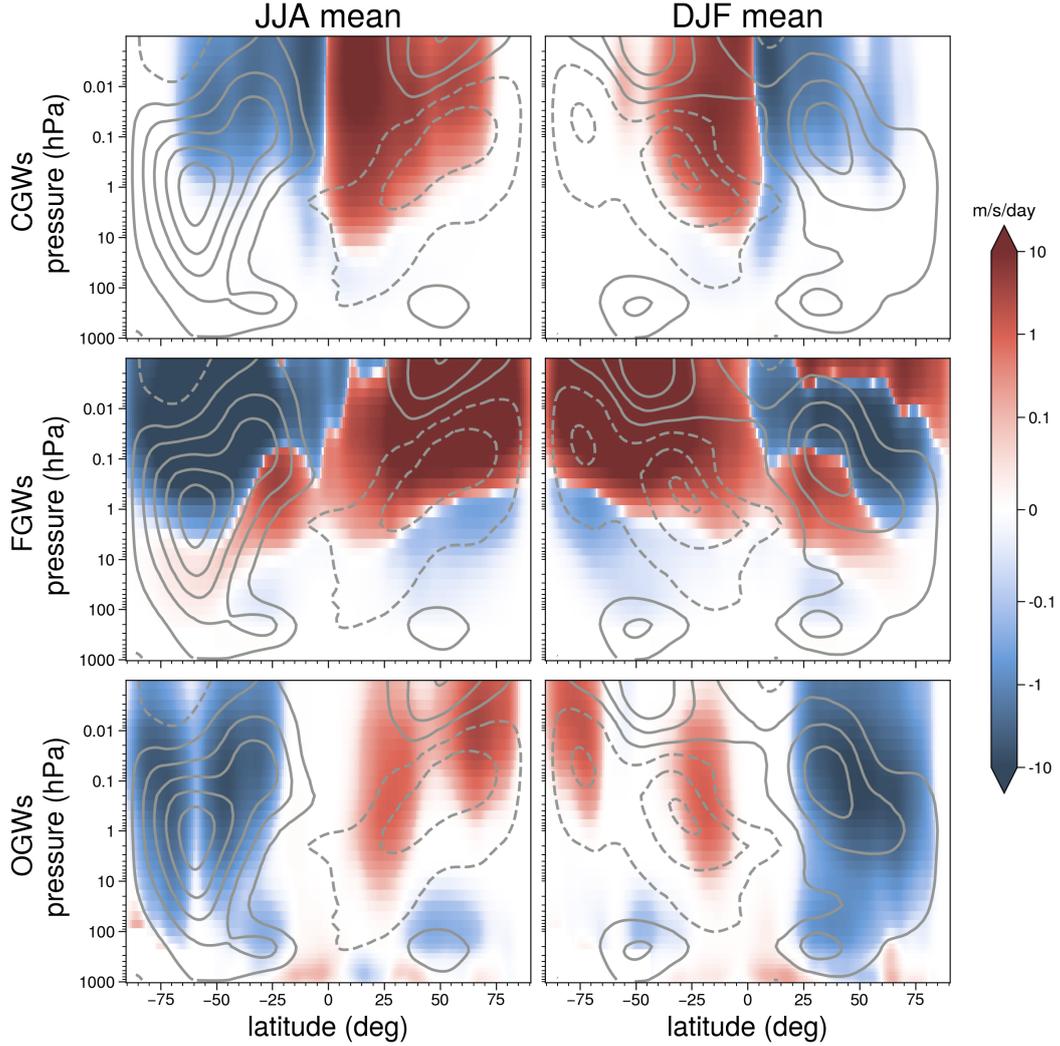


Figure 1. Climatology of zonal-mean GWD during summers (JJA) and winters (DJF) from all 3 sources in the control (pre-industrial) WACCM simulations. Top: CGWs; middle: FGWs; bottom: OGWs. The climatology of the zonal-mean zonal wind is also shown (grey lines), with an interval of 20 m/s. Dashed lines indicate negative values. Zero lines are omitted.

258 functions to produce an output, so as to learn a functional relationship between the input
 259 and output (Figure 2a). Deep MLPs have multiple layers of weights, which are optimized
 260 over many samples of input-output data pairs. Such MLPs are thus very powerful in terms
 261 of learning complicated functional relationships. Generally, we can write the governing
 262 equations of an MLP as

$$z^\ell = \sigma(W^\ell z^{\ell-1} + b^\ell), \quad (1)$$

263 where z^ℓ is the activation (output) of layer ℓ , W^ℓ is the weight matrix connecting layers ℓ
 264 and $\ell - 1$, and b^ℓ is the bias at layer ℓ , which allows the network to fit the data even when
 265 all input features are equal to 0. σ is the non-linear activation function.

In this study, we employ the same NN structure while training three distinct NNs, each for GWP originating from one of the three unique GW sources. The input layer contains the same input variables (see Appendix A) used by the WACCM GWPs across all vertical levels. There are 10 hidden layers in total (Figure 2a), and there are 500 neurons in each hidden layer. In the output layer, both zonal and meridional GWD are predicted. The activation function in each layer, σ , is chosen to be swish (Ramachandran et al., 2017), except for the output layer, where it is linear. During training, W^ℓ and b^ℓ are randomly initialized and learned by minimizing a loss function using an ADAM optimizer, with a fixed learning rate of $\alpha = 0.0001$. One of the loss functions used here is the common MSE, i.e.,

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n \left\| \text{NN}(x_i, \Theta) - y_i \right\|_2^2 \quad (2)$$

266 Here, n is the number of training samples and $\|\cdot\|_2$ is the L_2 norm. For training sample
 267 i , vector x_i contains all the inputs to the NN (Appendix A), vector y_i contains the true
 268 zonal and meridional GWD at each vertical level, and $\Theta = \{\theta_j\}_{j=1\dots p}$ denotes the trainable
 269 parameters, i.e., the weights ($p \approx 3 \times 10^6$).

270 **2.2.2 The UQ Methods and Metrics**

271 Although deterministic NNs are powerfully expressive and can exhibit high out-of-
 272 sample predictive skills, they do not provide estimates of the uncertainty associated with
 273 their predictions. As mentioned earlier, currently there is no rigorous method to estimate
 274 the uncertainty of an NN prediction. That said, a variety of techniques have been developed
 275 for UQ in NNs, though the validity and usefulness of the estimated uncertainty for scientific
 276 applications remain subjects of ongoing investigations (e.g., Psaros et al., 2023; Haynes et
 277 al., 2023). In this paper, we use three different and widely used approaches to perform UQ
 278 from the ML literature: Bayesian neural network (BNN), dropout neural network (DNN),
 279 and variational auto-encoder (VAE). A brief overview of these approaches is provided below.

280 *Bayesian neural network (BNN)*: A BNN combines the deterministic NN described
 281 earlier and in Figure 2a with Bayesian inference (Blundell et al., 2015). Simply speaking, a
 282 BNN estimates distributions of the weights, rather than point values (as in a deterministic
 283 NN). The posterior distributions in the BNN (i.e., the distributions of the weights and
 284 biases) are calculated using the Bayes rule. In this study, we follow the standard practice
 285 and assume that all variational forms of the posterior are normal distributions. Furthermore,
 286 to accelerate the training process, we use the normal distribution $\mathcal{N}(\mu, 1)$ for all the priors in
 287 the BNN (where μ is obtained from parameters of the trained deterministic NN). Note that
 288 while we are assuming normal distributions for the trainable parameters, the predictions
 289 generated by BNN can fit different distributions due to the use of nonlinear activation
 290 functions. The resulting distribution of the predictions during inference gives an estimate
 291 of their uncertainty.

292 *Dropout neural network (DNN)*: A DNN is developed by randomly eliminating all out-
 293 going connections from some of the nodes (Figure 2a) in each hidden layer of a deterministic

294 NN during the training and the inference (Srivastava et al., 2014). The fraction of nodes
 295 “dropped” on average in each layer is called the dropout ratio. Mathematically, Equation (1)
 296 can be reformulated for a DNN as:

$$z^\ell = \sigma(D^\ell W^\ell z^{\ell-1} + b^\ell), \quad (3)$$

297 where the dropout matrix D^ℓ is a square diagonal binary matrix of integers 0 or 1. The
 298 diagonal elements of D^ℓ follow a Bernoulli distribution where the probability of zero is the
 299 dropout ratio.

300 Dropout was initially developed as a regularization technique to prevent over-fitting in
 301 NNs. However, Gal & Ghahramani (2016) showed that training a NN with the dropout
 302 technique approximates a Bayesian NN. In this study, we use a dropout rate of 0.1, which
 303 is incorporated in all hidden layers, but we also investigate the sensitivity of the DNN to
 304 different dropout rates, as later shown in Appendix B. Note that the random dropping out
 305 is also used during inference, leading to a distribution for each prediction.

306 *Variational auto-encoder (VAE)*: A typical VAE (Kingma & Welling, 2014) consists
 307 of two NNs (Figure 2b): an encoder that transforms the input into a lower-dimensional
 308 latent space, parameterized by a normal probability distribution, and a decoder that inverts
 309 this transformation and produces the original input. The difference between the decoder’s
 310 output and the original input drives the learning process of the encoder and decoder, while
 311 the parameterized lower-dimensional latent space provides the uncertainty of this transfor-
 312 mation. The VAE was developed for generative reconstructions of data by simply drawing
 313 samples from the latent space. The VAE is basically a dimension-reduction method. Many
 314 variants, however, have been proposed for more specific purposes. In this study, following
 315 Foster et al. (2021), we add a third NN, as illustrated in Figure 2b, that randomly draws
 316 samples from the parameterized latent space as inputs, and predicts the zonal and merid-
 317 ional GWDs as outputs. The difference between the predicted GWDs and the true GWDs
 318 drives the learning of the third NN. Consequently, the loss for the entire network consists
 319 of three components: the loss between the reconstructed input and the original input, the
 320 Kullback–Leibler (KL) divergence between the distribution of the latent space and a nor-
 321 mal distribution, and the loss between the predicted GWDs by the third NN and the true
 322 GWDs.

323 For a specific input, each of these three UQ methods discussed above can be run multiple
 324 times, generating an ensemble of predictions with different realizations of the weights by
 325 drawing from the trained distribution. This is in contrast to the deterministic NN that
 326 provides just a single-valued prediction for a given input. These ensembles can then be used
 327 to quantify the uncertainty associated with that prediction. We expect that the RMSE of
 328 the ensemble mean should exhibit approximately a 1-1 relationship with the ensemble spread
 329 (i.e., the standard deviation of the ensemble members). To investigate this relationship, we
 330 use the spread-skill plot (Delle Monache et al., 2013). Detailed calculations behind the
 331 spread-skill plot can be found in Appendix C, where we also introduce two metrics: spread-

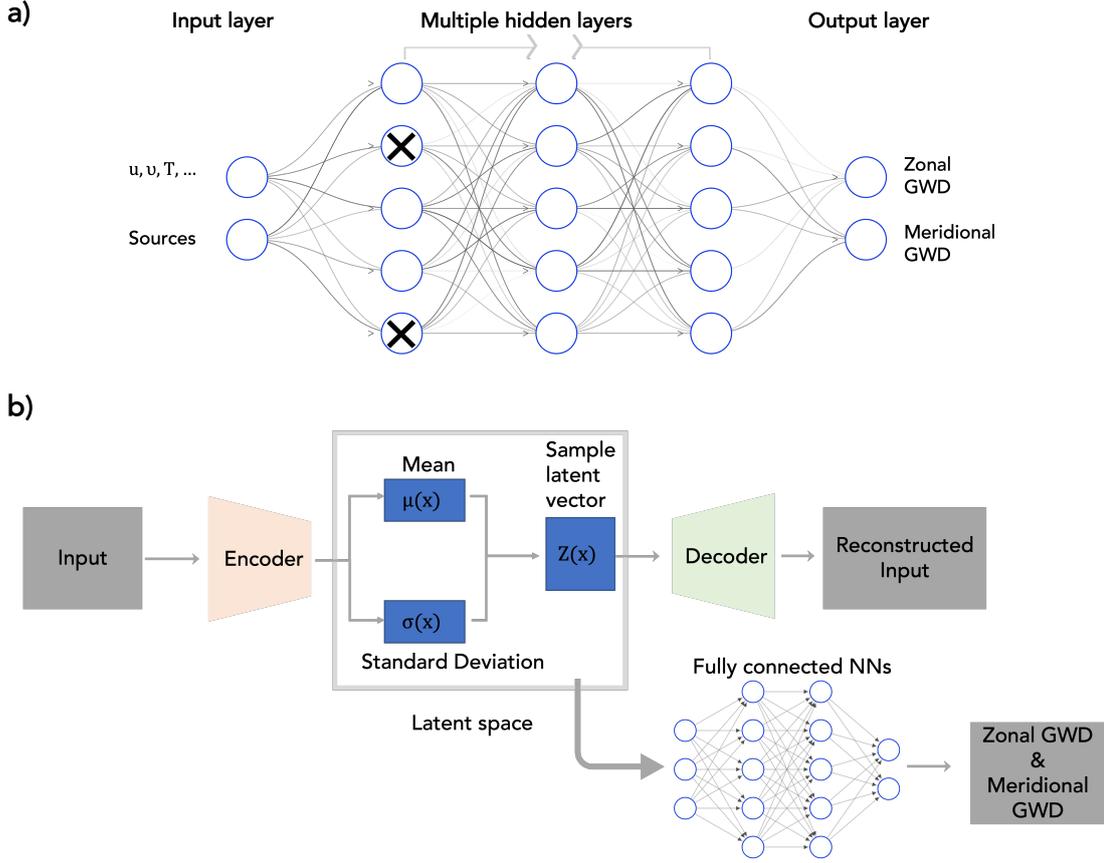


Figure 2. Schematics of the NN-based emulators and different training/re-training strategies used in this study. (a) Schematic for the MLP and DNN. The inputs of the NN are connected through successive layers of neurons (blue circles) to the output (GWDs). A fully connected MLP NN is trained from randomly initialized weights and biases in all layers. A DNN is the same but some connections are randomly eliminated during training and inference (black crosses). In TL, only some of the layers of a previously trained MLP are re-trained using new data. (b) Schematic for the VAE. A low-dimensional latent space is constructed and then used as the input for the additional fully connected NNs, which is similar to the one in (a).

332 skill reliability (SSREL) and overall spread-skill ratio (SSRAT), both of which summarize
 333 the information presented in the spread-skill plot.

334 2.3 Transfer Learning

335 Transfer learning refers to leveraging/reusing information (weights) from an already
 336 well-trained base NN to effectively build a new NN for a different system from which only a
 337 small amount of training data is available (Yosinski et al., 2014; Tan et al., 2018; Chattopad-
 338 hyay, Subel, & Hassanzadeh, 2020). For our purpose, which is improving OOD generalization
 339 to the warmer climate, the TL procedure is as follows. For any of the NNs described earlier
 340 (e.g., the one in Figure 2a), we train them from randomly initialized weights and biases
 341 with data from the control simulations. The NN will work well during inference for test

342 samples from the control but not from future (warmer climate) simulations (as shown in
 343 the Results section). To address this, TL is applied wherein most of the NN’s weights are
 344 kept constant, and only one or two hidden layers are re-trained using a limited dataset from
 345 the future simulation. Although this small dataset is insufficient for training an entire NN
 346 from random initialization, careful and correct selection of hidden layers for re-training, as
 347 discussed in Subel et al. (2023), allows the development of an NN that accurately adapts to
 348 the new system, i.e., the future climate conditions.

349 Here, we re-train the NN-based emulator that was initially trained on the control data
 350 with new data from only 1 month (30 consecutive days) of integration (1.4% of 6 years
 351 simulation for the initial training) of WACCM model under future forcing ($4\times\text{CO}_2$). We
 352 have explored different choices of layers to re-train with the same amount of new data and
 353 found that re-training the first hidden layer yields the best results, consistent with Subel et
 354 al. (2023). Therefore, the results with only re-training the first hidden layer are shown in
 355 Section 3 unless stated otherwise.

356 **3 Results**

357 **3.1 Data Imbalance**

358 As discussed earlier, the physics-based GWP schemes in WACCM are directly linked to
 359 their sources. This means they only produce non-zero values when their respective sources
 360 are active. For example, in a specific grid box, CGWs only register non-zero values when
 361 there is active convection within that box. The heterogeneous and sometimes intermittent
 362 nature of these sources leads to a dataset that is significantly imbalanced. Figure 3 shows
 363 global maps of the occurrence frequency of non-zero GWD for CGWs and FGWs. On
 364 average, only 7.6% of all GCM columns yield non-zero CGWs, primarily located in the
 365 tropics. Similarly, for FGWs, only 8.5% of all columns have non-zero outputs, but unlike
 366 CGWs, the majority of these are located in mid-to-high latitudes, particularly along the
 367 storm track region. For the OGWs in WACCM, data imbalance presents a greater challenge,
 368 to be discussed in a later section. While it is possible to simply separate zero and non-zero
 369 columns for emulation work where we know the truth, this approach falls short with real-
 370 world data, which is the main purpose of this study.

371 In addition to their sources, several other factors specific to GWD data exacerbate
 372 the data imbalance problem. In the case of each GCM column with non-zero GW activity,
 373 momentum fluxes are generally concentrated at a few critical height levels rather than being
 374 smoothly distributed throughout the entire column. This further restricts the effective
 375 occurrence frequency of non-zero values. Moreover, GWs exhibit significant intermittency,
 376 where a small portion of large-amplitude GWs often dominates the morphology of the
 377 observed global GW momentum flux distribution (Hertzog et al., 2012; Geller et al., 2013).
 378 Therefore, it is crucial for NNs to not only accurately identify the columns that produce
 379 GWDs but also to effectively learn and recognize rare and extreme GWDs.

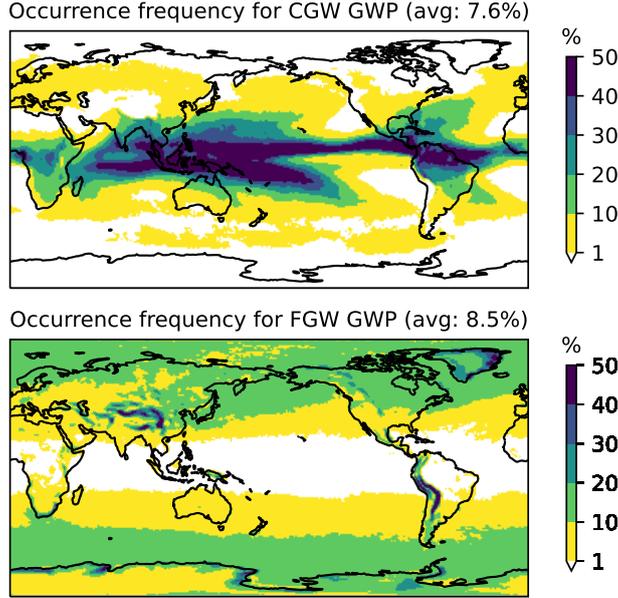


Figure 3. Distribution of occurrence frequency for CGWs (top) and FGWs (bottom) in the WACCM pre-industrial control simulations, based on the average of the 6-year training dataset.

380 Given the complexity of the GWD dataset, different normalization methods are con-
 381 sidered in this study. The first method, dubbed “NORM1”, is the typical normalization
 382 used in ML practices, which calculates elemental means and standard deviations for each
 383 feature (i.e., input variable at a given model level) and normalizes both inputs and outputs
 384 by these values (e.g., Espinosa et al. (2022)). With this approach, the same relative changes
 385 in wind at each level are treated equally in the input. The loss function in Equation (2) also
 386 penalizes the same relative error in GWD at each level equally. The second method, referred
 387 to as “NORM2“, is designed with the physics of GWD in mind. For the velocity inputs
 388 (u, v) and the tendency outputs (GWD), each column is normalized by one single value,
 389 which is the largest standard deviation from all model levels. Additionally, the mean values
 390 for these variables, are retained (e.g., $u_{norm2}(x, y, z, t) = u(x, y, z, t) / \max(std(u))$). Un-
 391 like NORM1, the original wind profile’s structure is preserved in NORM2, and large GWD
 392 values at certain heights maintain a relatively larger value after this normalization. For all
 393 other input variables, NORM2 is identical to NORM1. Compared to NORM1, NORM2
 394 places more emphasis on large GWD values and penalizes the NN more for missing these
 395 significant tendencies. These two normalization methods are also employed in Chantry et al.
 396 (2021), who found similar performance from these methods with the non-orographic GWPs.

397 Figure 4 shows the performance of the emulations for CGWs with the two normal-
 398 ization methods. When employing NORM1, the conventional approach seen in prior ML
 399 practices, and also our initial attempts, the emulator’s performance is poor. Although the
 400 NN demonstrates some skill, its predictions tend to cluster around zero. However, when the
 401 second normalization method (NORM2) is employed, the emulation results show significant

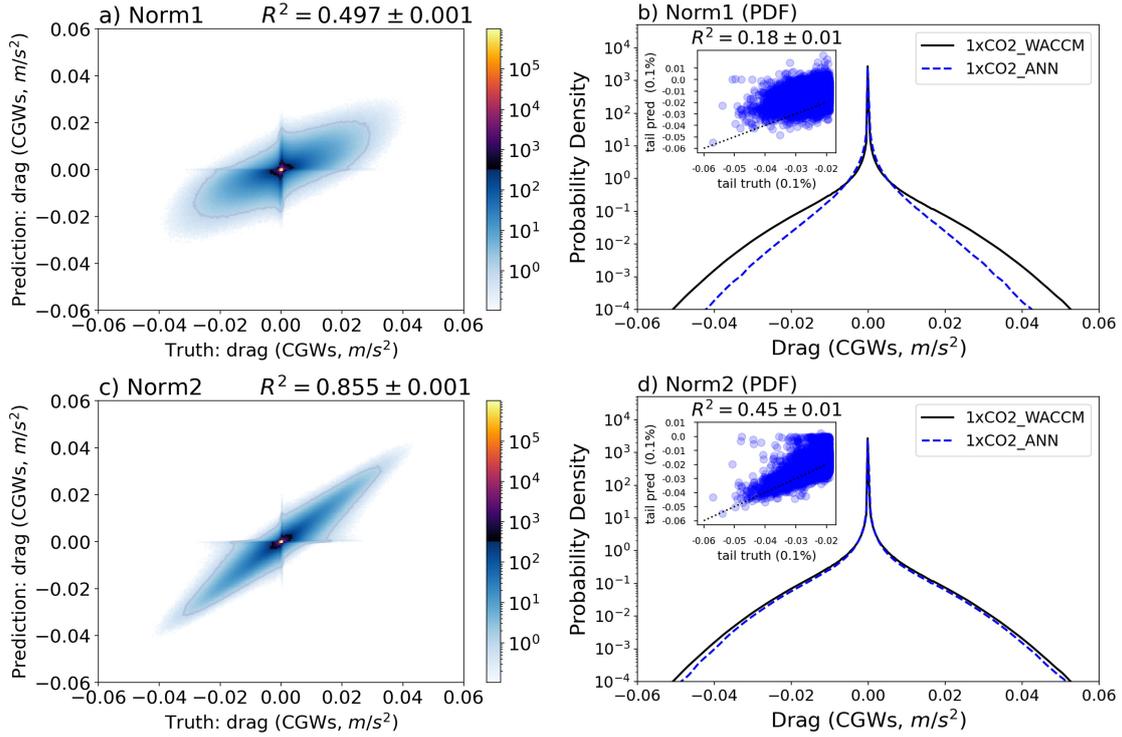


Figure 4. Data imbalance for GWD due to CGWs and the emulation results with two different normalization methods. a) A 2D histogram displaying the emulated GWD due to CGWs and the truth, with the training dataset normalized using NORM1; b) Distribution of the original convective GWD (black line) and the predicted values (blue line) with NORM2. The scatter plot in the corner represents the tail part only, including points with the top 0.1% amplitudes; c) Similar to a), but for NORM2; d) Similar to b), but for NORM2. The R^2 uncertainty range is estimated by dividing the test data into 10 segments, calculating the metric for each segment, and then computing the standard deviation (STD).

402 improvement, in contrast to the findings of Chantry et al. (2021). We attribute this improve-
 403 ment to the more pronounced data imbalance in our dataset, and it is likely a consequence
 404 of NORM2’s emphasis on modeling the large GWD values. Nonetheless, emulating the tail
 405 of the probability density function (PDF) (rare events) remains poor, as evidenced by the
 406 tails in Figure 4c, primarily due to the predominance of zero GWD columns in the training
 407 dataset. To more effectively address the data imbalance issue in these regression tasks, we
 408 further propose two approaches here:

- 409 1. Resampling the data (ReSAM): In this approach, we limit the number of training
 410 sample pairs with zero GWD to be equal to the number of samples with non-zero
 411 GWD. This significantly reduces the number of columns with zero GWD, thus mit-
 412 igating the data imbalance issue. Additionally, this sub-sampling reduces the total
 413 size of the training dataset, which, in turn, enhances the training speed (approx-
 414 imately sevenfold). While resampling methods have been well-established in the ML

415 literature, they have mainly been used for classification problems. Their application
 416 to regression problems in climate research has not been extensively explored.

417 2. Weighted loss function (WeLoss): Instead of assigning the same weight to all sample
 418 pairs in the loss function, we modify the weight for each column based on the PDF
 419 of its maximum GWD amplitude. This adjustment allows us to re-formulate the loss
 420 function defined in Equation (2) as

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{W}_i \{ \mathbf{NN}(x_i, \Theta) - y_i \} \right\|_2^2 \quad (4)$$

where

$$\mathbf{W}_i = \frac{1}{PDF(max(|y_i(z)|))} \quad (5)$$

421 Note that, in practice, we lack knowledge of the precise PDF for the maximum GWD
 422 within each column. Therefore, we employ a histogram with 20 bins as an alterna-
 423 tive. Given the fact that large-amplitude GW events are rare, the WeLoss approach
 424 incentivizes the NN to prioritize these significant events.

425 When we apply the ReSAM approach to balance the training dataset (after normal-
 426 ization with NORM1 or NORM2), the emulation results significantly improve, as shown in
 427 Figure 5. In fact, when considering the R-squared value between the NN prediction and the
 428 ground truth, the ReSAM approach with NORM2 yields the best results. However, as the
 429 training dataset is still predominantly composed of zeros and small GWD values due to the
 430 intermittence of the GWs, examining the emulation results for only large amplitude GW
 431 events (e.g., the top 0.1% in Figure 5d) reveals less satisfactory performance ($R^2 = 0.72$).
 432 Regarding the WeLoss approach, it has a more limited impact on improving the R-squared
 433 value of the emulation (as shown in Figure 5e). However, it proves valuable in capturing
 434 the tails of the PDF and, thus, rare events (as depicted in Figure 5f). Moreover, as ReSAM
 435 and WeLoss represent distinct operations, they can be effectively combined when construct-
 436 ing a NN. The result of this combined approach for emulating the CGWs can be found in
 437 Figures 5g and 5h. While the R-squared value for the entire distribution only marginally
 438 changes (0.925 vs. 0.931 with ReSAM only), the performance of the emulation for the tail
 439 part has been improved (R^2 increased to 0.77).

440 Similarly, Figure 6 presents the offline emulation results for the FGWs. The conclusions
 441 drawn for CGWs generally hold true. However, data imbalance in FGWs is less pronounced
 442 compared to CGWs, which simplifies the task of emulating FGWs. Even without any
 443 resampling or changes to the normalization or (see Figure 6a), we achieve reasonable emu-
 444 lation results ($R^2 = 0.9$). One contributing factor is the wider spatial distribution of FGWs
 445 compared to CGWs (refer to Figure 3). Additionally, the source of FGWs (frontogenesis
 446 function) in WACCM exhibits a much more continuous nature compared to precipitation
 447 and diabatic heating. As the data imbalance issue is less severe for FGWs, the performance
 448 with different normalization methods becomes more similar, echoing findings from Chantry
 449 et al. (2021) who emulated non-orographic GWs (including convective and frontal GWs)
 450 together.

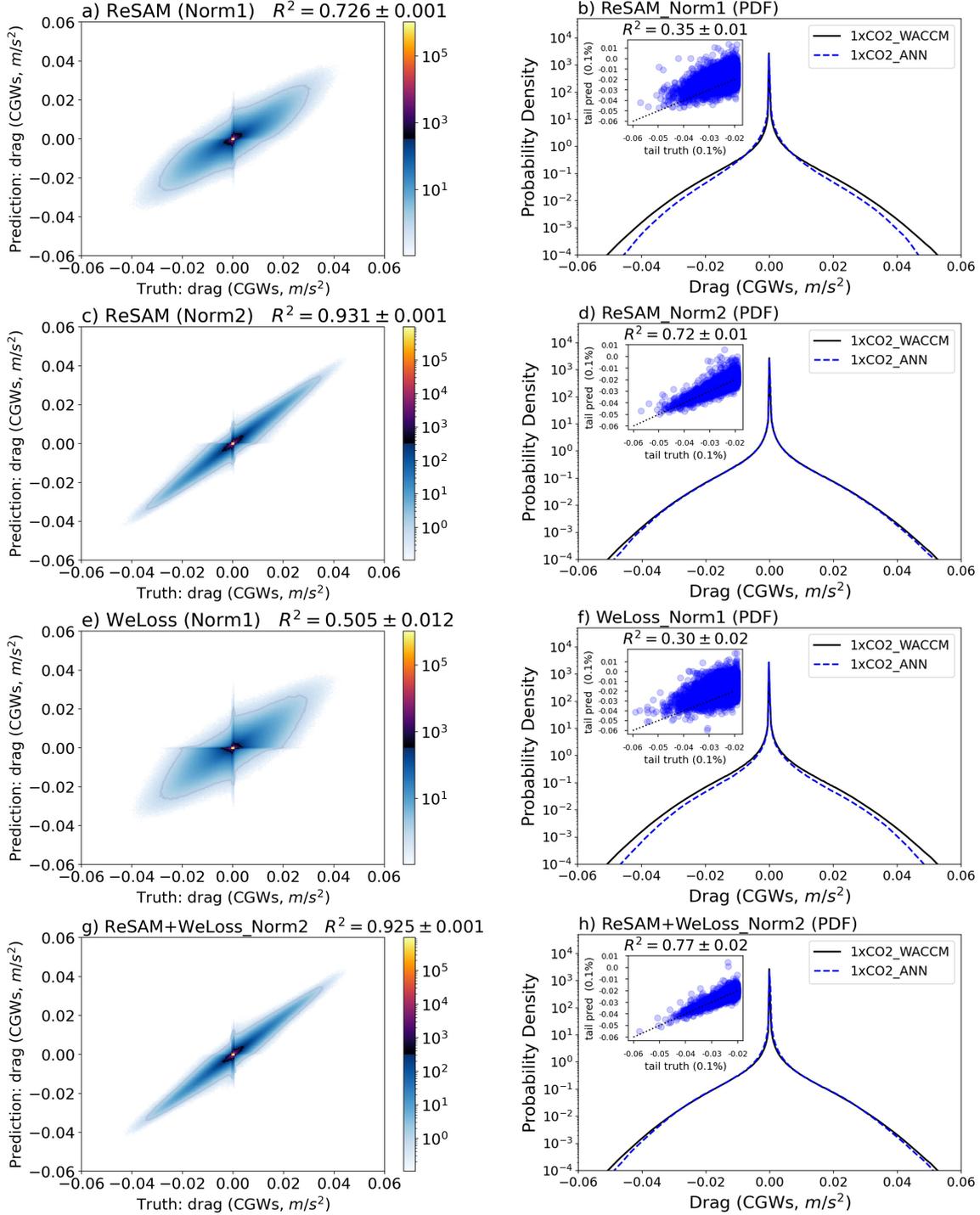


Figure 5. Similar to Figure 4, but for CGWs with the proposed ReSAM and WeLoss methods. a) A 2D histogram for the emulation with resampled data (ReSAM) after using Norm1; b) Distribution of the emulated GWD due to CGWs similar to Figure 4b, but with ReSAM applied; c) Similar to a), with training data normalized using Norm2; d) Similar to b), with training data normalized using Norm2; e) Similar to a), but with the WeLoss approach; f) Similar to b), but with the WeLoss approach; g) Similar to c), after applying both ReSAM and WeLoss methods together; h) Similar to d), after applying both ReSAM and WeLoss methods.

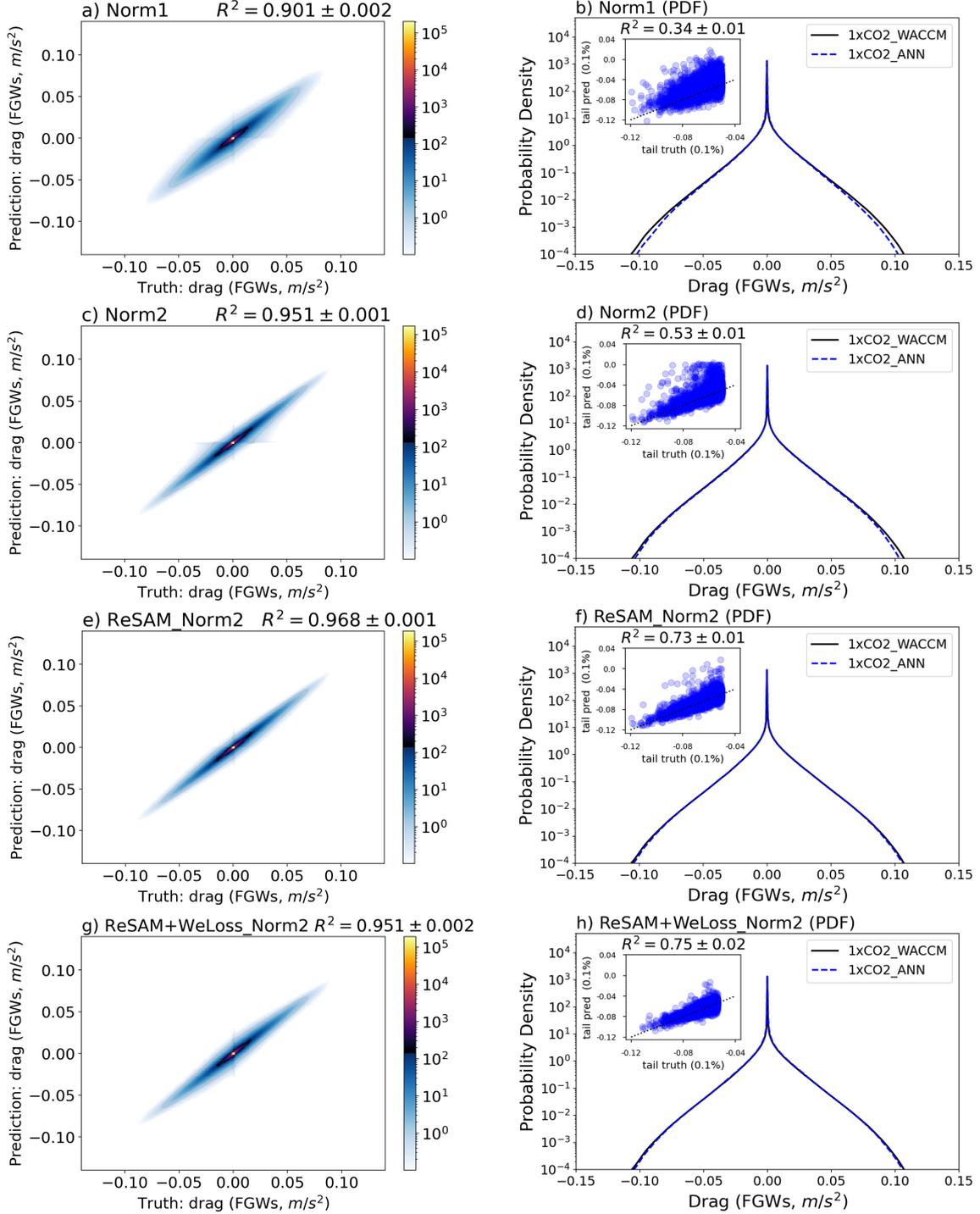


Figure 6. Similar to Figure 5, except for FGWs. a) A 2D histogram for the emulation using Norm1, without ReSAM or WeLoss; b) Distribution of the GWD due to FGWs with NORM1, similar to Figure 4b; c) Similar to a), with training data normalized using Norm2; d) Similar to b), with training data normalized using Norm2; e) Similar to c), but with the ReSAM approach; f) Similar to d), but with the ReSAM approach; g) Similar to Figure 5g, applying both ReSAM and WeLoss methods to the FGWs; h) Similar to Figure 5h, applying both ReSAM and WeLoss methods to the FGWs.

451 In summary, data imbalance can pose challenges when learning from data that closely
 452 resembles real-world data (further discussed in the subsequent section on emulating OGWs).
 453 Proper resampling techniques can significantly enhance the NNs' performance by improving
 454 dataset balance. Furthermore, modifying the loss function to penalize the NNs more for
 455 missing extreme values can further improve performance at the tails of the PDF. For the
 456 remainder of the paper, unless otherwise specified, we continue to employ the ReSAM
 457 approach and the standard loss function with NORM2 unless stated otherwise.

458 3.2 Uncertainty Quantification

459 As outlined in subsection 2.2.2, we employ three different methods (i.e., BNN, DNN,
 460 and VAE) to quantify the uncertainty of predictions during inference (testing). For this
 461 purpose, an ensemble of 1000 members is generated by running each UQ-equipped NN 1000
 462 times for each input from the testing set. Figure 7 presents sample profiles of zonal GWD
 463 derived from the deterministic NN (ANN) and the three UQ-equipped NNs, alongside the
 464 true GWD profiles from WACCM. Note that these examples have not been used in the
 465 training or validation process. It is evident from the figure that all three UQ-equipped
 466 NNs show reasonable skill in predicting the complex profiles of GWD due to CGWs and
 467 FGWs (also reflected in R-squared in Table 1), albeit with a slight decrease in accuracy
 468 compared to ANN. As discussed earlier, a valuable uncertainty estimate should correspond
 469 closely with the NN's test accuracy, providing insights into when to trust the NN's pre-
 470 diction during inference. Such a relationship can be seen in a few randomly chosen GWD
 471 profiles that's shown in Figure 7. In each pair of CGW and FGW profiles, the left column
 472 shows the estimated uncertainty is also low when the prediction error is low, indicating the
 473 NN's confidence in its accurate predictions. In contrast, the right column, which generally
 474 represents more complex profiles, exhibits the NN's less accurate predictions, and increased
 475 uncertainty, highlighted by the wider confidence intervals.

476 While Figure 7 demonstrates the performance of the UQ methods for just a few GWD
 477 profiles, the spread-skill plots shown in Figure 8 offer a broader perspective based on 60,000
 478 profiles, following the calculations detailed in Appendix C. It is evident from the plots
 479 that all three UQ methods produce reasonably informative uncertainty estimates, as their
 480 curves closely align with the 1-to-1 line. In the case of CGWs, all data points are above
 481 the 1:1 line, indicating a slight overconfidence (underdispersiveness) across all three UQ
 482 methods, with the DNN being slightly closer to the 1-to-1 line. For the FGWs, the DNN
 483 demonstrates slightly better performance, although it marginally drops below the 1-to-1
 484 line in the first few bins, indicating a slight underconfidence. Notably, it can be seen from
 485 the spread frequency inset that the vast majority of the data points are within the first few
 486 bins, for which both spread and skill values are small, and they are generally closer to the
 487 1-to-1 line.

488 It should also be noted that for the large values of model spread (\overline{SD}), there is only a
 489 very limited number of data points, as is evident from the inset histograms. Consequently,

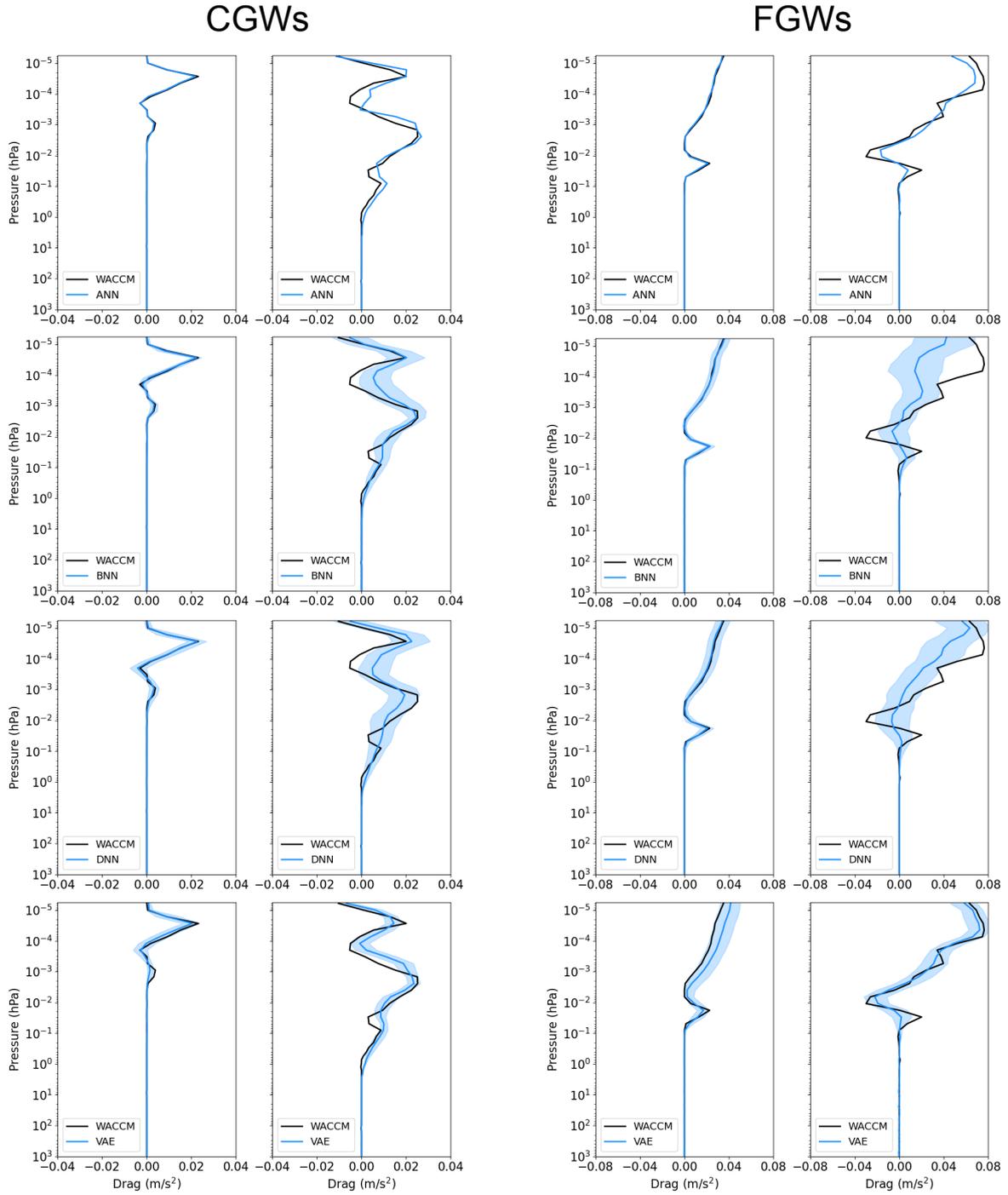


Figure 7. Sample profiles of zonal GWD as predicted by various NNs, as indicated. The true profile is shown by the black line, while the blue solid line represents the mean of 1000 ensemble members. The shaded region indicates the 95% confidence interval. In each pair of CGWs and FGWs profiles, the left column provides examples with low estimated uncertainty, corresponding to instances of low error. Conversely, the right column illustrates cases with high uncertainty when the error is high.

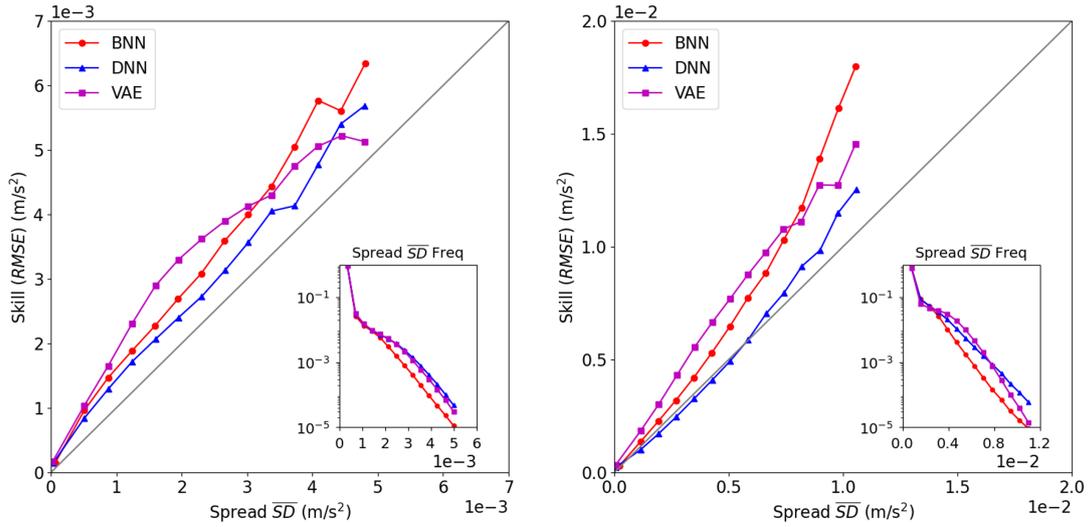


Figure 8. Spread-skill plot for GWD due to (left) CGWs, and (right) FGWs. The diagonal 1:1 line represents the perfect spread-skill line. Points above (below) this line correspond to spread values where the model is overconfident (underconfident). The inset histogram shows how often each spread value occurs. See Appendix C for a detailed discussion on the calculations of the spread-skill plot.

490 the standard deviation (STD) can become a misleading measure of spread because of the
 491 non-normal distributions.

492 To summarize the quality of the spread-skill plots for the three UQ methods, we explore
 493 the metrics introduced in subsection 2.2.2 and Appendix C (see Table 1). The R-squared
 494 value for the ensemble mean prediction is also given to show the accuracy of each UQ
 495 method. Based on SSREL, whose ideal value is zero, BNN shows the best performance for
 496 both CGWs and FGWs. However, if we check SSRAT, where 1 is the optimal number, DNN
 497 is the best among these three methods. This discrepancy can be explained by a closer look
 498 at the Equations (C2) and (C3). SSREL, which is a bin-weighted average difference, is most
 499 sensitive to the performance of the NN in the first bin, where the vast majority of the data
 500 points are located (see the inset histograms in Figure 8), while SSRAT is more influenced by
 501 larger values of spread and skill. Accordingly, the VAE shows the highest values of SSREL,
 502 which is indicative of its sub-optimal performance in the first bin, where there are small
 503 values of spread and skill.

504 In the results presented in Figure 8 and Table 1, each height level of a GWD profile is
 505 considered as an individual sample. A zonal GWD profile, with its 70 vertical levels, thus
 506 constitutes 70 distinct samples. While analyzing these samples offers insights into the NN's
 507 overall performance by averaging statistics across numerous profiles, our primary interest is
 508 often in the uncertainty associated with an individual GWD profile. This uncertainty can
 509 then aid in determining whether to trust/use the NN's prediction for that particular GWD

510 profile. Therefore, we use Equation (C4) to assess the relationship between uncertainty and
 511 test accuracy for each GWD profile. Furthermore, to estimate uncertainty, here we use the
 512 interquartile range (IQR) to reduce the influence of outliers.

513 Figure 9 shows the Gaussian kernel density of spread against RMSE for all 60,000
 514 profiles, as indicated by the color shading. The x -axis represents the IQR of each GWD
 515 profile, while the y -axis denotes its corresponding RMSE. A strong correlation between
 516 the two is observed across all three UQ methods. Consequently, GWD profiles with larger
 517 uncertainties often coincide with larger errors. Figure 9 also shows a close similarity between
 518 BNN and DNN. In contrast, VAE tends to provide marginally larger uncertainties, especially
 519 for FGWs. This is consistent with VAE’s slightly reduced accuracy as indicated in Table
 520 1. Overall, given the monotonic relationship between the uncertainty and test error, these
 521 results show that all three UQ methods provide useful and informative uncertainty for with-
 522 distribution test samples. A user can set a threshold on uncertainty based on their tolerance
 523 for error (RMSE) and decide whether they trust the NN for a given input sample.

524 The results presented so far show the performance of the UQ methods based on the
 525 testing data, i.e., data from the current climate. However, the effective performance of UQ
 526 methods can also be tested (perhaps more meaningfully) on OOD data, e.g., data from a
 527 warmer climate. This is particularly relevant for climate change studies. Accordingly, we
 528 evaluate the performance of these trained NNs with input data from the future climate, as
 529 depicted by the black lines in Figure 9. For FGWs, the spread-skill relationship remains
 530 largely similar, especially for BNN and DNN. This suggests that, based on their uncer-
 531 tainties, we can still reliably estimate the error in the NN predictions for FGWs for the
 532 warming climate. A similar pattern is observed for the VAE, though it exhibits increased
 533 uncertainties and higher errors with OOD data. As shown in a later section, for FGWs, the
 534 NNs generalize to the warmer climate without any further effort.

535 In contrast, for CGWs, given the same level of uncertainty, the error in NN predictions
 536 increases significantly for the OOD data compared to that from the current climate, which
 537 means the spread-skill relationship, especially for the BNN and DNN, fails to generalize to
 538 the OOD data. From this perspective, VAE performs better, showing that for the same
 539 level of uncertainty, the increase in error is not as substantial as in BNN and DNN. The
 540 VAE also yields considerably higher uncertainty estimates for future climate, which may aid
 541 in the detection of OOD data. The observed discrepancies in the performance of the NNs
 542 for CGWs and FGWs hint at different levels of their generalizability, a topic we will delve
 543 into more deeply in the following subsection.

544 In summary, while the three UQ methods provide credible and valuable uncertainty
 545 estimates for the current climate, the BNN and DNN are confidently wrong in estimating
 546 CGWs in a warmer climate although VAE shows some promising results. This problem is
 547 common among various UQ techniques as pointed out in the ML literature: they frequently
 548 show overconfidence when assessed with OOD data (e.g., Ovadia et al., 2019). The optimal
 549 UQ method selection depends on the specific metric of interest and the intended application.
 550 While BNN is more broadly used in the literature and gives the best accuracy, DNN could

Table 1. Evaluation scores for the three UQ methods. See Section 2 for more details.

	CGWs			FGWs		
	BNN	DNN	VAE	BNN	DNN	VAE
SSREL (1e-4)	1.29	1.48	2.14	1.20	1.69	5.21
SSRAT	0.73	0.82	0.72	0.69	0.93	0.69
R-squared	0.90	0.86	0.87	0.94	0.92	0.89

551 achieve similar performance and is often more practical given its simplicity. On the other
552 hand, VAE seems to perform better when applied to OOD data, at least in the one test case
553 here. These observations warrant further research in the future using multiple test cases
554 and climate-relevant applications. We also note here that each method has multiple tuning
555 hyperparameters to optimize its uncertainty quantification. Consequently, the discrepancies
556 among the three methods could potentially be mitigated with proper hyperparameter tuning
557 (as discussed in Appendix B).

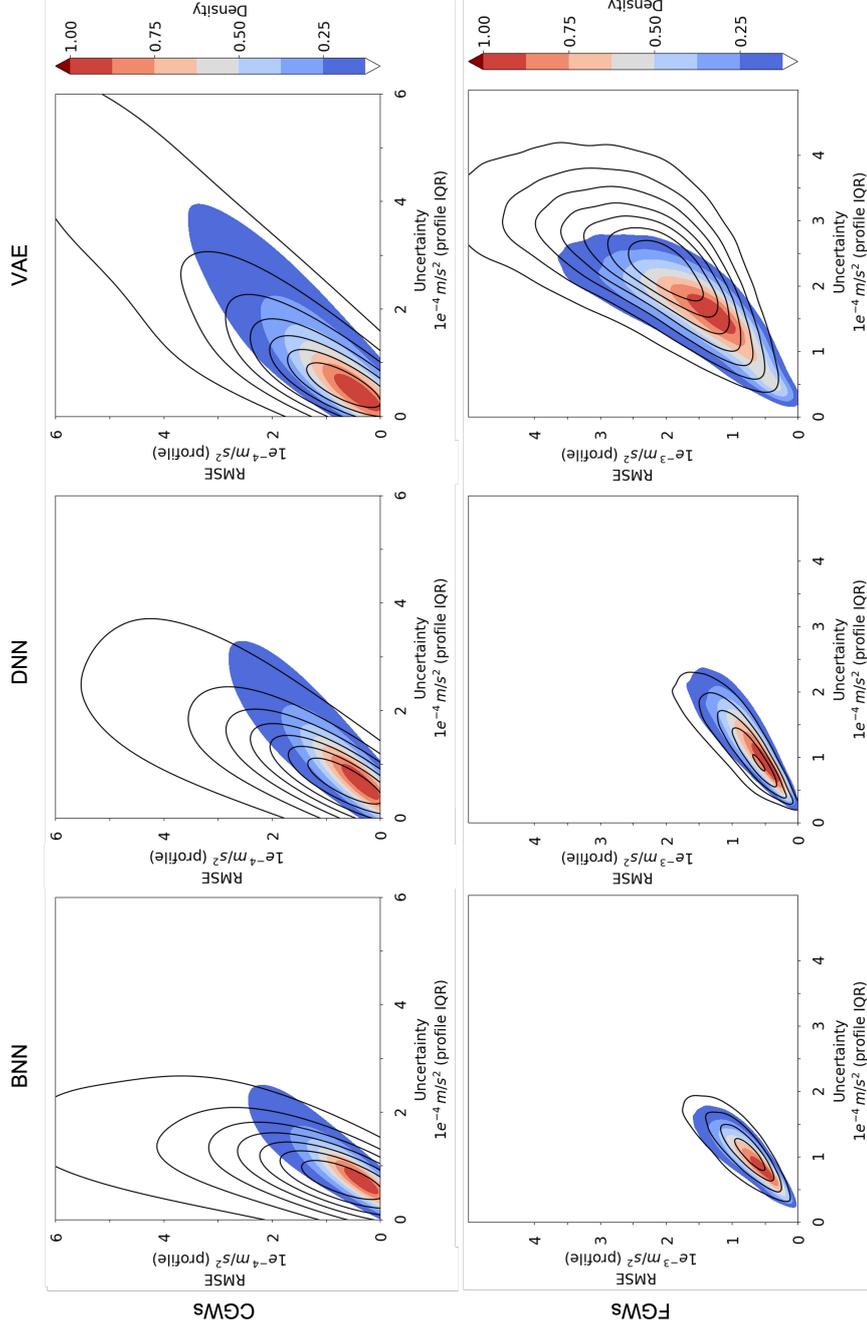


Figure 9. Gaussian kernel density for spread, defined as interquartile range (IQR) versus RMSE for (top) CGWs, and (bottom) FGWs, when validated against out-of-sample data from the current climate (color shading) and out-of-distribution data from the future, warmer climate (black lines). Results are shown for BNN (left), DNN (middle), and VAE (right). All NNs are exclusively trained on data from the current (control) climate.

3.3 Out-of-distribution (OOD) Generalization via Transfer Learning

As previously discussed, the GWP schemes in WACCM are coupled to their sources, which might change in a warmer climate. Specifically, under $4\times\text{CO}_2$ forcing, we expect changes in both the amplitude and the phase speed distribution of GWs, in particular for the CGWs, due to their built-in sensitivities to changes in the convection. Consequently, the physics scheme in WACCM produces slightly stronger GWD for CGWs, especially in the tail of the distribution. This intensified GWD results in a shorter quasi-biennial oscillation (QBO) period in WACCM. However, it is important to recognize that the response of the QBO to climate change differs across various general circulation models (Richter et al., 2022).

The intensification of the CGWs in future climate simulations presents an opportunity to study how NNs handle the OOD data. Our findings in the UQ section already suggest increased prediction errors when testing NNs with OOD data, which raises concerns about their applicability in climate change studies. To more thoroughly investigate this issue, we conduct additional evaluations on our ANNs, by applying them to data samples from future climate simulations, as illustrated in Figure 10. It is clear that the ANN for the CGWs does not generalize well, evidenced by a decrease in R^2 from 0.93 to 0.79. The ANN particularly struggles to capture the increase in GWD in the tail, with R^2 for the tails decreasing from 0.72 to 0.36. As a result, it seems unlikely that this emulator will accurately reproduce changes in the circulation under different climate conditions, such as the shorter QBO period resulting from future warming in WACCM.

In contrast to CGWs, the amplitude of FGWs shows a less marked increase in the future climate, and their PDF distribution closely resembles that of the control simulations. As a result, the ANN demonstrates better generalizability for FGWs when it is tested against future climate data, as seen in Figure 10d. There is only a slight decrease in the ANN's performance, with R^2 dropping from 0.97 to 0.95.

Two factors can contribute to the considerable OOD generalization errors in an NN when applied across two distinct systems. First, the input-output relationship might vary between the two systems. Second, the input variables in the new system could originate from a distribution different from that of the original system (regardless of whether the input-output relationship remains the same or changes). The former is hard to quantify in a high-dimensional dataset. The latter can be quantified using similarity distances. To help us better understand these differences between the OOD generalizability of CGWs and FGWs, we assess the similarity between their input and output data distributions from control and future climate simulations using the Mahalanobis distance (D). The Mahalanobis distance is a measure of the distance between a data point and a distribution (Ling & Templeton, 2015). Specifically, it is a multi-dimensional generalization of the idea of measuring how many standard deviations away a point is from the mean of the distribution. The application of Mahalanobis distance in understanding the source of OOD generalization errors in data-driven parameterization was previously demonstrated in Guan et al. (2022) for a simple turbulent system.

Table 2. Change of Mahalanobis distance based on the ratio of the average distance of the points that are more than 3 standard deviations away from the mean. The choice of the variables here is based on Appendix A, showing u, v, T , and source function contain most of the information needed for the NN.

Variables	u	v	T	Source (diabatic heating for CGWs, frontogenesis for FGWs)	Zonal drag	Meridional drag
Distance (Convection)	1.03	1.00	1.19	3.62	1.42	1.44
Distance (Front)	1.03	0.96	1.50	1.10	1.00	1.00

599 To use the Mahalanobis distance, we first calculate the mean and covariance matrix of
600 the training data from the control run. We then analyze the distribution of Mahalanobis
601 distances in this training data, setting a baseline value, referred to as D_{ctrl} . This baseline is
602 the average distance for data points that deviate by more than 3 standard deviations from
603 the mean. This choice aims to focus on outliers for which extrapolation is more challenging.
604 Following this, we apply the same process to the data points in the future climate dataset,
605 denoted as D_{warm} . Table 2 presents the ratio of D_{warm} for the warming scenario to D_{ctrl}
606 for the control scenario for selected variables. When this ratio is close to 1.0, it suggests
607 minimal changes in this variable’s distribution under a warming scenario. Note that the
608 NNs trained based only on these variables demonstrate performance comparable to NNs
609 trained on all variables (not shown), which is why we only focus on these few key variables.

610 The results reveal that among the various variables significantly contributing to the
611 emulation of CGWs, diabatic heating (source of CGWs) is the sole variable that exhibits
612 substantial changes from the control to the warming scenario. Conversely, changes in vari-
613 ables used to emulate FGWs are considerably smaller. This outcome suggests that the likely
614 reason for the better generalizability of FGWs is that the input distribution remains almost
615 unchanged (and the input-output relationship, which is the physics scheme, remains the
616 same too).

617 To improve the generalizability of the emulator for CGWs, we explore TL, a technique
618 introduced earlier and proven to be a powerful tool for improving the OOD generalizability
619 of data-driven parameterization in canonical turbulent flows (e.g., Guan et al., 2022; Subel
620 et al., 2023). Rather than re-training the entire NN for future climate scenarios, we only re-
621 train, following Subel et al. (2023), just a portion of the NN, thereby requiring only a small
622 fraction of the data. Figure 10e showcases the emulation results after only re-training the
623 first hidden layer of ANN using data from the first month of the WACCM simulation in the
624 $4\times\text{CO}_2$ scenario, which amounts to approximately 1% of the original training dataset. After
625 applying TL, the performance of the emulator in the warming scenario significantly improves,
626 with R^2 rising from 0.79 to 0.91, nearly matching its performance in the control simulations
627 ($R^2 = 0.93$). However, the improvement in the PDF tails is less pronounced, showing
628 only a modest increase in R^2 from 0.36 to 0.51. This is likely due to the limited number
629 of large-amplitude GW events within the one-month period. Instead of using more data

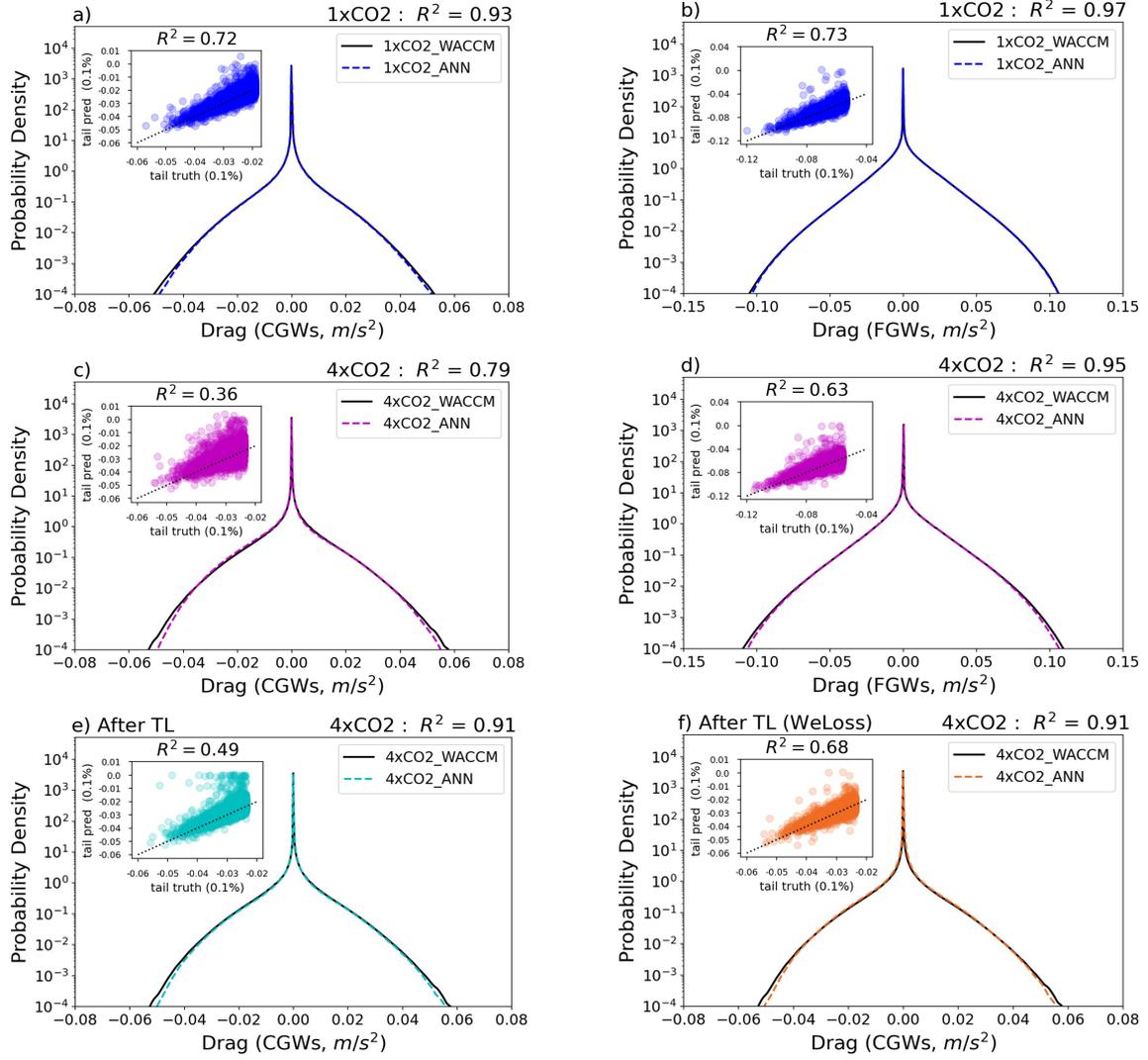


Figure 10. NN performance for pre-industrial and warming scenarios for different sources (a,c,e,f: CGWs ; b,d: FGWs). **a)** PDF of GWD due to CGWs in WACCM simulation and the predicted CGWs using NN emulator, scatter plot shows points for the tail part only. **b)** same as (a), but for FGWs. **c)** same as (a), but for the warming scenario, **d)** same as (b) but for the warming scenario. **e)** same as (c) but after applying transfer learning to the first hidden layer of the NN with 1-month WACCM simulation data under warming scenario ($\sim 1\%$ of the size of the training data) **f)** same as (e) but with the weighted loss function used when we conduct transfer learning (WeLoss).

630 from the future climate (which is challenging to obtain in a realistic situation), we leverage
 631 the WeLoss approach, described earlier, during re-training. This modification results in a
 632 significant improvement in the tail, with R^2 increasing from 0.51 to 0.68. Note that this
 633 improvement in the tail is critical, as inadequate learning of these rare but large-amplitude
 634 GWDs can result in significant errors and instabilities.

635 We would like to point out that during the TL experiments, we have examined the
 636 effects of re-training each individual hidden layer of the NN. Our findings indicate that
 637 re-training the first layer yields the best results, which aligns with the findings in Subel et
 638 al. (2023). Re-training the last layer only brings marginal improvements to the NN (not
 639 shown). Notably, our experiments involving re-training the first two layers did not result
 640 in further performance enhancements, suggesting that the number of neurons is not the
 641 primary factor contributing to the varied performance observed when re-training different
 642 layers.

643 Similar results regarding TL are also observed with other NNs used in this study. For
 644 instance, Figure 11 presents the same plot as Figure 10, but for the BNN. It is evident that
 645 BNN also struggles with generalization to OOD data, as could also be interpreted based
 646 on the results presented in section 3.2. It is also the case for DNN and VAE (not shown).
 647 Overall, when these NNs are tested against the $4\times\text{CO}_2$ future climate data, their accuracy
 648 is not better than the deterministic ANN. However, methods with UQ, especially the VAE
 649 (see Figure 9), could potentially indicate the increased uncertainty when testing with input
 650 data from the $4\times\text{CO}_2$ integration. These results underscore the necessity of re-training the
 651 NNs using TL.

652 **4 Emulation of Orographic GWs (OGWs)**

653 Similar to Chantry et al. (2021), our initial attempts to emulate OGWs did not succeed,
 654 primarily due to the presence of a pronounced data imbalance. Notably, the physics-based
 655 scheme responsible for OGW generation operates exclusively over terrestrial regions. How-
 656 ever, it is surprising that the issue of data imbalance continues to persist, even when we
 657 limit our NN training and testing exclusively to columns located over land (Figure 12a).
 658 Still, the emulated OGW drag often remains close to zero and completely fails to predict
 659 the rare events (Figure 12b), which poses a considerable hurdle for the emulator’s perfor-
 660 mance. Further investigations reveal that the key to this problem lies in the highly localized
 661 nature of orographic GWD, where significant drag is observed only at a handful of specific
 662 locations. Furthermore, even within these limited regions, GWD exhibits a significant in-
 663 termittent behavior. To help our understanding, we also conducted a K -means clustering
 664 analysis, categorizing GWD data for all land-based columns (Table 3). Among the 6 clus-
 665 ters, cluster 4 accounts for a staggering 97.51% of the dataset. Remarkably, all samples
 666 within this cluster exhibit exceptionally weak orographic GWD, as evidenced by the cluster
 667 center’s maximum GWD amplitude, which is two orders of magnitude smaller than that of
 668 other clusters.

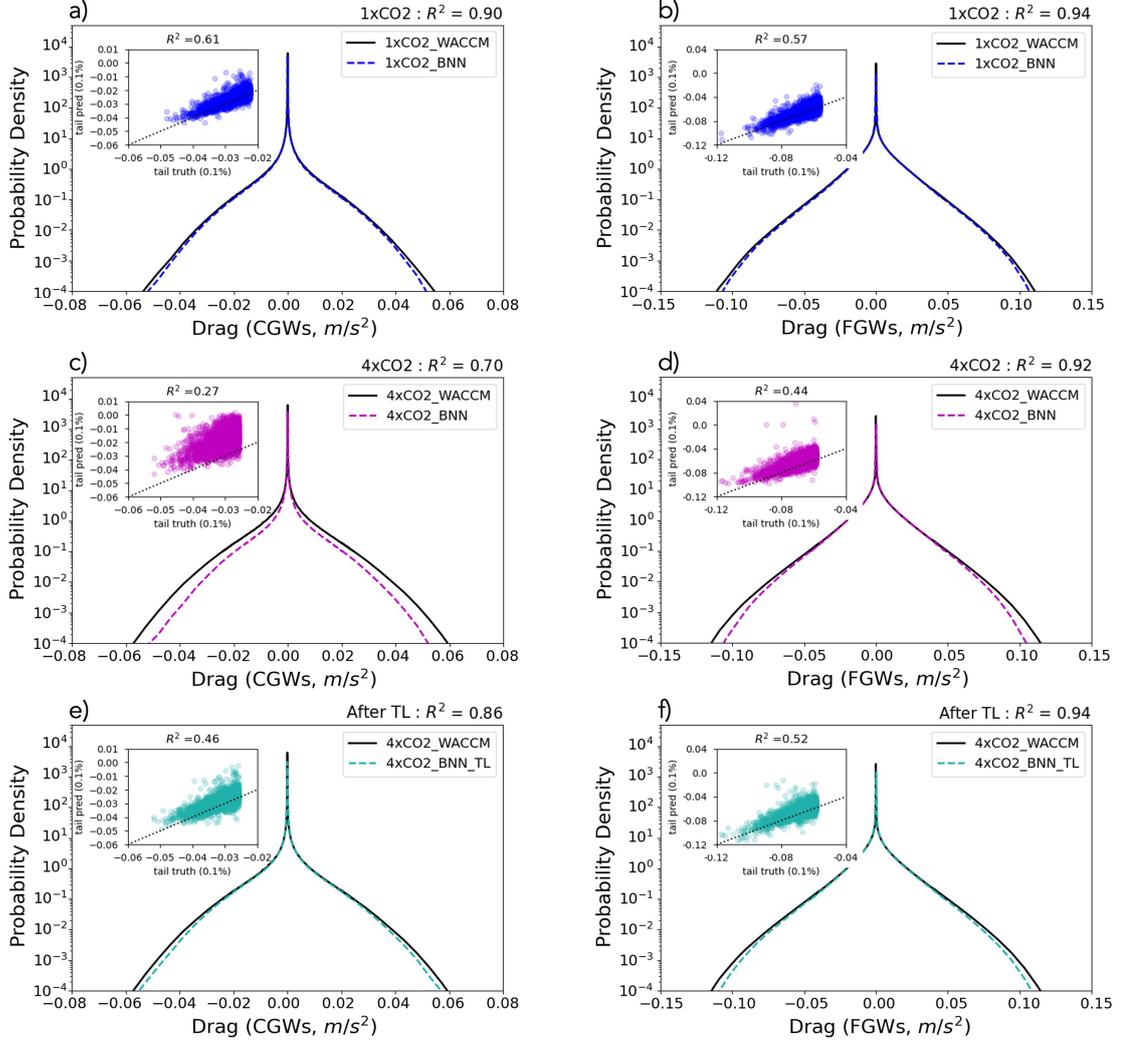


Figure 11. Panels (a) to (e) are the same as those in Figure 10 but for BNN. Panel (f) shows emulation for FGWs under warming scenario after applying transfer learning to the first hidden layer of the NN.

Table 3. Clustering analysis for OGWs. Analysis is done for all columns over land in the training data.

Cluster	Frequency (%) in the training data	Maximum GWD amplitude of cluster center
c1	0.18	8.7 e-3 m/s ²
c2	0.13	4.4 e-3 m/s ²
c3	0.93	3.6 e-3 m/s ²
c4	97.51	2.8 e-5 m/s ²
c5	0.15	2.1 e-3 m/s ²
c6	1.10	4.3 e-3 m/s ²

669 To overcome this persistent data imbalance in the OGWs, we first separate all columns
670 over land into large-drag columns (with column maximum greater than one STD of all
671 GWD from OGWs) and small-drag columns. We then perform subsampling on the latter
672 group only to create a more balanced dataset. To improve NN training, we also include all
673 columns from the 6-year simulation to augment the sample size of the large-drag columns.
674 Figures 12c and 12d illustrate the performance after re-balancing the dataset. Notably, the
675 result represents a substantial improvement, evidenced by an R^2 increase from 0.29 to 0.80,
676 and also a significant improvement in the accuracy for rare events. While we acknowledge
677 that this skill remains lower than what is achieved for CGWs and FGWs, it already signifies a
678 reasonable NN. Furthermore, we posit that by incorporating additional training data (either
679 by extending the WACCM model integration or simply augmenting the data with OGWs
680 scheme only), we can further improve our emulation results. The possibility of achieving
681 superior emulation outcomes through the adoption of an alternative NN architecture is also
682 possible, although such exploration is beyond the scope of this paper.

683 5 Summary and Discussion

684 Through the emulation of complex GWPs in a state-of-the-art atmospheric model
685 (WACCM), we have elucidated and explored solutions for three critical challenges in the
686 development of ML-based data-driven SGS schemes for climate applications: data imbalance,
687 UQ, and OOD generalizability under different climates. A brief summary is provided
688 below:

- 689 1. In the presence of non-stationary, and highly imbalanced datasets, such as those en-
690 countered in WACCM, specialized approaches (e.g., resampling and weighted loss
691 function) are essential to enhance the performance of data-driven models. Through
692 resampling, we have successfully trained a robust NN emulator for OGWs, a challeng-
693 ing task as demonstrated in Chantry et al. (2021). The effectiveness of the trained
694 emulator is also significantly influenced by the choice of the loss function used dur-
695 ing training. In our case, while a weighted loss function (WeLoss) does not improve
696 the overall R^2 score, it yields significant improvements in the emulation results for

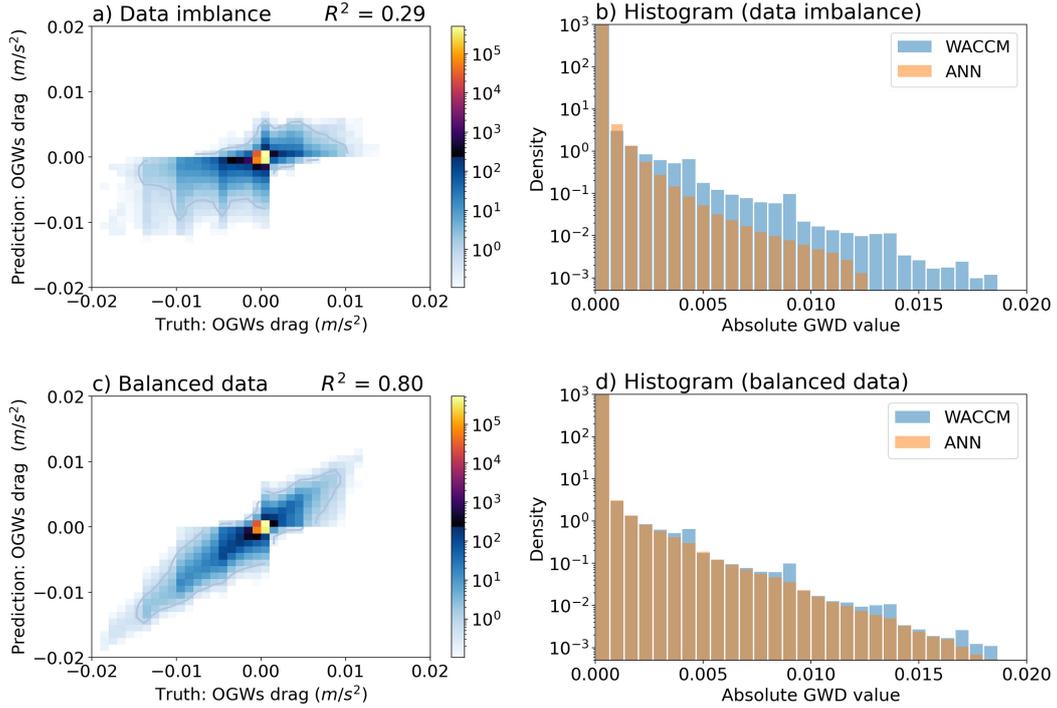


Figure 12. Performance of the emulator for OGWs when trained with all columns over land (panel a) & panel b)) and balanced training data with a balanced number of large-drag columns (column maximum > 1 STD of all GWD from OGWs) and small-drag columns (panel c) & panel d).

697 the PDF tails of the GWD. This finding aligns with those in Lopez-Gomez et al.
 698 (2022), where their custom loss function, tailored to emphasize extreme events, led
 699 to substantial improvements in predicting heatwaves.

700 2. All three UQ methods employed in this study provide reasonable uncertainty esti-
 701 mates for GWD prediction for the current climate. The spread-skill plots (refer to
 702 Figures 8 and 9) indicate that greater uncertainty corresponds to a larger prediction
 703 error. Yet, the reliability of UQ methods diminishes when they are challenged with
 704 OOD data. Both BNN and DNN used in this study tend to be overconfident in esti-
 705 mating CGWs in a warmer climate, thereby struggling to identify OOD samples.
 706 The VAE, on the other hand, yields rather promising results in providing useful UQ
 707 for OOD data. Given the variations in different methods, the metrics selected to
 708 assess the SGS model will play a significant role in determining the choice for the UQ
 709 methods. We also note that further optimization of tunable parameters within each
 710 UQ method could affect their performance (refer to Appendix C).

711 3. Our findings illustrate the challenges SGS schemes face in generalizing to OOD data
 712 and extrapolating to new climates. Nonetheless, the TL approach has proven highly
 713 effective in aiding an NN to extrapolate to different climates. For CGWs in WACCM,
 714 the physics-based scheme exhibits larger GWD under $4\times\text{CO}_2$ forcing, primarily due
 715 to an increase in diabatic heating from convection. With only one month of sim-

716 ulation data from this future warming scenario (representing approximately 1% of
 717 the original training dataset), we successfully reduce its OOD generalization error
 718 through re-training the first layer of the NN, following the findings of Subel et al.
 719 (2023). Additionally, we have illustrated the value of metrics like the Mahalanobis
 720 distance in assessing the potential OOD generalizability of NNs.

721 We would like to emphasize that these challenges are often intertwined. For instance,
 722 addressing data imbalance in CGWs is a prerequisite for obtaining an accurate NN model,
 723 which, in turn, impacts UQ and OOD generalizability assessments. Moreover, there exists
 724 a strong link between UQ and OOD generalizability evaluations: if the NN struggles with
 725 OOD generalization, performing poorly with OOD data, the reliability of UQ for such data
 726 (e.g., data from a warmer climate) also becomes questionable. This presents a substantial
 727 challenge for UQ methods, especially for climate change research where reliable UQ methods
 728 are crucial.

729 This study has primarily focused on offline skill assessment. We acknowledge that good
 730 offline performance (at least in terms of common metrics such as R^2) is not necessarily an
 731 indicator of stable and accurate online (coupled to climate model) performance (Ross et
 732 al., 2022; Guan et al., 2022), though more strict metrics such as R^2 of the PDF tails might
 733 better connect the offline and online performance (Pahlavan et al., 2023). However, for
 734 the purpose of this study, which is to provide a testbed to test ideas for data imbalance,
 735 UQ, and OOD generalization with transfer learning, the offline tests, particularly using
 736 the several metrics we have used, suffice. That said, the main reason that we have not
 737 provided online results is that coupling various complex NNs, with the same framework, to
 738 complex climate models (e.g., WACCM) without slowing down the model is a challenging
 739 and time-consuming task (Espinosa et al., 2022), and this is work in progress.

740 Emulating complex GWPs within the WACCM provided a unique opportunity to ad-
 741 dress three critical challenges in developing ML-based, data-driven SGS schemes for climate
 742 science applications. However, it is crucial to acknowledge that such emulated schemes
 743 essentially adopt the limitations inherent in the physics-based schemes. Addressing these
 744 limitations, the next step is to harness high-resolution data from GW-resolving simula-
 745 tions, which are carefully validated against observational data. A library of such high-
 746 resolution simulations, notably of convectively generated GWs using the Weather Research
 747 and Forecasting (WRF) model, is now established (Sun et al., 2023), alongside additional
 748 global high-resolution simulations (Wedi et al., 2020; Polichtchouk et al., 2023; Köhler et al.,
 749 2023). The next phase involves integrating the approaches outlined in this study with the
 750 data from these GW-resolving simulations to develop a stable, trustworthy, and generaliz-
 751 able data-driven GWP scheme. This scheme is then expected to overcome the limitations of
 752 physics-based GWPs and potentially incorporate features like the transient effect (Bölöni
 753 et al., 2021; Kim et al., 2021) and lateral propagation of GWs (e.g., Sato et al., 2009)—marking
 754 a significant advancement towards next-generation GWP schemes.

Appendix A Input/output variables for the physics-based GWP schemes and their emulators

We use the exact same inputs as those of each GWP scheme in the WACCM for the training of the NN-based emulator of that scheme. These inputs are listed in Table A1. As for the outputs, we only consider the zonal and meridional drag forcings. The GWPs in WACCM also estimate additional effects of the GWs that result in changes of temperature profile and vertical diffusion. These outputs are not considered in our emulations.

Table A1. List of the input and output variables for the NNs trained as emulators of the GWP schemes in WACCM. The numbers in parentheses in front of each variable are the number of vertical levels for that variable. Note that each input and output is a 1D column at a given latitude/longitude grid point. Diabatic heating in WACCM is provided by the cumulus scheme. The topography variables listed in the table are *mxdis* (height estimates for ridges), *hwdth* (width of ridges), *clngt* (length of ridges), *angll* (orientation of ridges), and *anixy* (anisotropy of ridges).

GWP	Input			Output
	pressure levels	surface level	forcing	
CGWs	$u(70)$,	lat (1),	diabatic heating (70)	zonal drag GWD _x (70), meridional drag GWD _y (70),
FGWs	$v(70)$, $T(70)$,	lon (1), $P_{surface}$ (1),	frontogenesis function (70)	
OGWs	$z(70)$, $\rho(71)$, Brunt–Väisälä frequency N (70), dry static energy DSE (70)		<i>mxdis</i> (16), <i>hwdth</i> (16), <i>clngt</i> (16), <i>angll</i> (16), <i>anixy</i> (16),	

From Table A1, one can guess that some input variables are correlated with each other. Consequently, it is plausible that the trained NNs may have spurious connections. Preliminary tests further support this notion, indicating that employing only u, v, T , and the forcing function as inputs yields comparable offline skill (results not presented here).

Appendix B Tuning UQ-equipped NNs

In addition to the hyperparameters of the deterministic NNs, designing an architecture for UQ often demands additional hyperparameter optimization. For instance, for the DNN, decisions need to be made regarding the number of neurons to drop out (dropout rate). While less common, one can also choose whether to apply dropout to all hidden layers or only selected ones. Variations in the dropout rate and the layers to which dropout is applied can influence the final configuration and performance of the DNN. Figure B1 illustrates these effects. As we increase the number of dropped neurons (whether through a higher dropout rate or by subjecting more layers to dropout), the uncertainty in the DNN predictions tends to rise. Yet, there is a persistent pattern in the relationship between spread (IQR) and RMSE across the various plots in Figure B1. Specifically, as spread increases, RMSE concurrently grows, consistent with the insights highlighted in Figure 9.

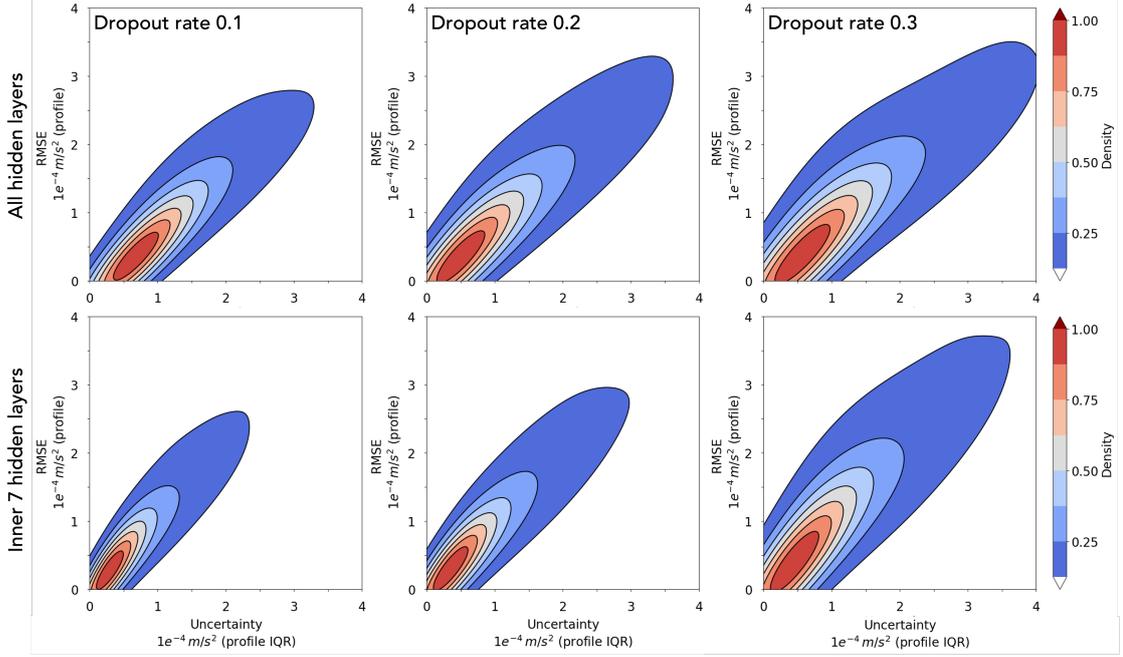


Figure B1. Similar to Figure 9 but for DNN only with different dropout rates, which are applied to different numbers of hidden layers.

778 In the case of BNN or VAE, even though there is no dropout rate, there are distinct
 779 tuning opportunities available. For instance, with the VAE, one might consider applying
 780 dropout to the NN emulator. Moreover, given that the loss function in VAE comprises three
 781 components, decisions can be made regarding which component to penalize more heavily,
 782 allowing for nuanced adjustments to its performance.

783 Appendix C The UQ metrics

784 Each point in the spread-skill plot corresponds to one specific bin of ensemble spread
 785 (\overline{SD}_k), which is defined as the average standard deviation of the ensemble members. We
 786 first separate the spread using a pre-selected number of bins (a subjective choice of 15 is
 787 used here). Then for the k^{th} bin:

$$\begin{cases} \text{RMSE}_k = \left[\frac{1}{N_k} \sum_{i=1}^{N_k} (\hat{y}_i - \bar{y}_i)^2 \right]^{\frac{1}{2}} \\ \overline{SD}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \left[\frac{1}{M-1} \sum_{j=1}^M (\bar{y}_i - y_{ij})^2 \right]^{\frac{1}{2}} \\ \bar{y}_i = \frac{1}{M} \sum_{j=1}^M y_{ij} \end{cases} \quad (\text{C1})$$

788 \hat{y}_i is the observed value for the i^{th} example, \bar{y}_i is the mean prediction for the i^{th} example, y_{ij}
 789 is the j^{th} prediction for the i^{th} example, N_k is the total number of examples in the k^{th} bin,
 790 and M is the ensemble size. Following Haynes et al. (2023), we summarize the quality of the

791 spread-skill plot by two measures: spread-skill reliability (SSREL) and overall spread-skill
 792 ratio (SSRAT). SSREL is the bin-weighted mean distance from the 1-to-1 line:

$$\text{SSREL} = \sum_{k=1}^K \frac{N_k}{N} |\text{RMSE}_k - \overline{\text{SD}}_k| \quad (\text{C2})$$

793 where N is the total number of examples, K is the total number of bins, and other variables
 794 are as in Equation C1. SSREL varies from $[0, \infty)$, and the ideal value is 0. On the other
 795 hand, SSRAT is averaged over the whole dataset:

$$\text{SSRAT} = \frac{\overline{\text{SD}}}{\text{RMSE}} \quad (\text{C3})$$

796 SSRAT also varies from $[0, \infty)$, and the ideal value is 1. $\text{SSRAT} > 1$ indicates the model is
 797 under-confident on average, while $\text{SSRAT} < 1$ indicates that the model is overconfident on
 798 average.

799 In Equation (C1), each level of a GWD profile is considered as an individual sample.
 800 As discussed earlier, while these samples help assess the model's overall performance, our
 801 main interest is often the uncertainty of individual GWD profiles. Such uncertainty informs
 802 the trustworthiness of the model's prediction for that specific profile. Accordingly, for each
 803 profile, we can compute:

$$\begin{cases} \text{RMSE}_{\text{profile}} = \left[\frac{1}{N_z} \sum_{z=1}^{N_z} (\hat{y}_z - \bar{y}_z)^2 \right]_{\text{profile}}^{\frac{1}{2}} \\ \text{IQR}_{\text{profile}} = \left[\frac{1}{N_z} \sum_{z=1}^{N_z} (y_{z,75th} - y_{z,25th})^2 \right]_{\text{profile}}^{\frac{1}{2}} \\ \bar{y}_z = \left[\frac{1}{M} \sum_{j=1}^M y_{zj} \right]_{\text{profile}} \end{cases} \quad (\text{C4})$$

804 where N_z is the number of vertical levels for each profile, and $\text{IQR}_{\text{profile}}$ is its interquartile
 805 range: $y_{z,25th}$ corresponds with the 25th percentile, and $y_{z,75th}$ corresponds with the 75th
 806 percentile.

807 Open Research

808 The data for all the analyses in the main text are available at [https://doi.org/10](https://doi.org/10.5281/zenodo.10019987)
 809 [.5281/zenodo.10019987](https://doi.org/10.5281/zenodo.10019987). The emulator code is available at [https://github.com/yqsun91/](https://github.com/yqsun91/WACCM-Emulation)
 810 [WACCM-Emulation](https://github.com/yqsun91/WACCM-Emulation). All the raw WACCM output data are available on request from authors.

811 Acknowledgments

812 We thank Andre Souza for insightful discussions. This work was supported by grants from
 813 the NSF OAC CSSI program (#2005123 , #2004512, #2004492, #2004572), and by the

814 generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program
 815 to PH, MJA, EG, and AS. PH is also supported by the Office of Naval Research (ONR)
 816 Young Investigator Award N00014-20-1-2722. SL is supported by the Office of Science, U.S.
 817 Department of Energy Biological and Environmental Research as part of the Regional and
 818 Global Climate Model Analysis program area. Computational resources were provided by
 819 NSF XSEDE (allocation ATM170020) and NCAR’s CISL (allocation URIC0009).

820 References

- 821 Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., ...
 822 Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques,
 823 applications and challenges. *Information Fusion*, *76*, 243-297. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S1566253521001081)
 824 www.sciencedirect.com/science/article/pii/S1566253521001081 doi: [https://doi](https://doi.org/10.1016/j.inffus.2021.05.008)
 825 [.org/10.1016/j.inffus.2021.05.008](https://doi.org/10.1016/j.inffus.2021.05.008)
- 826 Achatz, U. (2022). Gravity waves and their impact on the atmospheric flow. In *Atmo-*
 827 *spheric dynamics* (pp. 407–505). Berlin, Heidelberg: Springer Berlin Heidelberg. Re-
 828 trieved from https://doi.org/10.1007/978-3-662-63941-2_10 doi: 10.1007/978-3-
 829 -662-63941-2_10
- 830 Alexander, M. J., Geller, M., McLandress, C., Polavarapu, S., Preusse, P., Sassi, F., ...
 831 Watanabe, S. (2010). Recent developments in gravity-wave effects in climatemodels
 832 and the global distribution of gravity-wavemomentum flux from observations and models.
 833 *Quarterly Journal of the Royal Meteorological Society*, *136*. doi: 10.1002/qj.637
- 834 Amiramjadi, M., Plougonven, R., Mohebalhojeh, A. R., & Mirzaei, M. (2022). Using
 835 machine learning to estimate non-orographic gravity wave characteristics at source levels.
 836 *Journal of the Atmospheric Sciences*. doi: 10.1175/JAS-D-22-0021.1
- 837 Ando, S., & Huang, C. Y. (2017). Deep over-sampling framework for classifying imbalanced
 838 data. In *Machine learning and knowledge discovery in databases: European conference,*
 839 *ecml pkdd 2017, skopje, macedonia, september 18–22, 2017, proceedings, part i 10* (pp.
 840 770–785).
- 841 Bacmeister, J. T., Newman, P. A., Gary, B. L., & Chan, K. R. (1994). An algorithm for
 842 forecasting mountain wave-related turbulence in the stratosphere. *Weather Forecasting*,
 843 *9*. doi: 10.1175/1520-0434(1994)009<0241:AAFFMW>2.0.CO;2
- 844 Balaji, V. (2021). Climbing down charney’s ladder: machine learning and the post-dennard
 845 era of computational climate science. *Philosophical Transactions of the Royal Society A*,
 846 *379*(2194), 20200085.
- 847 Baldi, P., Sadowski, P., & Whiteson, D. (2014). Searching for exotic particles in high-energy
 848 physics with deep learning. *Nature communications*, *5*(1), 4308.
- 849 Ballnus, B., Hug, S., Hatz, K., Görlitz, L., Hasenauer, J., & Theis, F. J. (2017). Com-
 850 prehensive benchmarking of markov chain monte carlo methods for dynamical systems.
 851 *BMC Systems Biology*, *11*(1), 1–18.
- 852 Barnes, E. A., Barnes, R. J., & DeMaria, M. (2023). Sinh-arcsinh-normal distributions
 853 to add uncertainty to neural network regression tasks: Applications to tropical cyclone

- 854 intensity forecasts. *Environmental Data Science*, 2, e15.
- 855 Beck, A., Flad, D., & Munz, C.-D. (2019). Deep neural networks for data-driven LES
856 closure models. *Journal of Computational Physics*, 398, 108910.
- 857 Beljaars, A. C., Brown, A. R., & Wood, N. (2004). A new parametrization of turbulent
858 orographic form drag. *Quarterly Journal of the Royal Meteorological Society*, 130. doi:
859 10.1256/qj.03.73
- 860 Belochitski, A., & Krasnopolsky, V. (2021). Robustness of neural network emulations of ra-
861 diative transfer parameterizations in a state-of-the-art general circulation model. *Geosci-
862 entific Model Development*, 14(12), 7425–7437. Retrieved from [https://gmd.copernicus](https://gmd.copernicus.org/articles/14/7425/2021/)
863 [.org/articles/14/7425/2021/](https://gmd.copernicus.org/articles/14/7425/2021/) doi: 10.5194/gmd-14-7425-2021
- 864 Beucler, T., Ebert-Uphoff, I., Rasp, S., Pritchard, M., & Gentine, P. (2021, 05). *Machine
865 learning for clouds and climate (invited chapter for the agu geophysical monograph series
866 "clouds and climate")*. doi: 10.1002/essoar.10506925.1
- 867 Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015). *Weight uncertainty
868 in neural networks*.
- 869 Bolton, T., & Zanna, L. (2019). Applications of deep learning to ocean data inference
870 and subgrid parameterization. *Journal of Advances in Modeling Earth Systems*, 11(1),
871 376–399.
- 872 Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and
873 stabilizing machine-learning parametrizations of convection. *Journal of the Atmospheric
874 Sciences*, 77(12), 4357 - 4375. Retrieved from [https://journals.ametsoc.org/view/
875 journals/atms/77/12/jas-d-20-0082.1.xml](https://journals.ametsoc.org/view/journals/atms/77/12/jas-d-20-0082.1.xml) doi: [https://doi.org/10.1175/JAS-D-20-
876 -0082.1](https://doi.org/10.1175/JAS-D-20-0082.1)
- 877 Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network
878 parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Sys-
879 tems*, 11(8), 2728-2744. Retrieved from [https://agupubs.onlinelibrary.wiley.com/
880 doi/abs/10.1029/2019MS001711](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001711) doi: <https://doi.org/10.1029/2019MS001711>
- 881 Buda, M., Maki, A., & Mazurowski, M. A. (2018). A systematic study of the class imbalance
882 problem in convolutional neural networks. *Neural Networks*, 106, 249-259. Retrieved from
883 <https://www.sciencedirect.com/science/article/pii/S0893608018302107> doi:
884 <https://doi.org/10.1016/j.neunet.2018.07.011>
- 885 Bölöni, G., Kim, Y. H., Borchert, S., & Achatz, U. (2021). Toward transient subgrid-scale
886 gravity wave representation in atmospheric models. part i: Propagation model including
887 nondissipative wave mean-flow interactions. *Journal of the Atmospheric Sciences*, 78.
888 doi: 10.1175/JAS-D-20-0065.1
- 889 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Ma-
890 chine learning emulation of gravity wave drag in numerical weather forecasting. *Jour-
891 nal of Advances in Modeling Earth Systems*, 13(7), e2021MS002477. Retrieved
892 from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002477>
893 (e2021MS002477 2021MS002477) doi: <https://doi.org/10.1029/2021MS002477>
- 894 Chattopadhyay, A., Nabizadeh, E., & Hassanzadeh, P. (2020). Analog forecasting of
895 extreme-causing weather patterns using deep learning. *Journal of Advances in Model-*

- 896 *ing Earth Systems*, 12(2), e2019MS001958.
- 897 Chattopadhyay, A., Subel, A., & Hassanzadeh, P. (2020). Data-driven super-
898 parameterization using deep learning: Experimentation with multiscale Lorenz 96 sys-
899 tems and transfer learning. *Journal of Advances in Modeling Earth Systems*, 12(11),
900 e2020MS002084.
- 901 Chawla, N. V., Japkowicz, N., & Kotcz, A. (2004, jun). Editorial: Special issue on learning
902 from imbalanced data sets. *SIGKDD Explor. Newsl.*, 6(1), 1–6. Retrieved from [https://](https://doi.org/10.1145/1007730.1007733)
903 doi.org/10.1145/1007730.1007733 doi: 10.1145/1007730.1007733
- 904 Chen, N., & Majda, A. J. (2019). A new efficient parameter estimation algorithm for high-
905 dimensional complex nonlinear turbulent dynamical systems with partial observations.
906 *Journal of Computational Physics*, 397, 108836.
- 907 Clare, M. C., Sonnewald, M., Lguensat, R., Deshayes, J., & Balaji, V. (2022). Explainable
908 artificial intelligence for bayesian neural networks: Towards trustworthy predictions of
909 ocean dynamics. *arXiv preprint arXiv:2205.00202*.
- 910 Delle Monache, L., Eckel, F. A., Rife, D. L., Nagarajan, B., & Searight, K. (2013). Probabilistic
911 weather prediction with an analog ensemble. *Monthly Weather Review*, 141(10),
912 3498–3516.
- 913 Dong, W., Fritts, D. C., Liu, A. Z., Lund, T. S., Liu, H.-L., & Snively, J.
914 (2023). Accelerating atmospheric gravity wave simulations using machine learning:
915 Kelvin-helmholtz instability and mountain wave sources driving gravity wave break-
916 ing and secondary gravity wave generation. *Geophysical Research Letters*, 50(15),
917 e2023GL104668. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023GL104668)
918 [abs/10.1029/2023GL104668](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023GL104668) (e2023GL104668 2023GL104668) doi: [https://doi.org/](https://doi.org/10.1029/2023GL104668)
919 [10.1029/2023GL104668](https://doi.org/10.1029/2023GL104668)
- 920 Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J. (2022).
921 Machine Learning Gravity Wave Parameterization Generalizes to Capture the QBO and
922 Response to Increased CO₂. , 49(8), e98174. doi: 10.1029/2022GL098174
- 923 Finkel, J., Gerber, E. P., Abbot, D. S., & Weare, J. (2023). Revealing the statis-
924 tics of extreme events hidden in short weather forecast data. *AGU Advances*, 4(2),
925 e2023AV000881. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023AV000881)
926 [abs/10.1029/2023AV000881](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023AV000881) (e2023AV000881 2023AV000881) doi: [https://doi.org/](https://doi.org/10.1029/2023AV000881)
927 [10.1029/2023AV000881](https://doi.org/10.1029/2023AV000881)
- 928 Foster, D., Gagne, D. J., & Whitt, D. B. (2021, 12). Probabilistic machine learning esti-
929 mation of ocean mixed layer depth from dense satellite and sparse in situ observations.
930 *Journal of Advances in Modeling Earth Systems*, 13. doi: 10.1029/2021MS002474
- 931 Frezat, H., Sommer, J. L., Fablet, R., Balarac, G., & Lguensat, R. (2022). A posteriori learn-
932 ing for quasi-geostrophic turbulence parametrization. *arXiv preprint arXiv:2204.03911*.
933 doi: 10.1029/2022MS003124
- 934 Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine
935 learning for stochastic parameterization: Generative adversarial networks in the lorenz'96
936 model. *Journal of Advances in Modeling Earth Systems*, 12(3), e2019MS001896.
- 937 Gal, Y., & Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing

- 938 model uncertainty in deep learning. In *international conference on machine learning* (pp.
939 1050–1059).
- 940 Garcia, R. R., Smith, A. K., Kinnison, D. E., Álvaro de la Cámara, & Murphy, D. J. (2017).
941 Modification of the gravity wave parameterization in the whole atmosphere community
942 climate model: Motivation and results. *Journal of the Atmospheric Sciences*, *74*(1),
943 275 - 291. Retrieved from [https://journals.ametsoc.org/view/journals/atasc/74/
944 1/jas-d-16-0104.1.xml](https://journals.ametsoc.org/view/journals/atasc/74/1/jas-d-16-0104.1.xml) doi: <https://doi.org/10.1175/JAS-D-16-0104.1>
- 945 Geller, M. A., Alexander, M. J., Love, P. T., Bacmeister, J., Ern, M., Hertzog, A., ...
946 others (2013). A comparison between gravity wave momentum fluxes in observations and
947 climate models. *Journal of Climate*, *26*(17), 6383–6405.
- 948 Gettelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J., Morrison, H.,
949 & Gantos, G. (2021). Machine learning the warm rain process. *Journal of Advances in
950 Modeling Earth Systems*, *13*. doi: 10.1029/2020MS002268
- 951 Gettelman, A., Mills, M. J., Kinnison, D. E., Garcia, R. R., Smith, A. K., Marsh, D. R.,
952 ... Randel, W. J. (2019, 12). The whole atmosphere community climate model version
953 6 (waccm6). *Journal of Geophysical Research: Atmospheres*, *124*, 12380-12403. doi:
954 10.1029/2019JD030943
- 955 Gordon, E. M., & Barnes, E. A. (2022, 8). Incorporating uncertainty into a regression
956 neural network enables identification of decadal state-dependent predictability in cesm2.
957 *Geophysical Research Letters*, *49*. doi: 10.1029/2022GL098635
- 958 Guan, Y., Chattopadhyay, A., Subel, A., & Hassanzadeh, P. (2022). Stable a posteriori
959 LES of 2D turbulence using convolutional neural networks: Backscattering analysis and
960 generalization to higher Re via transfer learning. *Journal of Computational Physics*, *458*,
961 111090.
- 962 Guan, Y., Subel, A., Chattopadhyay, A., & Hassanzadeh, P. (2023). Learning physics-
963 constrained subgrid-scale closures in the small-data regime for stable and accurate les.
964 *Physica D: Nonlinear Phenomena*, *443*, 133568. doi: [https://doi.org/10.1016/j.physd
965 .2022.133568](https://doi.org/10.1016/j.physd.2022.133568)
- 966 Guillaumin, A. P., & Zanna, L. (2021). Stochastic-deep learning parameterization of
967 ocean momentum forcing. *Journal of Advances in Modeling Earth Systems*, *13*(9),
968 e2021MS002534.
- 969 Hardiman, S. C., Scaife, A. A., Niekerk, A. v., Prudden, R., Owen, A., Adams, S. V., ...
970 Madge, S. (2023). Machine learning for non-orographic gravity waves in a climate model.
971 *Artificial Intelligence for the Earth Systems*.
- 972 Haynes, K., Lagerquist, R., McGraw, M., Musgrave, K., & Ebert-Uphoff, I. (2023). Creat-
973 ing and evaluating uncertainty estimates with neural networks for environmental-science
974 applications. *Artificial Intelligence for the Earth Systems*, 1–58.
- 975 Hertzog, A., Alexander, M. J., & Plougonven, R. (2012). On the intermittency of gravity
976 wave momentum flux in the stratosphere. *Journal of the Atmospheric Sciences*, *69*(11),
977 3433–3448.
- 978 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ... others
979 (2017). The art and science of climate model tuning. *Bulletin of the American Meteorolo-*

- 980 *logical Society*, 98(3), 589–602.
- 981 Huang, C., Li, Y., Loy, C. C., & Tang, X. (2016). Learning deep representation for imbal-
982 anced classification. In *Proceedings of the IEEE conference on computer vision and pattern*
983 *recognition* (pp. 5375–5384).
- 984 Iglesias-Suarez, F., Gentine, P., Solino-Fernandez, B., Beucler, T., Pritchard, M., Runge,
985 J., & Eyring, V. (2023). *Causally-informed deep learning to improve climate models and*
986 *projections*.
- 987 Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study.
988 *Intelligent Data Analysis*, 6, 429-449. Retrieved from [https://doi.org/10.3233/IDA-](https://doi.org/10.3233/IDA-2002-6504)
989 [-2002-6504](https://doi.org/10.3233/IDA-2002-6504) (5) doi: 10.3233/IDA-2002-6504
- 990 Johnson, J. M., & Khoshgoftaar, T. M. (2019, Mar 19). Survey on deep learning with class
991 imbalance. *Journal of Big Data*, 6(1), 27. Retrieved from [https://doi.org/10.1186/s40537-](https://doi.org/10.1186/s40537-019-0192-5)
992 [s40537-019-0192-5](https://doi.org/10.1186/s40537-019-0192-5) doi: 10.1186/s40537-019-0192-5
- 993 Kim, Y., Bölöni, G., Borchert, S., Chun, H. Y., & Achatz, U. (2021). Toward transient
994 subgrid-scale gravity wave representation in atmospheric models. part ii: Wave intermit-
995 tency simulated with convective sources. *Journal of the Atmospheric Sciences*, 78. doi:
996 10.1175/JAS-D-20-0066.1
- 997 Kim, Y., Eckermann, S. D., & Chun, H. (2003). An overview of the past, present and
998 future of gravity-wave drag parametrization for numerical climate and weather prediction
999 models. *Atmosphere-Ocean*, 41, 65-98. Retrieved from [https://doi.org/10.3137/ao-](https://doi.org/10.3137/ao.410105)
1000 [.410105](https://doi.org/10.3137/ao.410105) doi: 10.3137/ao.410105
- 1001 Kingma, D. P., & Welling, M. (2014). Auto-encoding variational bayes. *2nd International*
1002 *Conference on Learning Representations, ICLR 2014 - Conference Track Proceedings*.
- 1003 Köhler, L., Green, B., & Stephan, C. C. (2023). Comparing loon superpressure balloon
1004 observations of gravity waves in the tropics with global storm-resolving models. *Journal*
1005 *of Geophysical Research: Atmospheres*, 128(15), e2023JD038549.
- 1006 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). New approach to
1007 calculation of atmospheric model physics: Accurate and fast neural network emulation
1008 of longwave radiation in a climate model. *Monthly Weather Review*, 133(5), 1370 -
1009 1383. Retrieved from [https://journals.ametsoc.org/view/journals/mwre/133/5/](https://journals.ametsoc.org/view/journals/mwre/133/5/mwr2923.1.xml)
1010 [mwr2923.1.xml](https://journals.ametsoc.org/view/journals/mwre/133/5/mwr2923.1.xml) doi: <https://doi.org/10.1175/MWR2923.1>
- 1011 Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., ... Courville,
1012 A. (2021). Out-of-distribution generalization via risk extrapolation (rex). In *International*
1013 *conference on machine learning* (pp. 5815–5826).
- 1014 Kruse, C. G., Alexander, M. J., Hoffmann, L., Niekerk, A. V., Polichtchouk, I., Bacmeister,
1015 J. T., ... Stein, O. (2022). Observed and modeled mountain waves from the surface to
1016 the mesosphere near the drake passage. *Journal of the Atmospheric Sciences*, 79. doi:
1017 10.1175/JAS-D-21-0252.1
- 1018 Li, X., Dai, Y., Ge, Y., Liu, J., Shan, Y., & Duan, L.-Y. (2022). Uncertainty modeling for
1019 out-of-distribution generalization. *arXiv preprint arXiv:2202.03958*.
- 1020 Ling, J., & Templeton, J. (2015, 08). Evaluation of machine learning algorithms for
1021 prediction of regions of high Reynolds averaged Navier Stokes uncertainty. *Physics of*

- 1022 *Fluids*, 27(8). Retrieved from <https://doi.org/10.1063/1.4927765> (085103) doi:
1023 10.1063/1.4927765
- 1024 Liu, Y., Racah, E., Prabhat, M., Correa, J., Khosrowshahi, A., Lavers, D., ... Collins,
1025 W. (2016, 05). Application of deep convolutional neural networks for detecting extreme
1026 weather in climate datasets.
- 1027 Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. (2022). Global extreme heat
1028 forecasting using neural weather models. *Artificial Intelligence for the Earth Systems*, 1
1029 - 41. Retrieved from [https://journals.ametsoc.org/view/journals/aies/aop/AIES-](https://journals.ametsoc.org/view/journals/aies/aop/AIES-D-22-0035.1/AIES-D-22-0035.1.xml)
1030 [D-22-0035.1/AIES-D-22-0035.1.xml](https://journals.ametsoc.org/view/journals/aies/aop/AIES-D-22-0035.1/AIES-D-22-0035.1.xml) doi: 10.1175/AIES-D-22-0035.1
- 1031 Lu, L., Jin, P., Pang, G., Zhang, Z., & Karniadakis, G. E. (2021). Learning nonlinear
1032 operators via deepnet based on the universal approximation theorem of operators. *Nature*
1033 *machine intelligence*, 3(3), 218–229.
- 1034 Maalouf, M., & Siddiqi, M. (2014). Weighted logistic regression for large-scale imbalanced
1035 and rare events data. *Knowledge-Based Systems*, 59, 142–148.
- 1036 Maalouf, M., & Trafalis, T. B. (2011). Robust weighted kernel logistic regression in imbal-
1037 anced and rare events data. *Computational Statistics & Data Analysis*, 55(1), 168–183.
- 1038 Maddox, W. J., Izmailov, P., Garipov, T., Vetrov, D. P., & Wilson, A. G. (2019). A
1039 simple baseline for bayesian uncertainty in deep learning. *Advances in Neural Information*
1040 *Processing Systems*, 32.
- 1041 Mamalakis, A., Barnes, E. A., & Ebert-Uphoff, I. (2022). Investigating the fidelity of
1042 explainable artificial intelligence methods for applications of convolutional neural networks
1043 in geoscience. *Artificial Intelligence for the Earth Systems*, 1(4), e220012.
- 1044 Matsuoka, D., Watanabe, S., Sato, K., Kawazoe, S., Yu, W., & Easterbrook, S. (2020).
1045 Application of Deep Learning to Estimate Atmospheric Gravity Wave Parameters in Re-
1046 analysis Data Sets. , 47(19), e89436. doi: 10.1029/2020GL089436
- 1047 Maulik, R., San, O., Rasheed, A., & Vedula, P. (2019). Subgrid modelling for two-
1048 dimensional turbulence using neural networks. *Journal of Fluid Mechanics*, 858, 122–
1049 144.
- 1050 McGovern, A., Lagerquist, R., Gagne, D. J., Jergensen, G. E., Elmore, K. L., Homeyer,
1051 C. R., & Smith, T. (2019). Making the black box more transparent: Understanding
1052 the physical implications of machine learning. *Bulletin of the American Meteorological*
1053 *Society*, 100. doi: 10.1175/BAMS-D-18-0195.1
- 1054 Miller, J. P., Taori, R., Raghunathan, A., Sagawa, S., Koh, P. W., Shankar, V., ... Schmidt,
1055 L. (2021). Accuracy on the line: on the strong correlation between out-of-distribution
1056 and in-distribution generalization. In *International conference on machine learning* (pp.
1057 7721–7735).
- 1058 Miloshevich, G., Cozian, B., Abry, P., Borgnat, P., & Bouchet, F. (2023, Apr). Probabilistic
1059 forecasts of extreme heatwaves using convolutional neural networks in a regime of lack
1060 of data. *Phys. Rev. Fluids*, 8, 040501. Retrieved from [https://link.aps.org/doi/](https://link.aps.org/doi/10.1103/PhysRevFluids.8.040501)
1061 [10.1103/PhysRevFluids.8.040501](https://link.aps.org/doi/10.1103/PhysRevFluids.8.040501) doi: 10.1103/PhysRevFluids.8.040501
- 1062 Nagarajan, V., Andreassen, A., & Neyshabur, B. (2020). Understanding the failure modes
1063 of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*.

- 1064 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameterize moist
 1065 convection: Potential for modeling of climate, climate change, and extreme events. *Journal*
 1066 *of Advances in Modeling Earth Systems*, 10(10), 2548-2563. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351)
 1067 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351 doi: [https://](https://doi.org/10.1029/2018MS001351)
 1068 doi.org/10.1029/2018MS001351
- 1069 Oh, S. M., Rehg, J. M., Balch, T., & Dellaert, F. (2005). Data-driven mcmc for learning and
 1070 inference in switching linear dynamic systems. In *Proceedings of the national conference*
 1071 *on artificial intelligence* (Vol. 20, p. 944).
- 1072 Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., ... Snoek, J. (2019).
 1073 Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset
 1074 shift. *Advances in neural information processing systems*, 32.
- 1075 Pahlavan, H. A., Hassanzadeh, P., & Alexander, M. J. (2023). Explainable offline-online
 1076 training of neural networks for parameterizations: A 1d gravity wave-qbo testbed in the
 1077 small-data regime. *arXiv preprint arXiv:2309.09024*.
- 1078 Palmer, T. (2019). Stochastic weather and climate models. *Nature Reviews Physics*, 1(7),
 1079 463–471.
- 1080 Polichtchouk, I., Van Niekerk, A., & Wedi, N. (2023). Resolved gravity waves in the
 1081 extratropical stratosphere: Effect of horizontal resolution increase from o (10) to o (1)
 1082 km. *Journal of the Atmospheric Sciences*, 80(2), 473–486.
- 1083 Prein, A. F., Langhans, W., Fossier, G., Ferrone, A., Ban, N., Goergen, K., ... others
 1084 (2015). A review on regional convection-permitting climate modeling: Demonstrations,
 1085 prospects, and challenges. *Reviews of geophysics*, 53(2), 323–361.
- 1086 Psaros, A. F., Meng, X., Zou, Z., Guo, L., & Karniadakis, G. E. (2023). Uncertainty
 1087 quantification in scientific machine learning: Methods, metrics, and comparisons. *Journal*
 1088 *of Computational Physics*, 477, 111902.
- 1089 Qi, D., & Majda, A. J. (2020). Using machine learning to predict extreme events in complex
 1090 systems. *Proceedings of the National Academy of Sciences*, 117(1), 52-59. Retrieved
 1091 from <https://www.pnas.org/doi/abs/10.1073/pnas.1917285117> doi: 10.1073/pnas
 1092 .1917285117
- 1093 Ramachandran, P., Zoph, B., & Le, Q. V. (2017). Searching for activation functions. *arXiv*
 1094 *preprint arXiv:1710.05941*.
- 1095 Rasp, S. (2020). Coupled online learning as a way to tackle instabilities and biases in neural
 1096 network parameterizations: general algorithms and lorenz 96 case study (v1.0). *Geosci-*
 1097 *entific Model Development*, 13(5), 2185–2196. Retrieved from [https://gmd.copernicus](https://gmd.copernicus.org/articles/13/2185/2020/)
 1098 [.org/articles/13/2185/2020/](https://gmd.copernicus.org/articles/13/2185/2020/) doi: 10.5194/gmd-13-2185-2020
- 1099 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid
 1100 processes in climate models. *Proceedings of the National Academy of Sciences*, 115(39),
 1101 9684–9689.
- 1102 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., &
 1103 Prabhat. (2019, Feb 01). Deep learning and process understanding for data-driven earth
 1104 system science. *Nature*, 566(7743), 195-204. Retrieved from [https://doi.org/10.1038/](https://doi.org/10.1038/s41586-019-0912-1)
 1105 [s41586-019-0912-1](https://doi.org/10.1038/s41586-019-0912-1) doi: 10.1038/s41586-019-0912-1

- 1106 Richter, J. H., Butchart, N., Kawatani, Y., Bushell, A. C., Holt, L., Serva, F., ... Yukimoto,
 1107 S. (2022, 4). Response of the quasi-biennial oscillation to a warming climate in global
 1108 climate models. *Quarterly Journal of the Royal Meteorological Society*, *148*, 1490-1518.
 1109 doi: 10.1002/qj.3749
- 1110 Richter, J. H., Sassi, F., & Garcia, R. R. (2010). Toward a physically based gravity
 1111 wave source parameterization in a general circulation model. *Journal of the Atmospheric*
 1112 *Sciences*, *67*. doi: 10.1175/2009JAS3112.1
- 1113 Ross, A. S., Li, Z., Perezhugin, P., Fernandez-Granda, C., & Zanna, L. (2022). Benchmark-
 1114 ing of machine learning ocean subgrid parameterizations in an idealized model.
- 1115 Sato, K., Watanabe, S., Kawatani, Y., Tomikawa, Y., Miyazaki, K., & Takahashi, M. (2009).
 1116 On the origins of mesospheric gravity waves. *Geophysical Research Letters*, *36*. doi:
 1117 10.1029/2009GL039908
- 1118 Schneider, T., Jeevanjee, N., & Socolow, R. (2021). Accelerating progress in climate science.
 1119 *Physics Today*, *74*(6), 44–51.
- 1120 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling 2.0: A
 1121 blueprint for models that learn from observations and targeted high-resolution simula-
 1122 tions. *Geophysical Research Letters*, *44*. doi: 10.1002/2017GL076101
- 1123 Scinocca, J. F., & McFarlane, N. A. (2000). The parametrization of drag induced by
 1124 stratified flow over anisotropic orography. *Quarterly Journal of the Royal Meteorological*
 1125 *Society*, *126*. doi: 10.1002/qj.49712656802
- 1126 Seifert, A., & Rasp, S. (2020). Potential and limitations of machine learning for modeling
 1127 warm-rain cloud microphysical processes. *Journal of Advances in Modeling Earth Sys-*
 1128 *tems*, *12*(12), e2020MS002301. Retrieved from [https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020MS002301)
 1129 [.com/doi/abs/10.1029/2020MS002301](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020MS002301) (e2020MS002301 10.1029/2020MS002301) doi:
 1130 <https://doi.org/10.1029/2020MS002301>
- 1131 Shamekh, S., Lamb, K. D., Huang, Y., & Gentine, P. (2023). Implicit learning of convective
 1132 organization explains precipitation stochasticity. *Proceedings of the National Academy*
 1133 *of Sciences*, *120*(20), e2216158120. Retrieved from [https://www.pnas.org/doi/abs/](https://www.pnas.org/doi/abs/10.1073/pnas.2216158120)
 1134 [10.1073/pnas.2216158120](https://www.pnas.org/doi/abs/10.1073/pnas.2216158120) doi: 10.1073/pnas.2216158120
- 1135 Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., & Cui, P. (2021). Towards out-of-
 1136 distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*.
- 1137 Song, H.-J., & Roh, S. (2021). Improved weather forecasting using neural network emulation
 1138 for radiation parameterization. *Journal of Advances in Modeling Earth Systems*, *13*(10),
 1139 e2021MS002609. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002609)
 1140 [abs/10.1029/2021MS002609](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002609) (e2021MS002609 2021MS002609) doi: [https://doi.org/](https://doi.org/10.1029/2021MS002609)
 1141 [10.1029/2021MS002609](https://doi.org/10.1029/2021MS002609)
- 1142 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014).
 1143 Dropout: a simple way to prevent neural networks from overfitting. *The journal of*
 1144 *machine learning research*, *15*(1), 1929–1958.
- 1145 Stensrud, D. J. (2007). *Parameterization schemes: Keys to understanding numerical*
 1146 *weather prediction models*. Cambridge University Press. Retrieved from [https://](https://www.cambridge.org/core/product/identifier/9780511812590/type/book)
 1147 www.cambridge.org/core/product/identifier/9780511812590/type/book doi: 10

- 1148 .1017/CBO9780511812590
- 1149 Subel, A., Chattopadhyay, A., Guan, Y., & Hassanzadeh, P. (2021). Data-driven subgrid-
1150 scale modeling of forced Burgers turbulence using deep learning with generalization to
1151 higher Reynolds numbers via transfer learning. *Physics of Fluids*, *33*(3), 031702.
- 1152 Subel, A., Guan, Y., Chattopadhyay, A., & Hassanzadeh, P. (2023). Explaining the physics
1153 of transfer learning in data-driven turbulence modeling. *PNAS nexus*, *2*(3), pgad015.
- 1154 Sun, Y., Hassanzadeh, P., Alexander, M., & Kruse, C. (2023, 12). Quantifying 3d grav-
1155 ity wave drag in a library of tropical convection-permitting simulations for data-driven
1156 parameterizations.
1157 doi: 10.1029/2022MS003585
- 1158 Sun, Y., Wong, A., & Kamel, M. S. (2009, 11). Classification of imbalanced data: a review.
1159 *International Journal of Pattern Recognition and Artificial Intelligence*, *23*, 687-719. doi:
1160 10.1142/S0218001409007326
- 1161 Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., & Liu, C. (2018). A survey on deep
1162 transfer learning. , 270–279.
- 1163 Wedi, N. P., Polichtchouk, I., Dueben, P., Anantharaj, V. G., Bauer, P., Boussetta, S., ...
1164 others (2020). A baseline for global weather and climate simulations at 1 km resolution.
1165 *Journal of Advances in Modeling Earth Systems*, *12*(11), e2020MS002192.
- 1166 Wu, D., Gao, L., Xiong, X., Chinazzi, M., Vespignani, A., Ma, Y.-A., & Yu, R. (2021).
1167 *Quantifying uncertainty in deep spatiotemporal forecasting*.
- 1168 Wu, G., & Chang, E. Y. (2003). Adaptive feature-space conformal transformation for
1169 imbalanced-data learning. In *Proceedings of the twentieth international conference on
1170 international conference on machine learning* (p. 816–823). AAAI Press.
- 1171 Ye, H., Xie, C., Cai, T., Li, R., Li, Z., & Wang, L. (2021). Towards a theoretical framework
1172 of out-of-distribution generalization. *Advances in Neural Information Processing Systems*,
1173 *34*, 23519–23531.
- 1174 Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in
1175 deep neural networks? *Advances in neural information processing systems*, *27*.
- 1176 Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid
1177 processes for climate modeling at a range of resolutions. *Nature communications*, *11*(1),
1178 1–10.
- 1179 Zhang, C., Perezhugin, P., Gultekin, C., Adcroft, A., Fernandez-Granda, C., & Zanna, L.
1180 (2023). *Implementation and evaluation of a machine learned mesoscale eddy parameteri-
1181 zation into a numerical ocean circulation model*.
- 1182 Zhang, D., Ahuja, K., Xu, Y., Wang, Y., & Courville, A. (2021). Can subnetwork structure
1183 be the key to out-of-distribution generalization? In *International conference on machine
1184 learning* (pp. 12356–12367).
- 1185 Zhu, Y., Zabarar, N., Koutsourelakis, P.-S., & Perdikaris, P. (2019). Physics-constrained
1186 deep learning for high-dimensional surrogate modeling and uncertainty quantification
1187 without labeled data. *Journal of Computational Physics*, *394*, 56–81.

Abstract

Neural networks (NNs) are increasingly used for data-driven subgrid-scale parameterization in weather and climate models. While NNs are powerful tools for learning complex nonlinear relationships from data, there are several challenges in using them for parameterizations. Three of these challenges are 1) data imbalance related to learning rare (often large-amplitude) samples; 2) uncertainty quantification (UQ) of the predictions to provide an accuracy indicator; and 3) generalization to other climates, e.g., those with higher radiative forcing. Here, we examine performance of methods for addressing these challenges using NN-based emulators of the Whole Atmosphere Community Climate Model (WACCM) physics-based gravity wave (GW) parameterizations as the test case. WACCM has complex, state-of-the-art parameterizations for orography-, convection- and frontal-driven GWs. Convection- and orography-driven GWs have significant data imbalance due to the absence of convection or orography in many grid points. We address data imbalance using resampling and/or weighted loss functions, enabling the successful emulation of parameterizations for all three sources. We demonstrate that three UQ methods (Bayesian NNs, variational auto-encoders, and dropouts) provide ensemble spreads that correspond to accuracy during testing, offering criteria on when a NN gives inaccurate predictions. Finally, we show that accuracy of these NNs decreases for a warmer climate ($4\times\text{CO}_2$). However, the generalization accuracy is significantly improved by applying transfer learning, e.g., re-training only one layer using $\sim 1\%$ new data from the warmer climate. The findings of this study offer insights for developing reliable and generalizable data-driven parameterizations for various processes, including (but not limited) to GWs.

Plain Language Summary

Scientists are increasingly using machine learning methods, especially neural networks (NNs), to improve weather and climate models. However, it can be challenging for a NN to learn rare, large-amplitude events, because they are infrequent in training data. Also, NNs need to express their confidence (certainty) about a prediction and work effectively across different climates, e.g., warmer climates due to increased CO_2 . Traditional NNs often struggle with these challenges. Here, we share insights gained from emulating the complex physics-based parameterization schemes for gravity waves in a state-of-the-art climate model. We propose specific strategies for addressing imbalanced data, uncertainty quantification (UQ), and making accurate predictions across various climates. For instance, to manage data balance, one such strategy involves amplifying the impact of infrequent events in the training data. We also demonstrate that several UQ methods could be useful in determining the accuracy of predictions. Furthermore, we show that NNs trained on simulations of the historical period do not perform as well in warmer climates. However, we improve the NNs' performance by employing transfer learning using limited data from warmer climates. This study provides lessons for developing robust and generalizable approaches for using NNs to improve models in the future.

59 1 Introduction

60 Small-scale processes such as moist convection, gravity waves, and turbulence are key
 61 players in the variability of the climate system and its response to increased greenhouse
 62 gases. However, as these processes cannot be resolved, entirely or partially, by the coarse-
 63 resolution general circulation models (GCMs), they need to be represented as functions of
 64 the resolved dynamics via subgrid-scale (SGS) parameterization schemes (e.g., Kim et al.,
 65 2003; Stensrud, 2007; Prein et al., 2015). Many of these parameterization schemes are based
 66 on heuristic approximations and simplifications, introducing large parametric and epistemic
 67 uncertainties in GCMs (Schneider et al., 2017; Hourdin et al., 2017; Palmer, 2019).

68 Recently, there has been a growing interest in developing data-driven SGS parameter-
 69 izations for different complex processes in the Earth system using machine learning (ML)
 70 techniques, particularly deep neural networks (NNs). Promising results have been demon-
 71 strated in a wide range of idealized applications, including prototype systems (Maulik et al.,
 72 2019; Gagne et al., 2020; Rasp, 2020; Chattopadhyay, Subel, & Hassanzadeh, 2020; Frezat
 73 et al., 2022; Guan et al., 2022; Pahlavan et al., 2023), ocean turbulent processes (Bolton
 74 & Zanna, 2019; C. Zhang et al., 2023), moist convection in the atmosphere (O’Gorman &
 75 Dwyer, 2018; Brenowitz & Bretherton, 2019; Yuval & O’Gorman, 2020; Beucler et al., 2021;
 76 Iglesias-Suarez et al., 2023), radiation (Krasnopolsky et al., 2005; Belochitski & Krasnopol-
 77 sky, 2021; Song & Roh, 2021), and microphysics (Seifert & Rasp, 2020; Gettelman et al.,
 78 2021). The ultimate promise of data-driven parameterizations, learned from observation-
 79 derived data and/or high-fidelity high-resolution simulations, is that they might have smaller
 80 parametric/structural errors, thus reducing the biases of GCMs and producing more reliable
 81 climate change projections (e.g., Schneider et al., 2017; Reichstein et al., 2019; Schneider et
 82 al., 2021).

83 However, there are major challenges in developing trustworthy, interpretable, stable,
 84 and generalizable data-driven parameterizations that can be used for such climate change
 85 projection efforts. Discussing and even listing all of these challenges is well beyond the
 86 scope of this paper. Well-known challenges such as interpretability and stability have been
 87 extensively discussed in a number of recent studies (e.g., McGovern et al., 2019; Beck et al.,
 88 2019; Brenowitz et al., 2020; Balaji, 2021; Clare et al., 2022; Mamalakis et al., 2022; Guan
 89 et al., 2022; Subel et al., 2023; Pahlavan et al., 2023). Here, we focus on three other key
 90 issues:

- 91 1. Data imbalance, and related to that, learning rare/extreme events,
- 92 2. Uncertainty quantification (UQ) of the NN-based SGS parameterization outputs,
- 93 3. Out-of-distribution (OOD) generalization (e.g., extrapolation to climates with higher
 94 radiative forcings).

95 Below we briefly discuss the importance of 1-3 and the current state-of-the-art methods
 96 in addressing them in the climate and ML literature. Data imbalance is a well-known prob-
 97 lem in the ML literature, especially in the context of classification tasks (e.g., Japkowicz &

98 Stephen, 2002; G. Wu & Chang, 2003; Chawla et al., 2004; Sun et al., 2009; Huang et al.,
 99 2016; Ando & Huang, 2017; Buda et al., 2018; Johnson & Khoshgoftaar, 2019). The prob-
 100 lem becomes particularly significant when one aims to learn rare/extreme events (Maalouf
 101 & Trafalis, 2011; Maalouf & Siddiqi, 2014; Baldi et al., 2014; Liu et al., 2016; O’Gorman &
 102 Dwyer, 2018; Qi & Majda, 2020; Chattopadhyay, Nabizadeh, & Hassanzadeh, 2020; Milo-
 103 shevich et al., 2023; Finkel et al., 2023; Shamekh et al., 2023). For example, suppose we
 104 aim to learn the binary classification of the 99 percentile of temperature anomalies using a
 105 NN. In this case, label 0 (no extreme) will constitute 99% of the training (or testing) set
 106 while label 1 (extreme) will be just 1%. With many common loss functions such as mean
 107 squared error (MSE) or root-mean squared-error (RMSE), training a NN will result in one
 108 that predicts 0 for any sample (extreme or no extreme) while having a seemingly high ac-
 109 curacy of 99% (of course, other metrics such as precision/recall will show the shortcoming,
 110 see Chattopadhyay, Nabizadeh, & Hassanzadeh (2020)). The most common remedy to this
 111 problem for classification tasks is resampling. An example is down-sampling non-extreme
 112 cases by a factor of 100, which effectively balances the dataset.

113 In addition to *classification* tasks, Data imbalance also presents a significant challenge
 114 in *regression* tasks required for parameterization schemes in climate models. As highlighted
 115 by Chantry et al. (2021), such imbalances contributed to the unsuccessful emulation of
 116 their orographic gravity wave parameterization (GWP) scheme, largely because orography
 117 affects the gravity wave (GW) drag in only a fraction of the grid columns. This challenge also
 118 persists in emulating GWP for non-orographic GWs, especially when GWs are intricately
 119 linked to their sources. For instance, the presence of zero convective GW drag at numerous
 120 grid points due to the absence of convection creates a notably imbalanced dataset. This
 121 issue will be explored further in the results section. In regression tasks, data imbalance
 122 may also manifest in the form of difficulty in learning large-amplitude (extreme) outputs,
 123 which are rare and constitute only a small fraction of the training set. In the case of GWs,
 124 Observations have shown that gravity wave amplitudes are highly intermittent such that
 125 the largest 10% events alone can contribute more than 50% of the total momentum flux
 126 (Hertzog et al., 2012), so the extreme events will contribute an outsized fraction of the
 127 total drag. Nonetheless, poorly learning these large-amplitude outputs, like drag forces,
 128 can result in instabilities (e.g., Guan et al., 2022). Addressing data imbalance in climate
 129 applications has received relatively limited attention. In this study, we propose several
 130 remedies based on resampling techniques and weighted loss functions, demonstrating their
 131 advantages in enabling successful emulations of all GWP schemes and improving the learning
 132 of rare/extreme events.

133 Quantifying the uncertainties in outputs from NN-based parameterization schemes is
 134 essential when employing these schemes, particularly for high-stakes decision-making tasks
 135 such as climate change projections. Crucially, during testing when we are unable to di-
 136 rectly determine a prediction’s accuracy, we need a UQ method that can provide a credible
 137 *confidence level* for each prediction, serving as a reliable indicator of its accuracy. During
 138 inference, the output of an NN can be inaccurate for various reasons, including poor approx-
 139 imation (e.g., due to poor NN architecture), poor within-distribution generalization (e.g.,

140 for inputs that are rare events), or poor optimization (collectively referred to as *epistemic*
 141 *uncertainty*), as well as because of OOD generalization errors due to input samples from
 142 a distribution different from that of the training set (Abdar et al., 2021; Lu et al., 2021;
 143 Krueger et al., 2021; Miller et al., 2021; Shen et al., 2021; D. Wu et al., 2021; Ye et al., 2021;
 144 D. Zhang et al., 2021; Subel et al., 2023). Quantifying the level of uncertainty would then
 145 allow us to avoid using a data-driven parameterization scheme when it is inaccurate due to
 146 one of the aforementioned reasons (Maddox et al., 2019; Zhu et al., 2019; Li et al., 2022;
 147 Psaros et al., 2023). In the context of data-driven parameterization in climate modeling,
 148 the two most challenging sources of uncertainty are rare/extreme events and OOD gener-
 149 alization errors. The latter is a concern, particularly when the GCM is used for climate
 150 change studies (see below for more discussions).

151 Developing UQ methods for NNs is an active area of research in the ML community,
 152 and there is not a generally applicable rigorous method yet. For instance, techniques like
 153 Markov-Chain Monte Carlo can be prohibitively expensive, especially when dealing with
 154 high-dimensional systems (Oh et al., 2005; Ballnus et al., 2017; Chen & Majda, 2019). For a
 155 comprehensive review in the context of scientific ML, refer to Psaros et al. (2023). The topic
 156 has also started to increasingly gain attention in the climate literature (Guillaumin & Zanna,
 157 2021; Gordon & Barnes, 2022; Haynes et al., 2023; Barnes et al., 2023). In this study, we will
 158 assess the performance of three common UQ methods (Bayesian, dropout, and variational
 159 NNs) by analyzing the relationships between uncertainty and accuracy during inference
 160 testing. We will also consider scenarios involving OOD generalization errors resulting from
 161 global warming.

162 As already mentioned above, OOD generalization (extrapolation to a test data distri-
 163 bution different from that of the training set) is a major challenge for applications involving
 164 non-stationarity, like a changing climate. Studies have already shown that the lack of OOD
 165 generalization in data-driven parameterizations leads to inaccurate and unstable simula-
 166 tion (Rasp et al., 2018; O’Gorman & Dwyer, 2018; Chattopadhyay, Subel, & Hassanzadeh,
 167 2020; Guan et al., 2022; Nagarajan et al., 2020). A general and powerful method for im-
 168 proving the OOD generalization capability of NNs is transfer learning (TL), which involves
 169 re-training a few or all of the layers of a NN using a small amount of data from the new
 170 system (Yosinski et al., 2014). This approach has already shown remarkable success in
 171 enabling data-driven parameterization schemes to extrapolate across the parameter space
 172 (e.g., to $100\times$ higher Reynolds number) in canonical test cases (Chattopadhyay, Subel, &
 173 Hassanzadeh, 2020; Subel et al., 2021; Guan et al., 2023; Subel et al., 2023; C. Zhang et al.,
 174 2023). In particular, Subel et al. (2023) introduced SpArK (Spectral Analysis of Regression
 175 Kernels and Activations) showing that re-training even one layer can lead to successful OOD
 176 generalization, although this optimal layer, unlike the rule of thumb in the ML literature,
 177 may not be the deepest but the shallowest hidden layer. Here, we further leverage these
 178 studies and show how TL can enable OOD generalization of data-driven parameterization
 179 schemes in state-of-the-art GCMs.

180 The methods used in this study and the learned lessons apply to a broad range of
 181 processes and applications in climate modeling. However, the results are presented for a
 182 single test case, that is based on the emulation of complex physics-based GWP schemes in
 183 version 6 of the Whole Atmosphere Community Climate Model (WACCM), a state-of-the-
 184 art GCM (Gettelman et al., 2019). Here, we use the emulations of current physics-based
 185 parameterization schemes as a stepping stone towards learning data-driven parameteriza-
 186 tions from observations and high-fidelity simulations by testing ideas for addressing items
 187 1-3 listed earlier. Furthermore, developing better representations of un- and under-resolved
 188 GWs in GCMs is an important problem on its own (Kim et al., 2003; Alexander et al., 2010;
 189 Achatz, 2022). A number of recent studies have taken the first steps in learning data-driven
 190 GWP from observations and high-resolution simulations (Matsuoka et al., 2020; Amiramjadi
 191 et al., 2022; Sun et al., 2023; Dong et al., 2023), though careful and time-consuming steps
 192 are needed in producing, analyzing, and using such data. Furthermore, two recent stud-
 193 ies focused on emulators of simpler GWP schemes in a forecast model and idealized GCM
 194 have readily shown the usefulness of lessons learned from emulators (Chantry et al., 2021;
 195 Espinosa et al., 2022; Hardiman et al., 2023). This further motivates the focus on using
 196 emulators for testing ideas for addressing data imbalance, UQ, and OOD generalization.

197 This paper is structured as follows. Section 2 introduces the WACCM simulations and
 198 the NN architectures used in this study. The findings, detailed in Section 3, emphasize
 199 the insights gained in addressing data imbalance and UQ, alongside OOD generalization of
 200 the emulators under warmer climate conditions. Consistent with Chantry et al. (2021), we
 201 find that using an NN to emulate the parameterization of orographic GWs is significantly
 202 more challenging than non-orographic GWs. This necessitated additional steps to achieve
 203 reasonable offline performance, as detailed in Section 4. To the best of our knowledge, this
 204 stands as the first NN-based emulation of orographic GWs to address the challenges in
 205 Chantry et al. (2021). Finally, we provide a concluding summary in Section 5.

206 **2 Data and Methods**

207 **2.1 The Whole Atmosphere Community Climate Model (WACCM)**

208 The NCAR’s WACCM version 6 introduced in Gettelman et al. (2019) is used in this
 209 study. WACCM has state-of-the-art GWP schemes for GWs from three different sources:
 210 orography (OGWs), convection (CGWs), and fronts (FGWs). These complex sources make
 211 the emulation of the GWP schemes in WACCM a challenging task. This is, therefore, a
 212 suitable test case to investigate ideas for learning rare events, UQ, and OOD generalization to
 213 benefit the future efforts for the much more complex task, that is learning data-driven GWP
 214 schemes from observations and/or high-resolution GW-resolving simulations (Amiramjadi
 215 et al., 2022; Sun et al., 2023).

216 The configuration of the WACCM used in this study is identical to the public version in
 217 Gettelman et al. (2019), with a horizontal resolution of $0.95^\circ \times 1.25^\circ$ and 70 vertical levels.
 218 The two non-orographic GWP schemes in WACCM both follow Richter et al. (2010), yet

219 allow separate specifications of FGW and CGW sources. For OGWs, WACCM uses an up-
 220 dated planetary boundary layer form drag scheme from Beljaars et al. (2004), near-surface
 221 nonlinear drag processes following Scinocca & McFarlane (2000), and a ridge-finding algo-
 222 rithm to define orographic sources based on Bacmeister et al. (1994). A full documentation
 223 of WACCM OGWs can also be found in Kruse et al. (2022).

224 We conduct two sets of simulations: A 10-year pre-industrial “control” run, and a
 225 10-year pseudo-global-warming “future” run with $4\times\text{CO}_2$ and uniform +4 K sea-surface
 226 temperature increases. In each run, we save, on the native grid, all the inputs and outputs
 227 for each of the three GWPs every 3 hours to capture the diurnal cycle. A complete list of
 228 these inputs/outputs, which are used in the training of the NN-based emulators, is presented
 229 in Appendix A.

230 We train separate NNs for emulating the three GWP schemes that have different
 231 sources. We use the first 6 years of the control run for training and the last 4 years for
 232 validation (years 7 and 8) and testing (years 9 and 10). With a grid resolution of $\sim 1^\circ$,
 233 there are 55,296 columns for each time snapshot, resulting in approximately 960 million
 234 input/output columns during the 6-year training period. Given the strong temporal cor-
 235 relation between the 3-hourly samples, we perform sub-sampling on both the training and
 236 validation data to reduce the dataset size. To accomplish this, we begin by shuffling all
 237 the input/output column pairs in time at each latitude/longitude grid point. Then, we
 238 randomly select 2,000 input/output pairs at each location for training and 500 pairs for
 239 validation.

240 To give the readers a general idea of the parameterized GWs and large-scale circulation
 241 in WACCM, Figure 1 shows the zonal-mean climatology for zonal GW drag/forcing, here-
 242 inafter referred to as GWD, arises from the divergence of gravity wave momentum transport
 243 (fluxes), from all three sources, computed from the 6-year training period in the control run.
 244 The zonal-mean zonal wind climatology is also shown. Seasonal dependency for both the
 245 GWD and the circulation is observed in the simulations. At levels below 100 hPa, the ten-
 246 dencies of non-orographic GW are relatively small compared to those from OGWs; however,
 247 their amplitudes increase significantly at higher altitudes. While the parameterized effect of
 248 GWs is generally to decelerate the zonal flow, there are exceptions, notably in regions like
 249 the equatorward flanks of the stratospheric polar night jets, where FGWs can accelerate the
 250 flow. For more information on the GWP schemes and circulations in WACCM, see Garcia
 251 et al. (2017) and Gettelman et al. (2019).

252 **2.2 The NNs and UQ**

253 ***2.2.1 The Deterministic Fully Connected NN***

254 Here we briefly describe the general structure of the NN-based regression models trained
 255 as emulators for GWP schemes. For the deterministic artificial NN, denoted as ANN in this
 256 study, we use multilayer perceptrons (MLP). MLPs, which are feedforward fully connected
 257 NNs, take inputs through successive layers of linear transformation and non-linear activation

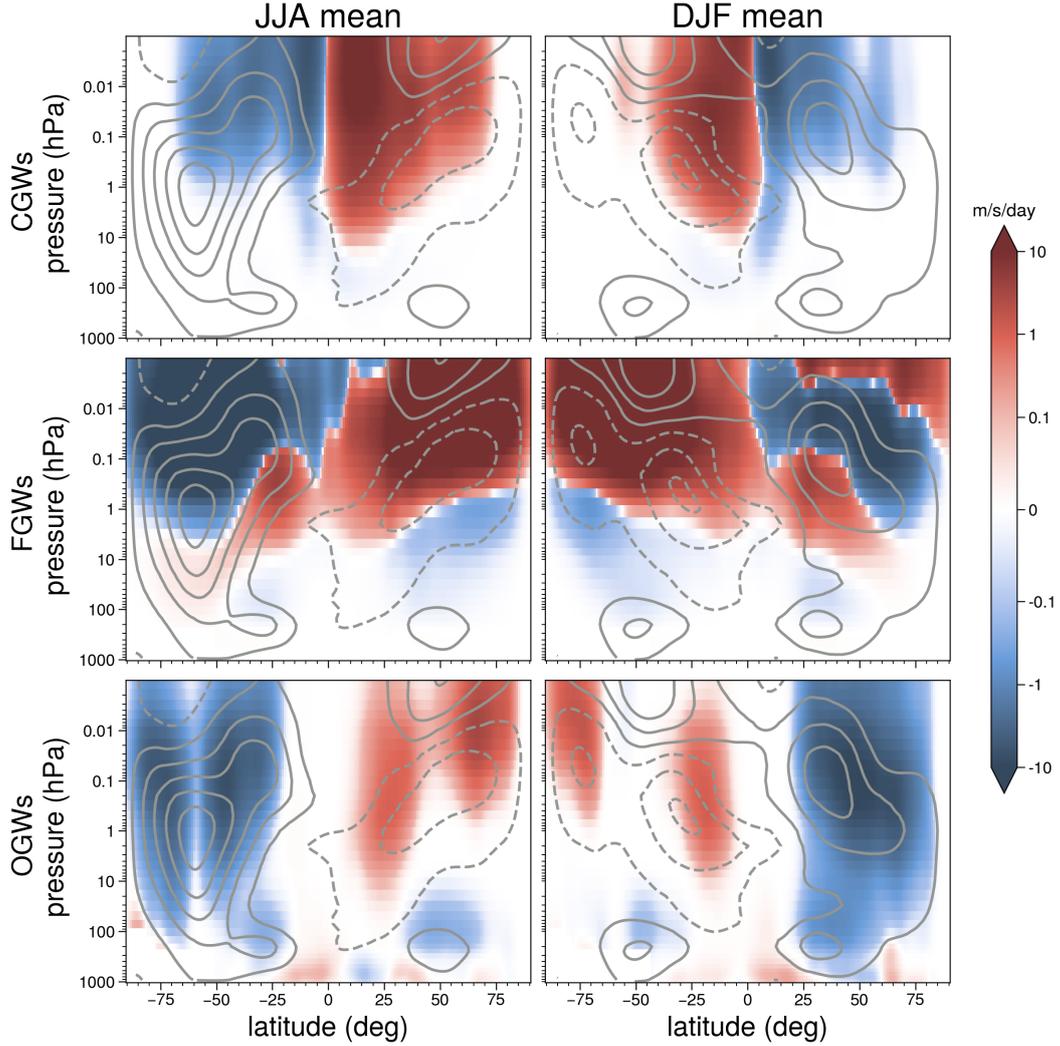


Figure 1. Climatology of zonal-mean GWD during summers (JJA) and winters (DJF) from all 3 sources in the control (pre-industrial) WACCM simulations. Top: CGWs; middle: FGWs; bottom: OGWs. The climatology of the zonal-mean zonal wind is also shown (grey lines), with an interval of 20 m/s. Dashed lines indicate negative values. Zero lines are omitted.

258 functions to produce an output, so as to learn a functional relationship between the input
 259 and output (Figure 2a). Deep MLPs have multiple layers of weights, which are optimized
 260 over many samples of input-output data pairs. Such MLPs are thus very powerful in terms
 261 of learning complicated functional relationships. Generally, we can write the governing
 262 equations of an MLP as

$$z^\ell = \sigma(W^\ell z^{\ell-1} + b^\ell), \quad (1)$$

263 where z^ℓ is the activation (output) of layer ℓ , W^ℓ is the weight matrix connecting layers ℓ
 264 and $\ell - 1$, and b^ℓ is the bias at layer ℓ , which allows the network to fit the data even when
 265 all input features are equal to 0. σ is the non-linear activation function.

In this study, we employ the same NN structure while training three distinct NNs, each for GWP originating from one of the three unique GW sources. The input layer contains the same input variables (see Appendix A) used by the WACCM GWPs across all vertical levels. There are 10 hidden layers in total (Figure 2a), and there are 500 neurons in each hidden layer. In the output layer, both zonal and meridional GWD are predicted. The activation function in each layer, σ , is chosen to be swish (Ramachandran et al., 2017), except for the output layer, where it is linear. During training, W^ℓ and b^ℓ are randomly initialized and learned by minimizing a loss function using an ADAM optimizer, with a fixed learning rate of $\alpha = 0.0001$. One of the loss functions used here is the common MSE, i.e.,

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n \left\| \text{NN}(x_i, \Theta) - y_i \right\|_2^2 \quad (2)$$

266 Here, n is the number of training samples and $\|\cdot\|_2$ is the L_2 norm. For training sample
 267 i , vector x_i contains all the inputs to the NN (Appendix A), vector y_i contains the true
 268 zonal and meridional GWD at each vertical level, and $\Theta = \{\theta_j\}_{j=1\dots p}$ denotes the trainable
 269 parameters, i.e., the weights ($p \approx 3 \times 10^6$).

270 **2.2.2 The UQ Methods and Metrics**

271 Although deterministic NNs are powerfully expressive and can exhibit high out-of-
 272 sample predictive skills, they do not provide estimates of the uncertainty associated with
 273 their predictions. As mentioned earlier, currently there is no rigorous method to estimate
 274 the uncertainty of an NN prediction. That said, a variety of techniques have been developed
 275 for UQ in NNs, though the validity and usefulness of the estimated uncertainty for scientific
 276 applications remain subjects of ongoing investigations (e.g., Psaros et al., 2023; Haynes et
 277 al., 2023). In this paper, we use three different and widely used approaches to perform UQ
 278 from the ML literature: Bayesian neural network (BNN), dropout neural network (DNN),
 279 and variational auto-encoder (VAE). A brief overview of these approaches is provided below.

280 *Bayesian neural network (BNN):* A BNN combines the deterministic NN described
 281 earlier and in Figure 2a with Bayesian inference (Blundell et al., 2015). Simply speaking, a
 282 BNN estimates distributions of the weights, rather than point values (as in a deterministic
 283 NN). The posterior distributions in the BNN (i.e., the distributions of the weights and
 284 biases) are calculated using the Bayes rule. In this study, we follow the standard practice
 285 and assume that all variational forms of the posterior are normal distributions. Furthermore,
 286 to accelerate the training process, we use the normal distribution $\mathcal{N}(\mu, 1)$ for all the priors in
 287 the BNN (where μ is obtained from parameters of the trained deterministic NN). Note that
 288 while we are assuming normal distributions for the trainable parameters, the predictions
 289 generated by BNN can fit different distributions due to the use of nonlinear activation
 290 functions. The resulting distribution of the predictions during inference gives an estimate
 291 of their uncertainty.

292 *Dropout neural network (DNN):* A DNN is developed by randomly eliminating all out-
 293 going connections from some of the nodes (Figure 2a) in each hidden layer of a deterministic

294 NN during the training and the inference (Srivastava et al., 2014). The fraction of nodes
 295 “dropped” on average in each layer is called the dropout ratio. Mathematically, Equation (1)
 296 can be reformulated for a DNN as:

$$z^\ell = \sigma(D^\ell W^\ell z^{\ell-1} + b^\ell), \quad (3)$$

297 where the dropout matrix D^ℓ is a square diagonal binary matrix of integers 0 or 1. The
 298 diagonal elements of D^ℓ follow a Bernoulli distribution where the probability of zero is the
 299 dropout ratio.

300 Dropout was initially developed as a regularization technique to prevent over-fitting in
 301 NNs. However, Gal & Ghahramani (2016) showed that training a NN with the dropout
 302 technique approximates a Bayesian NN. In this study, we use a dropout rate of 0.1, which
 303 is incorporated in all hidden layers, but we also investigate the sensitivity of the DNN to
 304 different dropout rates, as later shown in Appendix B. Note that the random dropping out
 305 is also used during inference, leading to a distribution for each prediction.

306 *Variational auto-encoder (VAE)*: A typical VAE (Kingma & Welling, 2014) consists
 307 of two NNs (Figure 2b): an encoder that transforms the input into a lower-dimensional
 308 latent space, parameterized by a normal probability distribution, and a decoder that inverts
 309 this transformation and produces the original input. The difference between the decoder’s
 310 output and the original input drives the learning process of the encoder and decoder, while
 311 the parameterized lower-dimensional latent space provides the uncertainty of this transfor-
 312 mation. The VAE was developed for generative reconstructions of data by simply drawing
 313 samples from the latent space. The VAE is basically a dimension-reduction method. Many
 314 variants, however, have been proposed for more specific purposes. In this study, following
 315 Foster et al. (2021), we add a third NN, as illustrated in Figure 2b, that randomly draws
 316 samples from the parameterized latent space as inputs, and predicts the zonal and merid-
 317 ional GWDs as outputs. The difference between the predicted GWDs and the true GWDs
 318 drives the learning of the third NN. Consequently, the loss for the entire network consists
 319 of three components: the loss between the reconstructed input and the original input, the
 320 Kullback–Leibler (KL) divergence between the distribution of the latent space and a nor-
 321 mal distribution, and the loss between the predicted GWDs by the third NN and the true
 322 GWDs.

323 For a specific input, each of these three UQ methods discussed above can be run multiple
 324 times, generating an ensemble of predictions with different realizations of the weights by
 325 drawing from the trained distribution. This is in contrast to the deterministic NN that
 326 provides just a single-valued prediction for a given input. These ensembles can then be used
 327 to quantify the uncertainty associated with that prediction. We expect that the RMSE of
 328 the ensemble mean should exhibit approximately a 1-1 relationship with the ensemble spread
 329 (i.e., the standard deviation of the ensemble members). To investigate this relationship, we
 330 use the spread-skill plot (Delle Monache et al., 2013). Detailed calculations behind the
 331 spread-skill plot can be found in Appendix C, where we also introduce two metrics: spread-

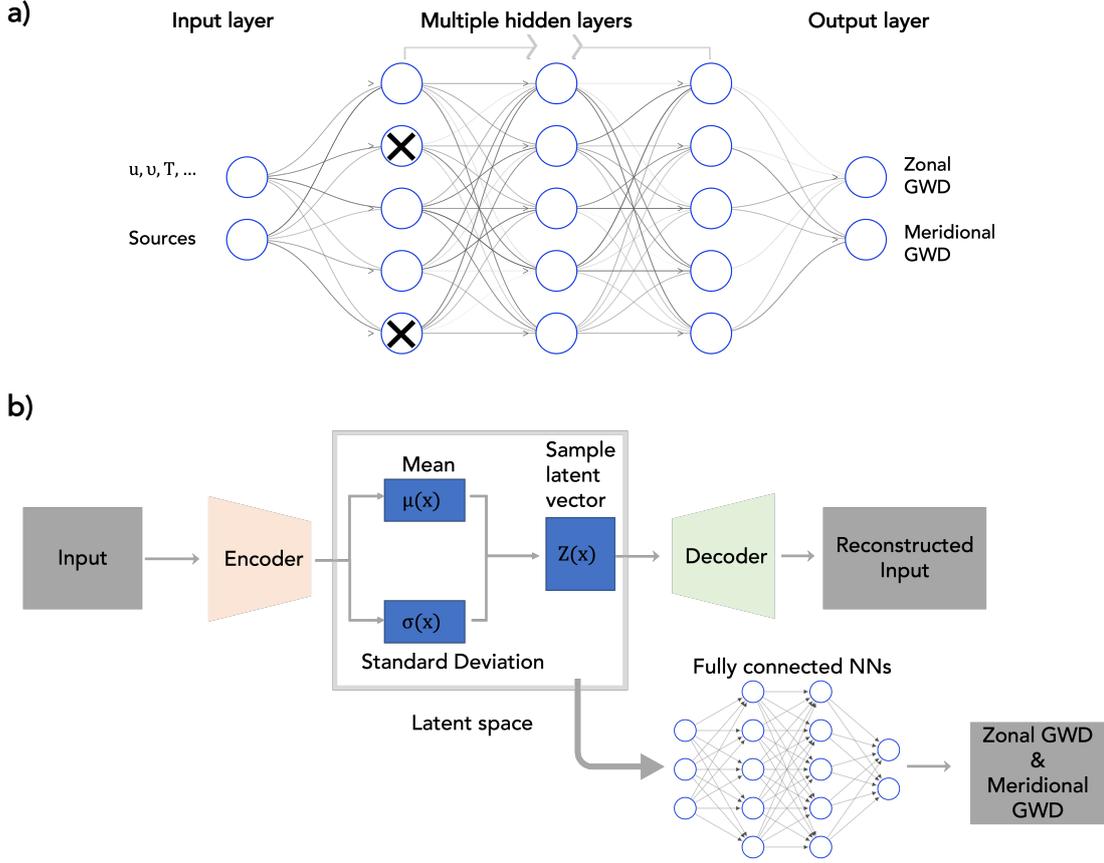


Figure 2. Schematics of the NN-based emulators and different training/re-training strategies used in this study. (a) Schematic for the MLP and DNN. The inputs of the NN are connected through successive layers of neurons (blue circles) to the output (GWDs). A fully connected MLP NN is trained from randomly initialized weights and biases in all layers. A DNN is the same but some connections are randomly eliminated during training and inference (black crosses). In TL, only some of the layers of a previously trained MLP are re-trained using new data. (b) Schematic for the VAE. A low-dimensional latent space is constructed and then used as the input for the additional fully connected NNs, which is similar to the one in (a).

332 skill reliability (SSREL) and overall spread-skill ratio (SSRAT), both of which summarize
 333 the information presented in the spread-skill plot.

334 2.3 Transfer Learning

335 Transfer learning refers to leveraging/reusing information (weights) from an already
 336 well-trained base NN to effectively build a new NN for a different system from which only a
 337 small amount of training data is available (Yosinski et al., 2014; Tan et al., 2018; Chattopad-
 338 hyay, Subel, & Hassanzadeh, 2020). For our purpose, which is improving OOD generalization
 339 to the warmer climate, the TL procedure is as follows. For any of the NNs described earlier
 340 (e.g., the one in Figure 2a), we train them from randomly initialized weights and biases
 341 with data from the control simulations. The NN will work well during inference for test

342 samples from the control but not from future (warmer climate) simulations (as shown in
 343 the Results section). To address this, TL is applied wherein most of the NN’s weights are
 344 kept constant, and only one or two hidden layers are re-trained using a limited dataset from
 345 the future simulation. Although this small dataset is insufficient for training an entire NN
 346 from random initialization, careful and correct selection of hidden layers for re-training, as
 347 discussed in Subel et al. (2023), allows the development of an NN that accurately adapts to
 348 the new system, i.e., the future climate conditions.

349 Here, we re-train the NN-based emulator that was initially trained on the control data
 350 with new data from only 1 month (30 consecutive days) of integration (1.4% of 6 years
 351 simulation for the initial training) of WACCM model under future forcing ($4\times\text{CO}_2$). We
 352 have explored different choices of layers to re-train with the same amount of new data and
 353 found that re-training the first hidden layer yields the best results, consistent with Subel et
 354 al. (2023). Therefore, the results with only re-training the first hidden layer are shown in
 355 Section 3 unless stated otherwise.

356 **3 Results**

357 **3.1 Data Imbalance**

358 As discussed earlier, the physics-based GWP schemes in WACCM are directly linked to
 359 their sources. This means they only produce non-zero values when their respective sources
 360 are active. For example, in a specific grid box, CGWs only register non-zero values when
 361 there is active convection within that box. The heterogeneous and sometimes intermittent
 362 nature of these sources leads to a dataset that is significantly imbalanced. Figure 3 shows
 363 global maps of the occurrence frequency of non-zero GWD for CGWs and FGWs. On
 364 average, only 7.6% of all GCM columns yield non-zero CGWs, primarily located in the
 365 tropics. Similarly, for FGWs, only 8.5% of all columns have non-zero outputs, but unlike
 366 CGWs, the majority of these are located in mid-to-high latitudes, particularly along the
 367 storm track region. For the OGWs in WACCM, data imbalance presents a greater challenge,
 368 to be discussed in a later section. While it is possible to simply separate zero and non-zero
 369 columns for emulation work where we know the truth, this approach falls short with real-
 370 world data, which is the main purpose of this study.

371 In addition to their sources, several other factors specific to GWD data exacerbate
 372 the data imbalance problem. In the case of each GCM column with non-zero GW activity,
 373 momentum fluxes are generally concentrated at a few critical height levels rather than being
 374 smoothly distributed throughout the entire column. This further restricts the effective
 375 occurrence frequency of non-zero values. Moreover, GWs exhibit significant intermittency,
 376 where a small portion of large-amplitude GWs often dominates the morphology of the
 377 observed global GW momentum flux distribution (Hertzog et al., 2012; Geller et al., 2013).
 378 Therefore, it is crucial for NNs to not only accurately identify the columns that produce
 379 GWDs but also to effectively learn and recognize rare and extreme GWDs.

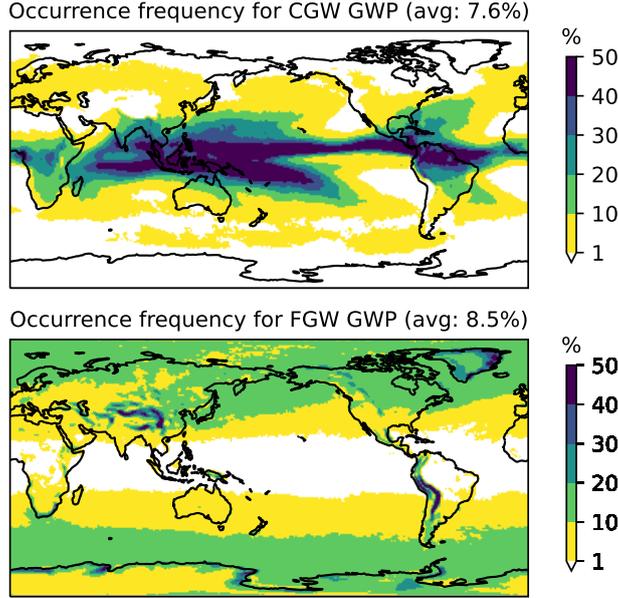


Figure 3. Distribution of occurrence frequency for CGWs (top) and FGWs (bottom) in the WACCM pre-industrial control simulations, based on the average of the 6-year training dataset.

380 Given the complexity of the GWD dataset, different normalization methods are con-
 381 sidered in this study. The first method, dubbed “NORM1”, is the typical normalization
 382 used in ML practices, which calculates elemental means and standard deviations for each
 383 feature (i.e., input variable at a given model level) and normalizes both inputs and outputs
 384 by these values (e.g., Espinosa et al. (2022)). With this approach, the same relative changes
 385 in wind at each level are treated equally in the input. The loss function in Equation (2) also
 386 penalizes the same relative error in GWD at each level equally. The second method, referred
 387 to as “NORM2“, is designed with the physics of GWD in mind. For the velocity inputs
 388 (u, v) and the tendency outputs (GWD), each column is normalized by one single value,
 389 which is the largest standard deviation from all model levels. Additionally, the mean values
 390 for these variables, are retained (e.g., $u_{norm2}(x, y, z, t) = u(x, y, z, t) / \max(std(u))$). Un-
 391 like NORM1, the original wind profile’s structure is preserved in NORM2, and large GWD
 392 values at certain heights maintain a relatively larger value after this normalization. For all
 393 other input variables, NORM2 is identical to NORM1. Compared to NORM1, NORM2
 394 places more emphasis on large GWD values and penalizes the NN more for missing these
 395 significant tendencies. These two normalization methods are also employed in Chantry et al.
 396 (2021), who found similar performance from these methods with the non-orographic GWPs.

397 Figure 4 shows the performance of the emulations for CGWs with the two normal-
 398 ization methods. When employing NORM1, the conventional approach seen in prior ML
 399 practices, and also our initial attempts, the emulator’s performance is poor. Although the
 400 NN demonstrates some skill, its predictions tend to cluster around zero. However, when the
 401 second normalization method (NORM2) is employed, the emulation results show significant

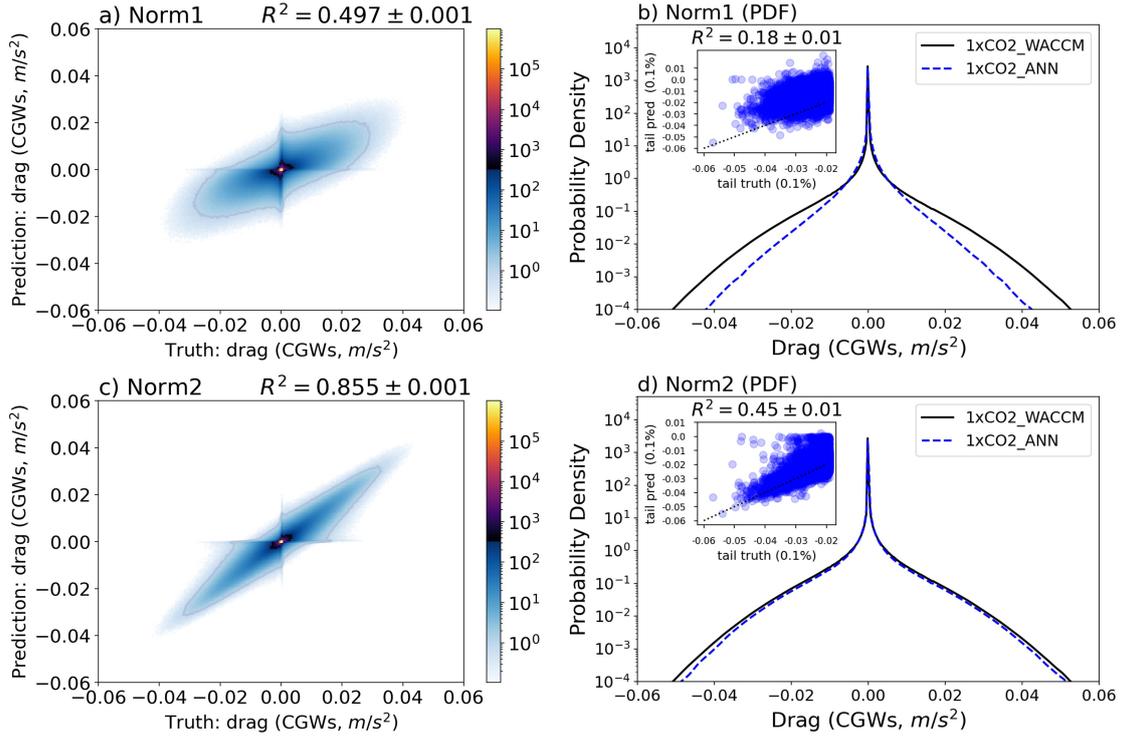


Figure 4. Data imbalance for GWD due to CGWs and the emulation results with two different normalization methods. a) A 2D histogram displaying the emulated GWD due to CGWs and the truth, with the training dataset normalized using NORM1; b) Distribution of the original convective GWD (black line) and the predicted values (blue line) with NORM2. The scatter plot in the corner represents the tail part only, including points with the top 0.1% amplitudes; c) Similar to a), but for NORM2; d) Similar to b), but for NORM2. The R^2 uncertainty range is estimated by dividing the test data into 10 segments, calculating the metric for each segment, and then computing the standard deviation (STD).

402 improvement, in contrast to the findings of Chantry et al. (2021). We attribute this improve-
 403 ment to the more pronounced data imbalance in our dataset, and it is likely a consequence
 404 of NORM2’s emphasis on modeling the large GWD values. Nonetheless, emulating the tail
 405 of the probability density function (PDF) (rare events) remains poor, as evidenced by the
 406 tails in Figure 4c, primarily due to the predominance of zero GWD columns in the training
 407 dataset. To more effectively address the data imbalance issue in these regression tasks, we
 408 further propose two approaches here:

- 409 1. Resampling the data (ReSAM): In this approach, we limit the number of training
 410 sample pairs with zero GWD to be equal to the number of samples with non-zero
 411 GWD. This significantly reduces the number of columns with zero GWD, thus mit-
 412 igating the data imbalance issue. Additionally, this sub-sampling reduces the total
 413 size of the training dataset, which, in turn, enhances the training speed (approx-
 414 imately sevenfold). While resampling methods have been well-established in the ML

415 literature, they have mainly been used for classification problems. Their application
 416 to regression problems in climate research has not been extensively explored.

417 2. Weighted loss function (WeLoss): Instead of assigning the same weight to all sample
 418 pairs in the loss function, we modify the weight for each column based on the PDF
 419 of its maximum GWD amplitude. This adjustment allows us to re-formulate the loss
 420 function defined in Equation (2) as

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{W}_i \{ \mathbf{NN}(x_i, \Theta) - y_i \} \right\|_2^2 \quad (4)$$

where

$$\mathbf{W}_i = \frac{1}{PDF(max(|y_i(z)|))} \quad (5)$$

421 Note that, in practice, we lack knowledge of the precise PDF for the maximum GWD
 422 within each column. Therefore, we employ a histogram with 20 bins as an alterna-
 423 tive. Given the fact that large-amplitude GW events are rare, the WeLoss approach
 424 incentivizes the NN to prioritize these significant events.

425 When we apply the ReSAM approach to balance the training dataset (after normal-
 426 ization with NORM1 or NORM2), the emulation results significantly improve, as shown in
 427 Figure 5. In fact, when considering the R-squared value between the NN prediction and the
 428 ground truth, the ReSAM approach with NORM2 yields the best results. However, as the
 429 training dataset is still predominantly composed of zeros and small GWD values due to the
 430 intermittence of the GWs, examining the emulation results for only large amplitude GW
 431 events (e.g., the top 0.1% in Figure 5d) reveals less satisfactory performance ($R^2 = 0.72$).
 432 Regarding the WeLoss approach, it has a more limited impact on improving the R-squared
 433 value of the emulation (as shown in Figure 5e). However, it proves valuable in capturing
 434 the tails of the PDF and, thus, rare events (as depicted in Figure 5f). Moreover, as ReSAM
 435 and WeLoss represent distinct operations, they can be effectively combined when construct-
 436 ing a NN. The result of this combined approach for emulating the CGWs can be found in
 437 Figures 5g and 5h. While the R-squared value for the entire distribution only marginally
 438 changes (0.925 vs. 0.931 with ReSAM only), the performance of the emulation for the tail
 439 part has been improved (R^2 increased to 0.77).

440 Similarly, Figure 6 presents the offline emulation results for the FGWs. The conclusions
 441 drawn for CGWs generally hold true. However, data imbalance in FGWs is less pronounced
 442 compared to CGWs, which simplifies the task of emulating FGWs. Even without any
 443 resampling or changes to the normalization or (see Figure 6a), we achieve reasonable emu-
 444 lation results ($R^2 = 0.9$). One contributing factor is the wider spatial distribution of FGWs
 445 compared to CGWs (refer to Figure 3). Additionally, the source of FGWs (frontogenesis
 446 function) in WACCM exhibits a much more continuous nature compared to precipitation
 447 and diabatic heating. As the data imbalance issue is less severe for FGWs, the performance
 448 with different normalization methods becomes more similar, echoing findings from Chantry
 449 et al. (2021) who emulated non-orographic GWs (including convective and frontal GWs)
 450 together.

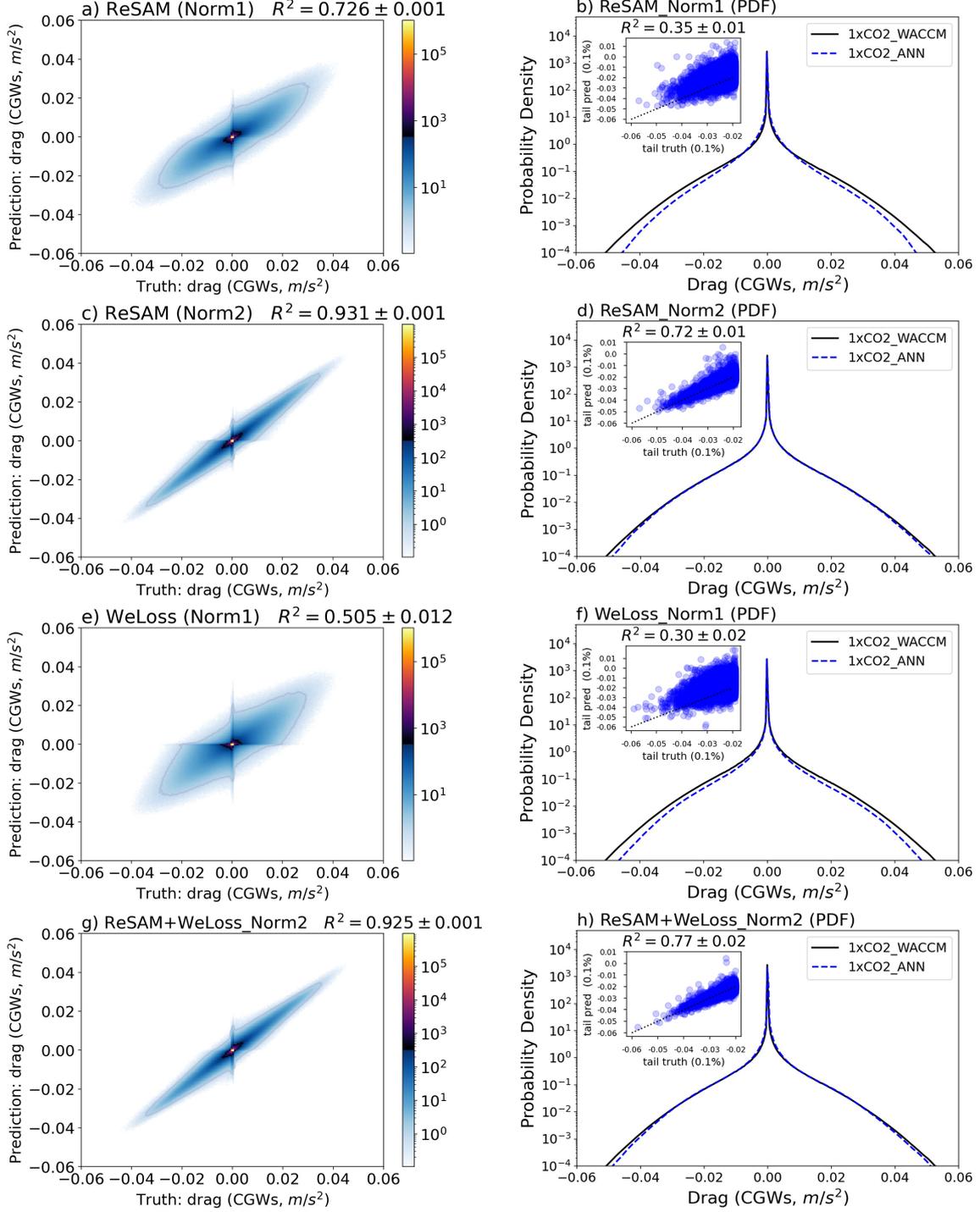


Figure 5. Similar to Figure 4, but for CGWs with the proposed ReSAM and WeLoss methods. a) A 2D histogram for the emulation with resampled data (ReSAM) after using Norm1; b) Distribution of the emulated GWD due to CGWs similar to Figure 4b, but with ReSAM applied; c) Similar to a), with training data normalized using Norm2; d) Similar to b), with training data normalized using Norm2; e) Similar to a), but with the WeLoss approach; f) Similar to b), but with the WeLoss approach; g) Similar to c), after applying both ReSAM and WeLoss methods together; h) Similar to d), after applying both ReSAM and WeLoss methods.

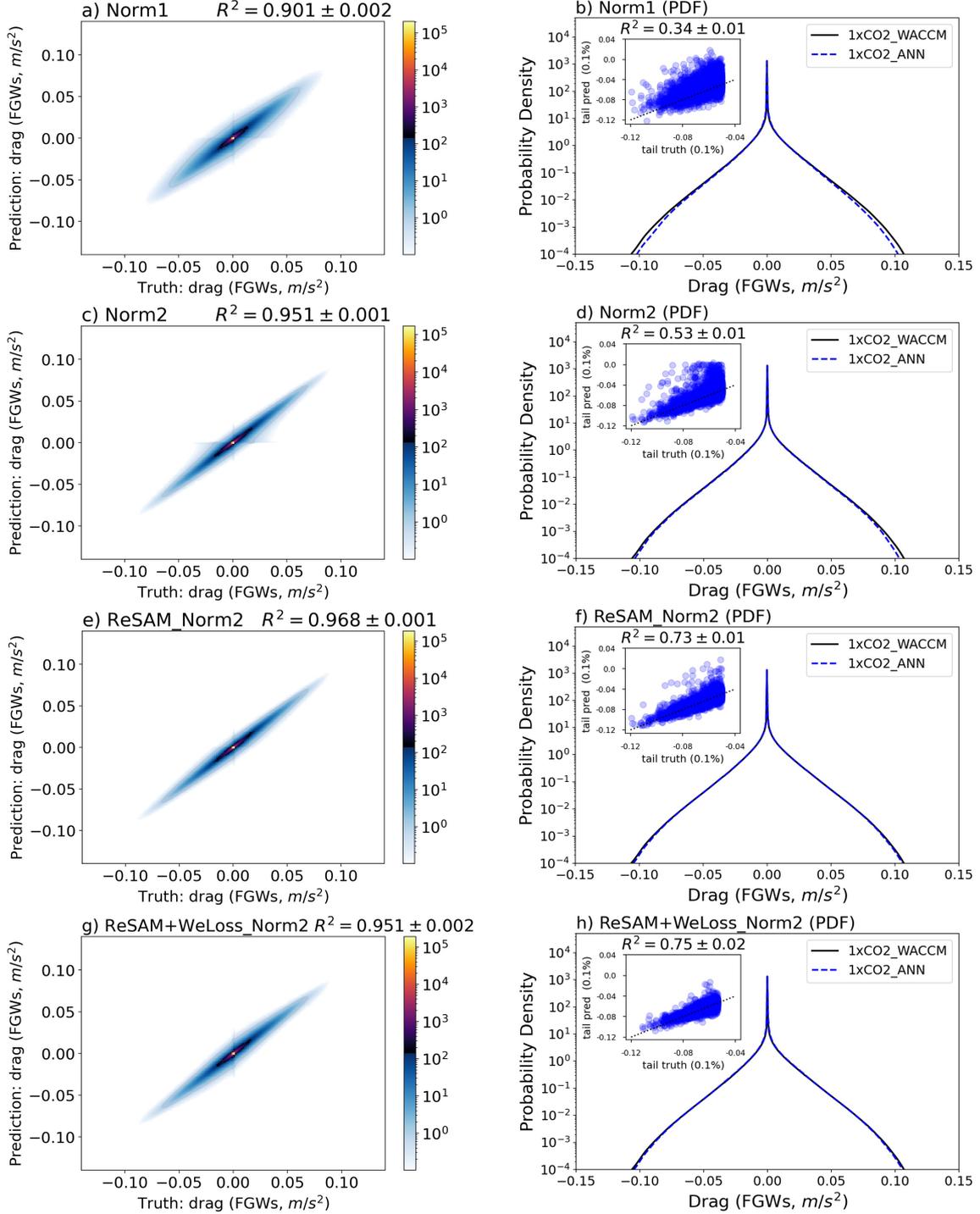


Figure 6. Similar to Figure 5, except for FGWs. a) A 2D histogram for the emulation using Norm1, without ReSAM or WeLoss; b) Distribution of the GWD due to FGWs with NORM1, similar to Figure 4b; c) Similar to a), with training data normalized using Norm2; d) Similar to b), with training data normalized using Norm2; e) Similar to c), but with the ReSAM approach; f) Similar to d), but with the ReSAM approach; g) Similar to Figure 5g, applying both ReSAM and WeLoss methods to the FGWs; h) Similar to Figure 5h, applying both ReSAM and WeLoss methods to the FGWs.

451 In summary, data imbalance can pose challenges when learning from data that closely
 452 resembles real-world data (further discussed in the subsequent section on emulating OGWs).
 453 Proper resampling techniques can significantly enhance the NNs' performance by improving
 454 dataset balance. Furthermore, modifying the loss function to penalize the NNs more for
 455 missing extreme values can further improve performance at the tails of the PDF. For the
 456 remainder of the paper, unless otherwise specified, we continue to employ the ReSAM
 457 approach and the standard loss function with NORM2 unless stated otherwise.

458 3.2 Uncertainty Quantification

459 As outlined in subsection 2.2.2, we employ three different methods (i.e., BNN, DNN,
 460 and VAE) to quantify the uncertainty of predictions during inference (testing). For this
 461 purpose, an ensemble of 1000 members is generated by running each UQ-equipped NN 1000
 462 times for each input from the testing set. Figure 7 presents sample profiles of zonal GWD
 463 derived from the deterministic NN (ANN) and the three UQ-equipped NNs, alongside the
 464 true GWD profiles from WACCM. Note that these examples have not been used in the
 465 training or validation process. It is evident from the figure that all three UQ-equipped
 466 NNs show reasonable skill in predicting the complex profiles of GWD due to CGWs and
 467 FGWs (also reflected in R-squared in Table 1), albeit with a slight decrease in accuracy
 468 compared to ANN. As discussed earlier, a valuable uncertainty estimate should correspond
 469 closely with the NN's test accuracy, providing insights into when to trust the NN's pre-
 470 diction during inference. Such a relationship can be seen in a few randomly chosen GWD
 471 profiles that's shown in Figure 7. In each pair of CGW and FGW profiles, the left column
 472 shows the estimated uncertainty is also low when the prediction error is low, indicating the
 473 NN's confidence in its accurate predictions. In contrast, the right column, which generally
 474 represents more complex profiles, exhibits the NN's less accurate predictions, and increased
 475 uncertainty, highlighted by the wider confidence intervals.

476 While Figure 7 demonstrates the performance of the UQ methods for just a few GWD
 477 profiles, the spread-skill plots shown in Figure 8 offer a broader perspective based on 60,000
 478 profiles, following the calculations detailed in Appendix C. It is evident from the plots
 479 that all three UQ methods produce reasonably informative uncertainty estimates, as their
 480 curves closely align with the 1-to-1 line. In the case of CGWs, all data points are above
 481 the 1:1 line, indicating a slight overconfidence (underdispersiveness) across all three UQ
 482 methods, with the DNN being slightly closer to the 1-to-1 line. For the FGWs, the DNN
 483 demonstrates slightly better performance, although it marginally drops below the 1-to-1
 484 line in the first few bins, indicating a slight underconfidence. Notably, it can be seen from
 485 the spread frequency inset that the vast majority of the data points are within the first few
 486 bins, for which both spread and skill values are small, and they are generally closer to the
 487 1-to-1 line.

488 It should also be noted that for the large values of model spread (\overline{SD}), there is only a
 489 very limited number of data points, as is evident from the inset histograms. Consequently,

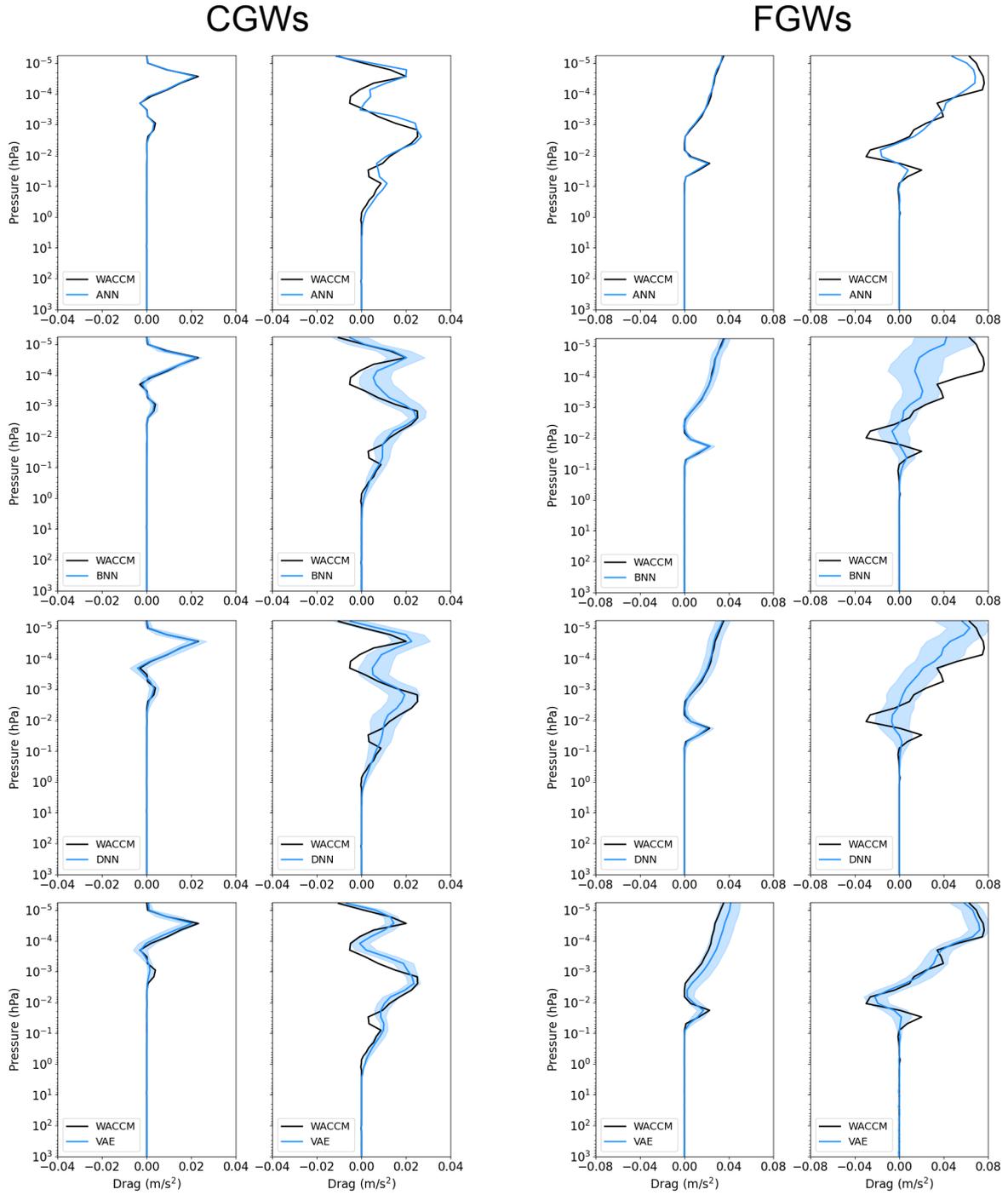


Figure 7. Sample profiles of zonal GWD as predicted by various NNs, as indicated. The true profile is shown by the black line, while the blue solid line represents the mean of 1000 ensemble members. The shaded region indicates the 95% confidence interval. In each pair of CGWs and FGWs profiles, the left column provides examples with low estimated uncertainty, corresponding to instances of low error. Conversely, the right column illustrates cases with high uncertainty when the error is high.

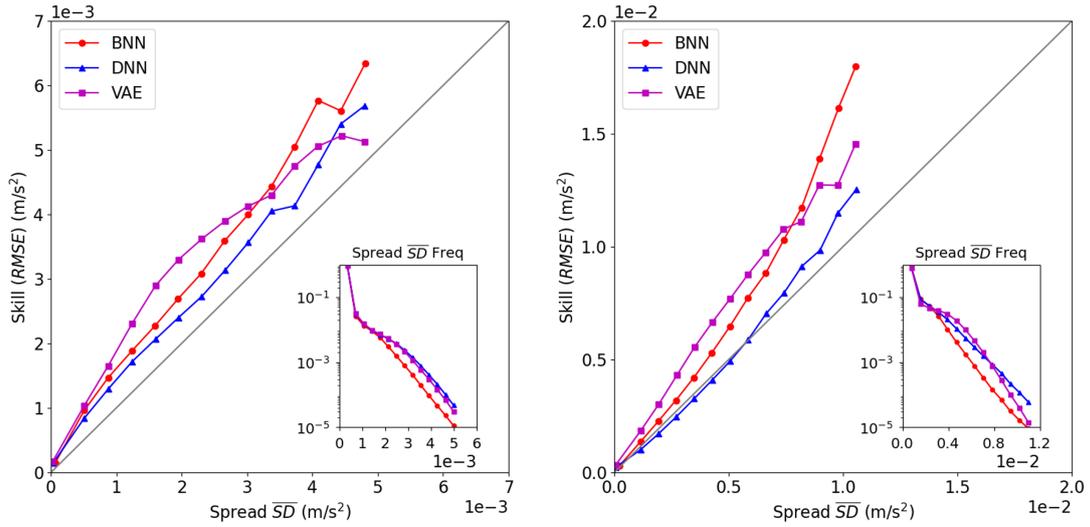


Figure 8. Spread-skill plot for GWD due to (left) CGWs, and (right) FGWs. The diagonal 1:1 line represents the perfect spread-skill line. Points above (below) this line correspond to spread values where the model is overconfident (underconfident). The inset histogram shows how often each spread value occurs. See Appendix C for a detailed discussion on the calculations of the spread-skill plot.

490 the standard deviation (STD) can become a misleading measure of spread because of the
 491 non-normal distributions.

492 To summarize the quality of the spread-skill plots for the three UQ methods, we explore
 493 the metrics introduced in subsection 2.2.2 and Appendix C (see Table 1). The R-squared
 494 value for the ensemble mean prediction is also given to show the accuracy of each UQ
 495 method. Based on SSREL, whose ideal value is zero, BNN shows the best performance for
 496 both CGWs and FGWs. However, if we check SSRAT, where 1 is the optimal number, DNN
 497 is the best among these three methods. This discrepancy can be explained by a closer look
 498 at the Equations (C2) and (C3). SSREL, which is a bin-weighted average difference, is most
 499 sensitive to the performance of the NN in the first bin, where the vast majority of the data
 500 points are located (see the inset histograms in Figure 8), while SSRAT is more influenced by
 501 larger values of spread and skill. Accordingly, the VAE shows the highest values of SSREL,
 502 which is indicative of its sub-optimal performance in the first bin, where there are small
 503 values of spread and skill.

504 In the results presented in Figure 8 and Table 1, each height level of a GWD profile is
 505 considered as an individual sample. A zonal GWD profile, with its 70 vertical levels, thus
 506 constitutes 70 distinct samples. While analyzing these samples offers insights into the NN's
 507 overall performance by averaging statistics across numerous profiles, our primary interest is
 508 often in the uncertainty associated with an individual GWD profile. This uncertainty can
 509 then aid in determining whether to trust/use the NN's prediction for that particular GWD

510 profile. Therefore, we use Equation (C4) to assess the relationship between uncertainty and
 511 test accuracy for each GWD profile. Furthermore, to estimate uncertainty, here we use the
 512 interquartile range (IQR) to reduce the influence of outliers.

513 Figure 9 shows the Gaussian kernel density of spread against RMSE for all 60,000
 514 profiles, as indicated by the color shading. The x -axis represents the IQR of each GWD
 515 profile, while the y -axis denotes its corresponding RMSE. A strong correlation between
 516 the two is observed across all three UQ methods. Consequently, GWD profiles with larger
 517 uncertainties often coincide with larger errors. Figure 9 also shows a close similarity between
 518 BNN and DNN. In contrast, VAE tends to provide marginally larger uncertainties, especially
 519 for FGWs. This is consistent with VAE’s slightly reduced accuracy as indicated in Table
 520 1. Overall, given the monotonic relationship between the uncertainty and test error, these
 521 results show that all three UQ methods provide useful and informative uncertainty for with-
 522 distribution test samples. A user can set a threshold on uncertainty based on their tolerance
 523 for error (RMSE) and decide whether they trust the NN for a given input sample.

524 The results presented so far show the performance of the UQ methods based on the
 525 testing data, i.e., data from the current climate. However, the effective performance of UQ
 526 methods can also be tested (perhaps more meaningfully) on OOD data, e.g., data from a
 527 warmer climate. This is particularly relevant for climate change studies. Accordingly, we
 528 evaluate the performance of these trained NNs with input data from the future climate, as
 529 depicted by the black lines in Figure 9. For FGWs, the spread-skill relationship remains
 530 largely similar, especially for BNN and DNN. This suggests that, based on their uncer-
 531 tainties, we can still reliably estimate the error in the NN predictions for FGWs for the
 532 warming climate. A similar pattern is observed for the VAE, though it exhibits increased
 533 uncertainties and higher errors with OOD data. As shown in a later section, for FGWs, the
 534 NNs generalize to the warmer climate without any further effort.

535 In contrast, for CGWs, given the same level of uncertainty, the error in NN predictions
 536 increases significantly for the OOD data compared to that from the current climate, which
 537 means the spread-skill relationship, especially for the BNN and DNN, fails to generalize to
 538 the OOD data. From this perspective, VAE performs better, showing that for the same
 539 level of uncertainty, the increase in error is not as substantial as in BNN and DNN. The
 540 VAE also yields considerably higher uncertainty estimates for future climate, which may aid
 541 in the detection of OOD data. The observed discrepancies in the performance of the NNs
 542 for CGWs and FGWs hint at different levels of their generalizability, a topic we will delve
 543 into more deeply in the following subsection.

544 In summary, while the three UQ methods provide credible and valuable uncertainty
 545 estimates for the current climate, the BNN and DNN are confidently wrong in estimating
 546 CGWs in a warmer climate although VAE shows some promising results. This problem is
 547 common among various UQ techniques as pointed out in the ML literature: they frequently
 548 show overconfidence when assessed with OOD data (e.g., Ovadia et al., 2019). The optimal
 549 UQ method selection depends on the specific metric of interest and the intended application.
 550 While BNN is more broadly used in the literature and gives the best accuracy, DNN could

Table 1. Evaluation scores for the three UQ methods. See Section 2 for more details.

	CGWs			FGWs		
	BNN	DNN	VAE	BNN	DNN	VAE
SSREL (1e-4)	1.29	1.48	2.14	1.20	1.69	5.21
SSRAT	0.73	0.82	0.72	0.69	0.93	0.69
R-squared	0.90	0.86	0.87	0.94	0.92	0.89

551 achieve similar performance and is often more practical given its simplicity. On the other
552 hand, VAE seems to perform better when applied to OOD data, at least in the one test case
553 here. These observations warrant further research in the future using multiple test cases
554 and climate-relevant applications. We also note here that each method has multiple tuning
555 hyperparameters to optimize its uncertainty quantification. Consequently, the discrepancies
556 among the three methods could potentially be mitigated with proper hyperparameter tuning
557 (as discussed in Appendix B).

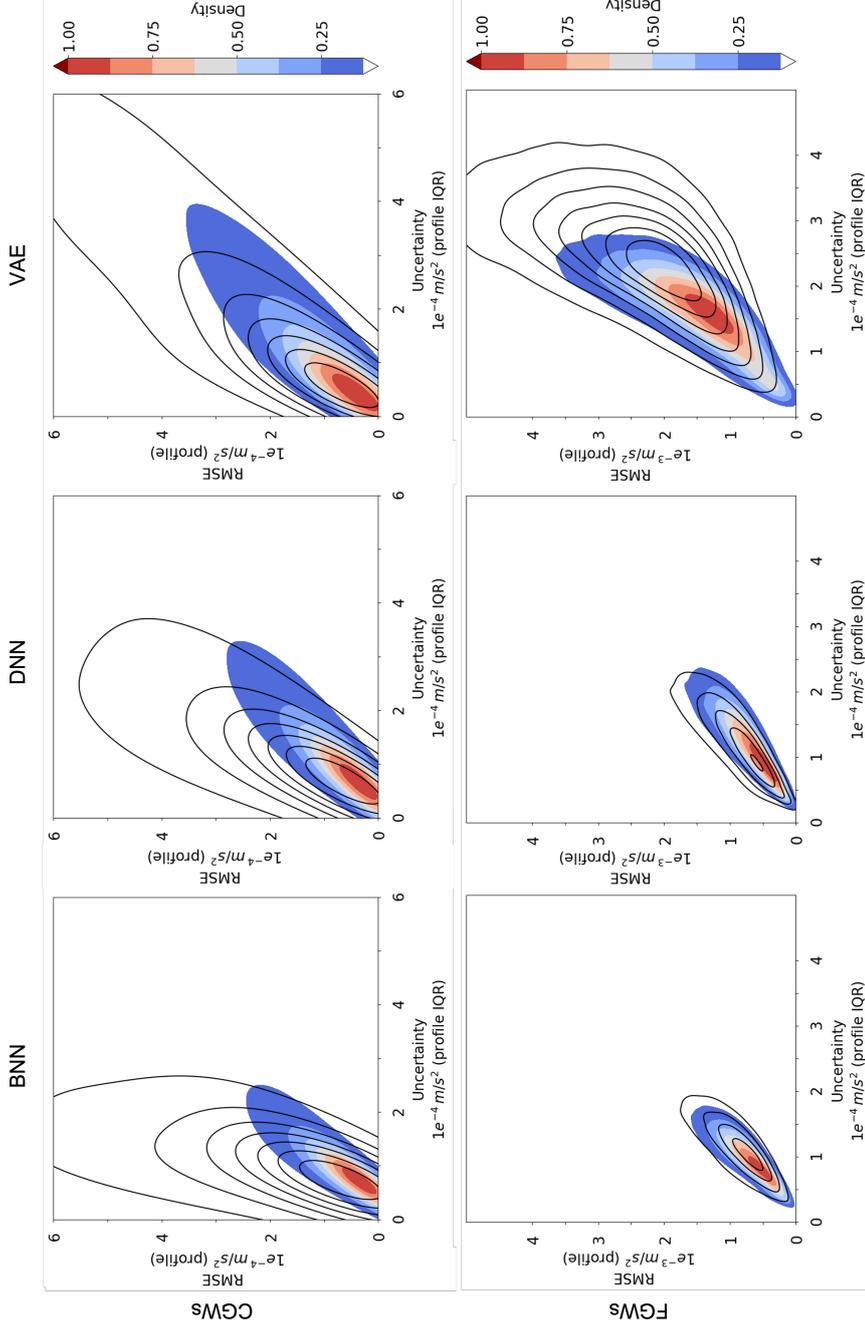


Figure 9. Gaussian kernel density for spread, defined as interquartile range (IQR) versus RMSE for (top) CGWs, and (bottom) FGWs, when validated against out-of-sample data from the current climate (color shading) and out-of-distribution data from the future, warmer climate (black lines). Results are shown for BNN (left), DNN (middle), and VAE (right). All NNs are exclusively trained on data from the current (control) climate.

3.3 Out-of-distribution (OOD) Generalization via Transfer Learning

As previously discussed, the GWP schemes in WACCM are coupled to their sources, which might change in a warmer climate. Specifically, under $4\times\text{CO}_2$ forcing, we expect changes in both the amplitude and the phase speed distribution of GWs, in particular for the CGWs, due to their built-in sensitivities to changes in the convection. Consequently, the physics scheme in WACCM produces slightly stronger GWD for CGWs, especially in the tail of the distribution. This intensified GWD results in a shorter quasi-biennial oscillation (QBO) period in WACCM. However, it is important to recognize that the response of the QBO to climate change differs across various general circulation models (Richter et al., 2022).

The intensification of the CGWs in future climate simulations presents an opportunity to study how NNs handle the OOD data. Our findings in the UQ section already suggest increased prediction errors when testing NNs with OOD data, which raises concerns about their applicability in climate change studies. To more thoroughly investigate this issue, we conduct additional evaluations on our ANNs, by applying them to data samples from future climate simulations, as illustrated in Figure 10. It is clear that the ANN for the CGWs does not generalize well, evidenced by a decrease in R^2 from 0.93 to 0.79. The ANN particularly struggles to capture the increase in GWD in the tail, with R^2 for the tails decreasing from 0.72 to 0.36. As a result, it seems unlikely that this emulator will accurately reproduce changes in the circulation under different climate conditions, such as the shorter QBO period resulting from future warming in WACCM.

In contrast to CGWs, the amplitude of FGWs shows a less marked increase in the future climate, and their PDF distribution closely resembles that of the control simulations. As a result, the ANN demonstrates better generalizability for FGWs when it is tested against future climate data, as seen in Figure 10d. There is only a slight decrease in the ANN's performance, with R^2 dropping from 0.97 to 0.95.

Two factors can contribute to the considerable OOD generalization errors in an NN when applied across two distinct systems. First, the input-output relationship might vary between the two systems. Second, the input variables in the new system could originate from a distribution different from that of the original system (regardless of whether the input-output relationship remains the same or changes). The former is hard to quantify in a high-dimensional dataset. The latter can be quantified using similarity distances. To help us better understand these differences between the OOD generalizability of CGWs and FGWs, we assess the similarity between their input and output data distributions from control and future climate simulations using the Mahalanobis distance (D). The Mahalanobis distance is a measure of the distance between a data point and a distribution (Ling & Templeton, 2015). Specifically, it is a multi-dimensional generalization of the idea of measuring how many standard deviations away a point is from the mean of the distribution. The application of Mahalanobis distance in understanding the source of OOD generalization errors in data-driven parameterization was previously demonstrated in Guan et al. (2022) for a simple turbulent system.

Table 2. Change of Mahalanobis distance based on the ratio of the average distance of the points that are more than 3 standard deviations away from the mean. The choice of the variables here is based on Appendix A, showing u, v, T , and source function contain most of the information needed for the NN.

Variables	u	v	T	Source (diabatic heating for CGWs, frontogenesis for FGWs)	Zonal drag	Meridional drag
Distance (Convection)	1.03	1.00	1.19	3.62	1.42	1.44
Distance (Front)	1.03	0.96	1.50	1.10	1.00	1.00

599 To use the Mahalanobis distance, we first calculate the mean and covariance matrix of
600 the training data from the control run. We then analyze the distribution of Mahalanobis
601 distances in this training data, setting a baseline value, referred to as D_{ctrl} . This baseline is
602 the average distance for data points that deviate by more than 3 standard deviations from
603 the mean. This choice aims to focus on outliers for which extrapolation is more challenging.
604 Following this, we apply the same process to the data points in the future climate dataset,
605 denoted as D_{warm} . Table 2 presents the ratio of D_{warm} for the warming scenario to D_{ctrl}
606 for the control scenario for selected variables. When this ratio is close to 1.0, it suggests
607 minimal changes in this variable’s distribution under a warming scenario. Note that the
608 NNs trained based only on these variables demonstrate performance comparable to NNs
609 trained on all variables (not shown), which is why we only focus on these few key variables.

610 The results reveal that among the various variables significantly contributing to the
611 emulation of CGWs, diabatic heating (source of CGWs) is the sole variable that exhibits
612 substantial changes from the control to the warming scenario. Conversely, changes in vari-
613 ables used to emulate FGWs are considerably smaller. This outcome suggests that the likely
614 reason for the better generalizability of FGWs is that the input distribution remains almost
615 unchanged (and the input-output relationship, which is the physics scheme, remains the
616 same too).

617 To improve the generalizability of the emulator for CGWs, we explore TL, a technique
618 introduced earlier and proven to be a powerful tool for improving the OOD generalizability
619 of data-driven parameterization in canonical turbulent flows (e.g., Guan et al., 2022; Subel
620 et al., 2023). Rather than re-training the entire NN for future climate scenarios, we only re-
621 train, following Subel et al. (2023), just a portion of the NN, thereby requiring only a small
622 fraction of the data. Figure 10e showcases the emulation results after only re-training the
623 first hidden layer of ANN using data from the first month of the WACCM simulation in the
624 $4\times\text{CO}_2$ scenario, which amounts to approximately 1% of the original training dataset. After
625 applying TL, the performance of the emulator in the warming scenario significantly improves,
626 with R^2 rising from 0.79 to 0.91, nearly matching its performance in the control simulations
627 ($R^2 = 0.93$). However, the improvement in the PDF tails is less pronounced, showing
628 only a modest increase in R^2 from 0.36 to 0.51. This is likely due to the limited number
629 of large-amplitude GW events within the one-month period. Instead of using more data

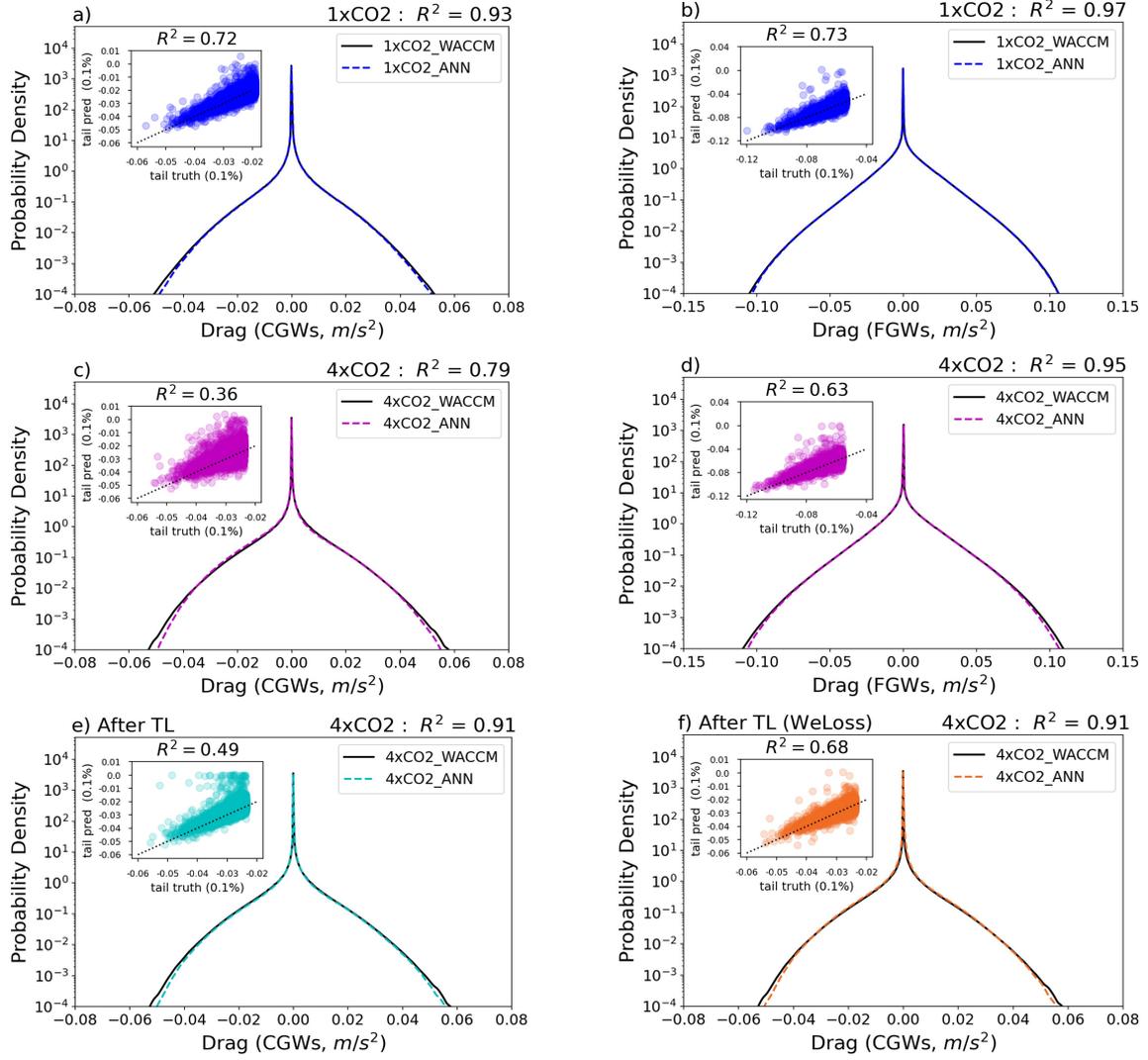


Figure 10. NN performance for pre-industrial and warming scenarios for different sources (a,c,e,f: CGWs ; b,d: FGWs). **a)** PDF of GWD due to CGWs in WACCM simulation and the predicted CGWs using NN emulator, scatter plot shows points for the tail part only. **b)** same as (a), but for FGWs. **c)** same as (a), but for the warming scenario, **d)** same as (b) but for the warming scenario. **e)** same as (c) but after applying transfer learning to the first hidden layer of the NN with 1-month WACCM simulation data under warming scenario ($\sim 1\%$ of the size of the training data) **f)** same as (e) but with the weighted loss function used when we conduct transfer learning (WeLoss).

630 from the future climate (which is challenging to obtain in a realistic situation), we leverage
 631 the WeLoss approach, described earlier, during re-training. This modification results in a
 632 significant improvement in the tail, with R^2 increasing from 0.51 to 0.68. Note that this
 633 improvement in the tail is critical, as inadequate learning of these rare but large-amplitude
 634 GWDs can result in significant errors and instabilities.

635 We would like to point out that during the TL experiments, we have examined the
 636 effects of re-training each individual hidden layer of the NN. Our findings indicate that
 637 re-training the first layer yields the best results, which aligns with the findings in Subel et
 638 al. (2023). Re-training the last layer only brings marginal improvements to the NN (not
 639 shown). Notably, our experiments involving re-training the first two layers did not result
 640 in further performance enhancements, suggesting that the number of neurons is not the
 641 primary factor contributing to the varied performance observed when re-training different
 642 layers.

643 Similar results regarding TL are also observed with other NNs used in this study. For
 644 instance, Figure 11 presents the same plot as Figure 10, but for the BNN. It is evident that
 645 BNN also struggles with generalization to OOD data, as could also be interpreted based
 646 on the results presented in section 3.2. It is also the case for DNN and VAE (not shown).
 647 Overall, when these NNs are tested against the $4\times\text{CO}_2$ future climate data, their accuracy
 648 is not better than the deterministic ANN. However, methods with UQ, especially the VAE
 649 (see Figure 9), could potentially indicate the increased uncertainty when testing with input
 650 data from the $4\times\text{CO}_2$ integration. These results underscore the necessity of re-training the
 651 NNs using TL.

652 **4 Emulation of Orographic GWs (OGWs)**

653 Similar to Chantry et al. (2021), our initial attempts to emulate OGWs did not succeed,
 654 primarily due to the presence of a pronounced data imbalance. Notably, the physics-based
 655 scheme responsible for OGW generation operates exclusively over terrestrial regions. How-
 656 ever, it is surprising that the issue of data imbalance continues to persist, even when we
 657 limit our NN training and testing exclusively to columns located over land (Figure 12a).
 658 Still, the emulated OGW drag often remains close to zero and completely fails to predict
 659 the rare events (Figure 12b), which poses a considerable hurdle for the emulator’s perfor-
 660 mance. Further investigations reveal that the key to this problem lies in the highly localized
 661 nature of orographic GWD, where significant drag is observed only at a handful of specific
 662 locations. Furthermore, even within these limited regions, GWD exhibits a significant in-
 663 termittent behavior. To help our understanding, we also conducted a K -means clustering
 664 analysis, categorizing GWD data for all land-based columns (Table 3). Among the 6 clus-
 665 ters, cluster 4 accounts for a staggering 97.51% of the dataset. Remarkably, all samples
 666 within this cluster exhibit exceptionally weak orographic GWD, as evidenced by the cluster
 667 center’s maximum GWD amplitude, which is two orders of magnitude smaller than that of
 668 other clusters.

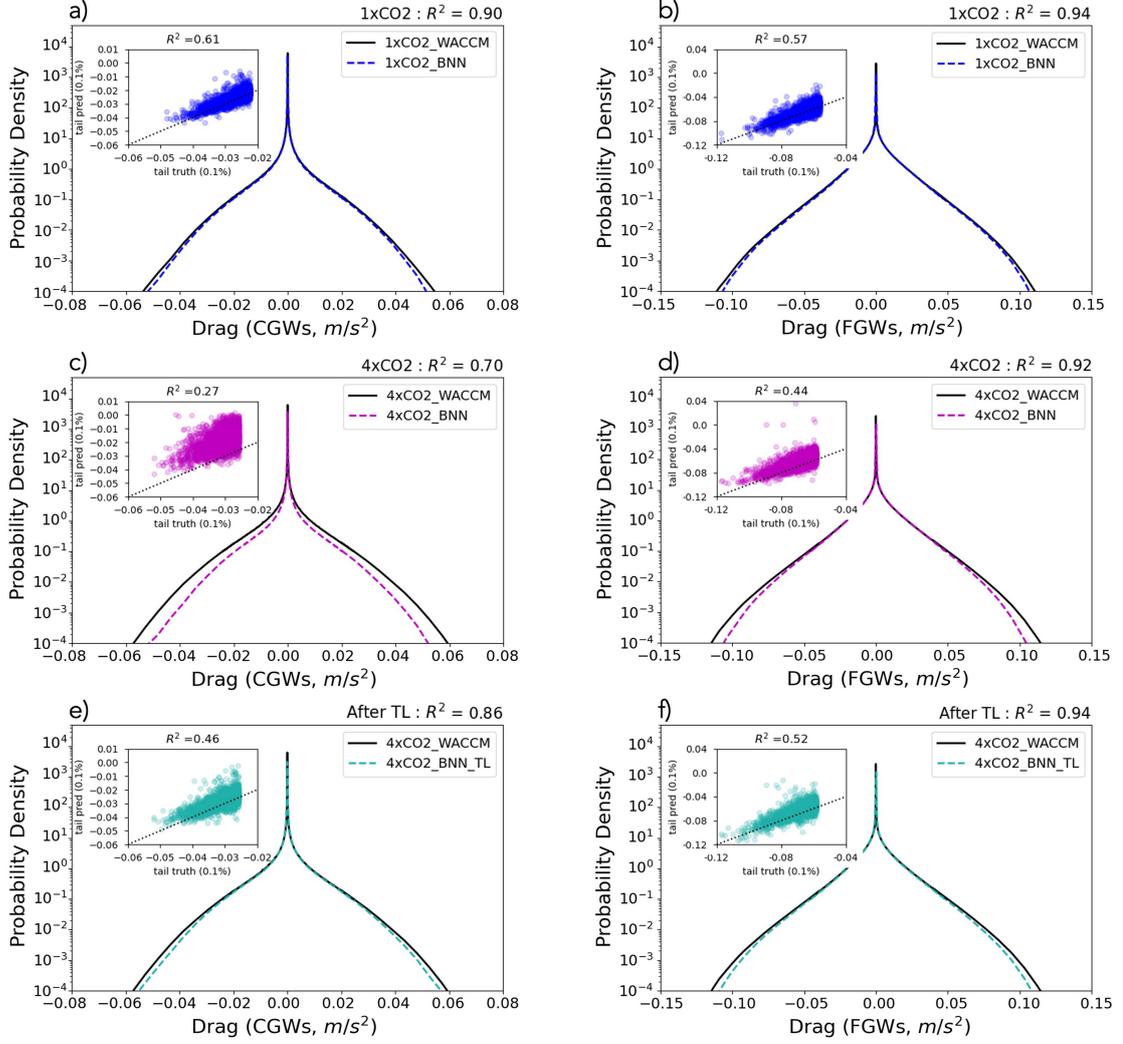


Figure 11. Panels (a) to (e) are the same as those in Figure 10 but for BNN. Panel (f) shows emulation for FGWs under warming scenario after applying transfer learning to the first hidden layer of the NN.

Table 3. Clustering analysis for OGWs. Analysis is done for all columns over land in the training data.

Cluster	Frequency (%) in the training data	Maximum GWD amplitude of cluster center
c1	0.18	8.7 e-3 m/s ²
c2	0.13	4.4 e-3 m/s ²
c3	0.93	3.6 e-3 m/s ²
c4	97.51	2.8 e-5 m/s ²
c5	0.15	2.1 e-3 m/s ²
c6	1.10	4.3 e-3 m/s ²

669 To overcome this persistent data imbalance in the OGWs, we first separate all columns
670 over land into large-drag columns (with column maximum greater than one STD of all
671 GWD from OGWs) and small-drag columns. We then perform subsampling on the latter
672 group only to create a more balanced dataset. To improve NN training, we also include all
673 columns from the 6-year simulation to augment the sample size of the large-drag columns.
674 Figures 12c and 12d illustrate the performance after re-balancing the dataset. Notably, the
675 result represents a substantial improvement, evidenced by an R^2 increase from 0.29 to 0.80,
676 and also a significant improvement in the accuracy for rare events. While we acknowledge
677 that this skill remains lower than what is achieved for CGWs and FGWs, it already signifies a
678 reasonable NN. Furthermore, we posit that by incorporating additional training data (either
679 by extending the WACCM model integration or simply augmenting the data with OGWs
680 scheme only), we can further improve our emulation results. The possibility of achieving
681 superior emulation outcomes through the adoption of an alternative NN architecture is also
682 possible, although such exploration is beyond the scope of this paper.

683 5 Summary and Discussion

684 Through the emulation of complex GWPs in a state-of-the-art atmospheric model
685 (WACCM), we have elucidated and explored solutions for three critical challenges in the
686 development of ML-based data-driven SGS schemes for climate applications: data imbalance,
687 UQ, and OOD generalizability under different climates. A brief summary is provided
688 below:

- 689 1. In the presence of non-stationary, and highly imbalanced datasets, such as those en-
690 countered in WACCM, specialized approaches (e.g., resampling and weighted loss
691 function) are essential to enhance the performance of data-driven models. Through
692 resampling, we have successfully trained a robust NN emulator for OGWs, a challeng-
693 ing task as demonstrated in Chantry et al. (2021). The effectiveness of the trained
694 emulator is also significantly influenced by the choice of the loss function used dur-
695 ing training. In our case, while a weighted loss function (WeLoss) does not improve
696 the overall R^2 score, it yields significant improvements in the emulation results for

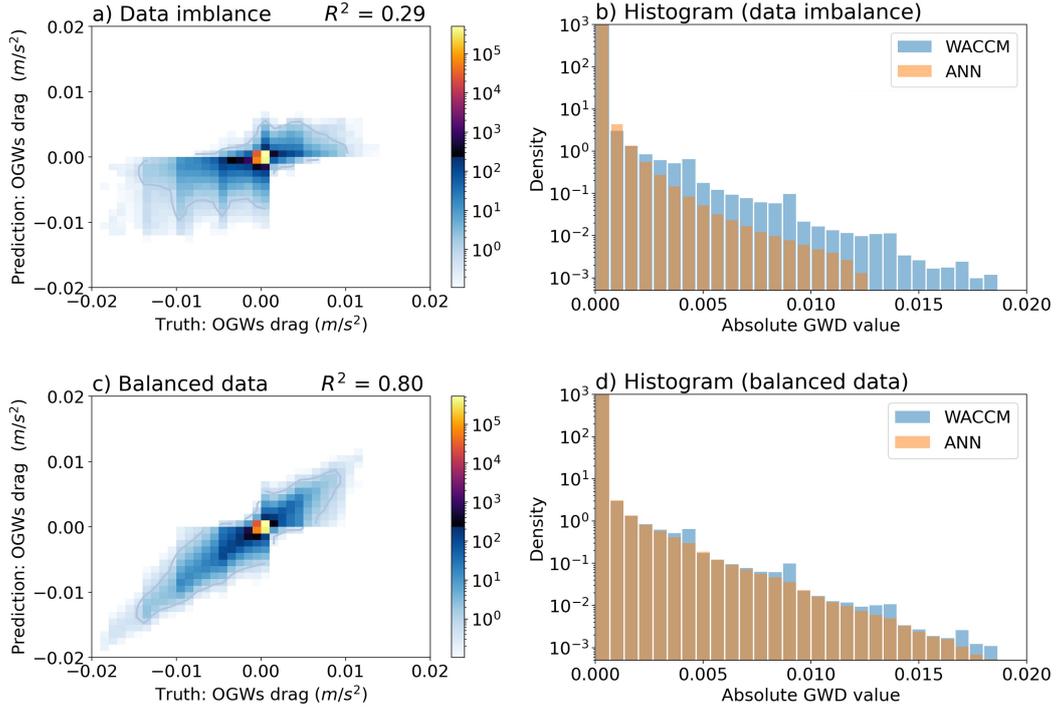


Figure 12. Performance of the emulator for OGWs when trained with all columns over land (panel a) & panel b)) and balanced training data with a balanced number of large-drag columns (column maximum > 1 STD of all GWD from OGWs) and small-drag columns (panel c) & panel d).

697 the PDF tails of the GWD. This finding aligns with those in Lopez-Gomez et al.
 698 (2022), where their custom loss function, tailored to emphasize extreme events, led
 699 to substantial improvements in predicting heatwaves.

700 2. All three UQ methods employed in this study provide reasonable uncertainty esti-
 701 mates for GWD prediction for the current climate. The spread-skill plots (refer to
 702 Figures 8 and 9) indicate that greater uncertainty corresponds to a larger prediction
 703 error. Yet, the reliability of UQ methods diminishes when they are challenged with
 704 OOD data. Both BNN and DNN used in this study tend to be overconfident in esti-
 705 mating CGWs in a warmer climate, thereby struggling to identify OOD samples.
 706 The VAE, on the other hand, yields rather promising results in providing useful UQ
 707 for OOD data. Given the variations in different methods, the metrics selected to
 708 assess the SGS model will play a significant role in determining the choice for the UQ
 709 methods. We also note that further optimization of tunable parameters within each
 710 UQ method could affect their performance (refer to Appendix C).

711 3. Our findings illustrate the challenges SGS schemes face in generalizing to OOD data
 712 and extrapolating to new climates. Nonetheless, the TL approach has proven highly
 713 effective in aiding an NN to extrapolate to different climates. For CGWs in WACCM,
 714 the physics-based scheme exhibits larger GWD under $4\times\text{CO}_2$ forcing, primarily due
 715 to an increase in diabatic heating from convection. With only one month of sim-

716 ulation data from this future warming scenario (representing approximately 1% of
 717 the original training dataset), we successfully reduce its OOD generalization error
 718 through re-training the first layer of the NN, following the findings of Subel et al.
 719 (2023). Additionally, we have illustrated the value of metrics like the Mahalanobis
 720 distance in assessing the potential OOD generalizability of NNs.

721 We would like to emphasize that these challenges are often intertwined. For instance,
 722 addressing data imbalance in CGWs is a prerequisite for obtaining an accurate NN model,
 723 which, in turn, impacts UQ and OOD generalizability assessments. Moreover, there exists
 724 a strong link between UQ and OOD generalizability evaluations: if the NN struggles with
 725 OOD generalization, performing poorly with OOD data, the reliability of UQ for such data
 726 (e.g., data from a warmer climate) also becomes questionable. This presents a substantial
 727 challenge for UQ methods, especially for climate change research where reliable UQ methods
 728 are crucial.

729 This study has primarily focused on offline skill assessment. We acknowledge that good
 730 offline performance (at least in terms of common metrics such as R^2) is not necessarily an
 731 indicator of stable and accurate online (coupled to climate model) performance (Ross et
 732 al., 2022; Guan et al., 2022), though more strict metrics such as R^2 of the PDF tails might
 733 better connect the offline and online performance (Pahlavan et al., 2023). However, for
 734 the purpose of this study, which is to provide a testbed to test ideas for data imbalance,
 735 UQ, and OOD generalization with transfer learning, the offline tests, particularly using
 736 the several metrics we have used, suffice. That said, the main reason that we have not
 737 provided online results is that coupling various complex NNs, with the same framework, to
 738 complex climate models (e.g., WACCM) without slowing down the model is a challenging
 739 and time-consuming task (Espinosa et al., 2022), and this is work in progress.

740 Emulating complex GWPs within the WACCM provided a unique opportunity to ad-
 741 dress three critical challenges in developing ML-based, data-driven SGS schemes for climate
 742 science applications. However, it is crucial to acknowledge that such emulated schemes
 743 essentially adopt the limitations inherent in the physics-based schemes. Addressing these
 744 limitations, the next step is to harness high-resolution data from GW-resolving simula-
 745 tions, which are carefully validated against observational data. A library of such high-
 746 resolution simulations, notably of convectively generated GWs using the Weather Research
 747 and Forecasting (WRF) model, is now established (Sun et al., 2023), alongside additional
 748 global high-resolution simulations (Wedi et al., 2020; Polichtchouk et al., 2023; Köhler et al.,
 749 2023). The next phase involves integrating the approaches outlined in this study with the
 750 data from these GW-resolving simulations to develop a stable, trustworthy, and generaliz-
 751 able data-driven GWP scheme. This scheme is then expected to overcome the limitations of
 752 physics-based GWPs and potentially incorporate features like the transient effect (Bölöni
 753 et al., 2021; Kim et al., 2021) and lateral propagation of GWs (e.g., Sato et al., 2009)—marking
 754 a significant advancement towards next-generation GWP schemes.

Appendix A Input/output variables for the physics-based GWP schemes and their emulators

We use the exact same inputs as those of each GWP scheme in the WACCM for the training of the NN-based emulator of that scheme. These inputs are listed in Table A1. As for the outputs, we only consider the zonal and meridional drag forcings. The GWPs in WACCM also estimate additional effects of the GWs that result in changes of temperature profile and vertical diffusion. These outputs are not considered in our emulations.

Table A1. List of the input and output variables for the NNs trained as emulators of the GWP schemes in WACCM. The numbers in parentheses in front of each variable are the number of vertical levels for that variable. Note that each input and output is a 1D column at a given latitude/longitude grid point. Diabatic heating in WACCM is provided by the cumulus scheme. The topography variables listed in the table are *mxdis* (height estimates for ridges), *hwdth* (width of ridges), *clngt* (length of ridges), *angll* (orientation of ridges), and *anixy* (anisotropy of ridges).

GWP	Input			Output
	pressure levels	surface level	forcing	
CGWs	$u(70)$,	lat (1),	diabatic heating (70)	zonal drag GWD _x (70), meridional drag GWD _y (70),
FGWs	$v(70)$, $T(70)$,	lon (1), $P_{surface}$ (1),	frontogenesis function (70)	
OGWs	$z(70)$, $\rho(71)$, Brunt–Väisälä frequency N (70), dry static energy DSE (70)		<i>mxdis</i> (16), <i>hwdth</i> (16), <i>clngt</i> (16), <i>angll</i> (16), <i>anixy</i> (16),	

From Table A1, one can guess that some input variables are correlated with each other. Consequently, it is plausible that the trained NNs may have spurious connections. Preliminary tests further support this notion, indicating that employing only u, v, T , and the forcing function as inputs yields comparable offline skill (results not presented here).

Appendix B Tuning UQ-equipped NNs

In addition to the hyperparameters of the deterministic NNs, designing an architecture for UQ often demands additional hyperparameter optimization. For instance, for the DNN, decisions need to be made regarding the number of neurons to drop out (dropout rate). While less common, one can also choose whether to apply dropout to all hidden layers or only selected ones. Variations in the dropout rate and the layers to which dropout is applied can influence the final configuration and performance of the DNN. Figure B1 illustrates these effects. As we increase the number of dropped neurons (whether through a higher dropout rate or by subjecting more layers to dropout), the uncertainty in the DNN predictions tends to rise. Yet, there is a persistent pattern in the relationship between spread (IQR) and RMSE across the various plots in Figure B1. Specifically, as spread increases, RMSE concurrently grows, consistent with the insights highlighted in Figure 9.

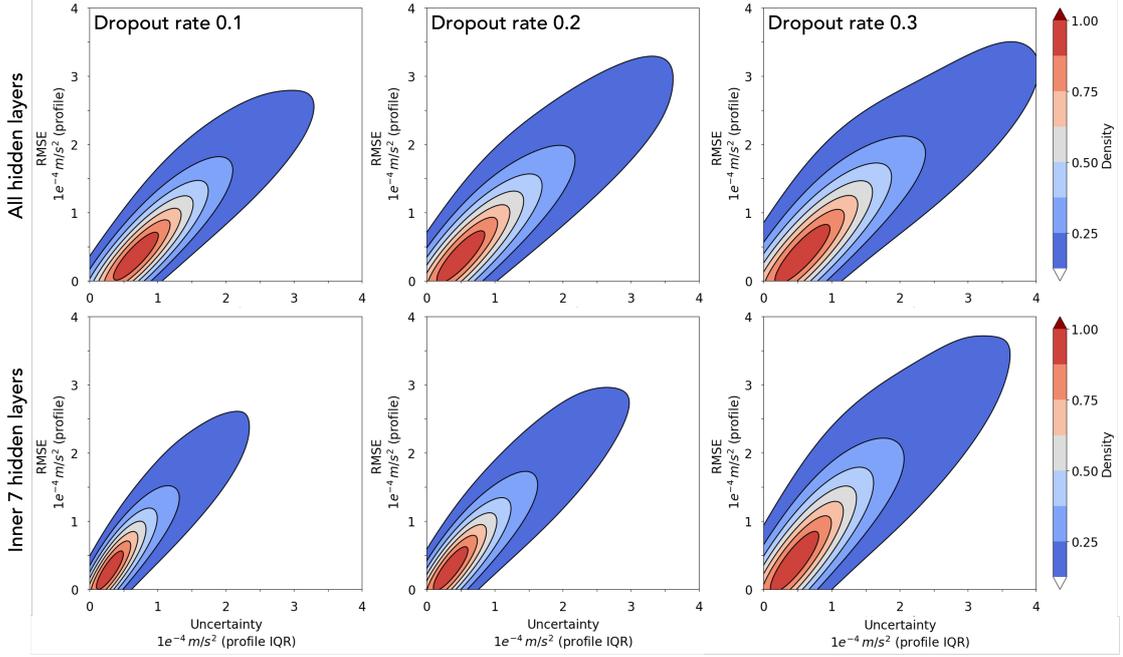


Figure B1. Similar to Figure 9 but for DNN only with different dropout rates, which are applied to different numbers of hidden layers.

778 In the case of BNN or VAE, even though there is no dropout rate, there are distinct
 779 tuning opportunities available. For instance, with the VAE, one might consider applying
 780 dropout to the NN emulator. Moreover, given that the loss function in VAE comprises three
 781 components, decisions can be made regarding which component to penalize more heavily,
 782 allowing for nuanced adjustments to its performance.

783 Appendix C The UQ metrics

784 Each point in the spread-skill plot corresponds to one specific bin of ensemble spread
 785 (\overline{SD}_k), which is defined as the average standard deviation of the ensemble members. We
 786 first separate the spread using a pre-selected number of bins (a subjective choice of 15 is
 787 used here). Then for the k^{th} bin:

$$\begin{cases} \text{RMSE}_k = \left[\frac{1}{N_k} \sum_{i=1}^{N_k} (\hat{y}_i - \bar{y}_i)^2 \right]^{\frac{1}{2}} \\ \overline{SD}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \left[\frac{1}{M-1} \sum_{j=1}^M (\bar{y}_i - y_{ij})^2 \right]^{\frac{1}{2}} \\ \bar{y}_i = \frac{1}{M} \sum_{j=1}^M y_{ij} \end{cases} \quad (\text{C1})$$

788 \hat{y}_i is the observed value for the i^{th} example, \bar{y}_i is the mean prediction for the i^{th} example, y_{ij}
 789 is the j^{th} prediction for the i^{th} example, N_k is the total number of examples in the k^{th} bin,
 790 and M is the ensemble size. Following Haynes et al. (2023), we summarize the quality of the

791 spread-skill plot by two measures: spread-skill reliability (SSREL) and overall spread-skill
 792 ratio (SSRAT). SSREL is the bin-weighted mean distance from the 1-to-1 line:

$$\text{SSREL} = \sum_{k=1}^K \frac{N_k}{N} |\text{RMSE}_k - \overline{\text{SD}}_k| \quad (\text{C2})$$

793 where N is the total number of examples, K is the total number of bins, and other variables
 794 are as in Equation C1. SSREL varies from $[0, \infty)$, and the ideal value is 0. On the other
 795 hand, SSRAT is averaged over the whole dataset:

$$\text{SSRAT} = \frac{\overline{\text{SD}}}{\text{RMSE}} \quad (\text{C3})$$

796 SSRAT also varies from $[0, \infty)$, and the ideal value is 1. $\text{SSRAT} > 1$ indicates the model is
 797 under-confident on average, while $\text{SSRAT} < 1$ indicates that the model is overconfident on
 798 average.

799 In Equation (C1), each level of a GWD profile is considered as an individual sample.
 800 As discussed earlier, while these samples help assess the model's overall performance, our
 801 main interest is often the uncertainty of individual GWD profiles. Such uncertainty informs
 802 the trustworthiness of the model's prediction for that specific profile. Accordingly, for each
 803 profile, we can compute:

$$\begin{cases} \text{RMSE}_{\text{profile}} = \left[\frac{1}{N_z} \sum_{z=1}^{N_z} (\hat{y}_z - \bar{y}_z)^2 \right]_{\text{profile}}^{\frac{1}{2}} \\ \text{IQR}_{\text{profile}} = \left[\frac{1}{N_z} \sum_{z=1}^{N_z} (y_{z,75th} - y_{z,25th})^2 \right]_{\text{profile}}^{\frac{1}{2}} \\ \bar{y}_z = \left[\frac{1}{M} \sum_{j=1}^M y_{zj} \right]_{\text{profile}} \end{cases} \quad (\text{C4})$$

804 where N_z is the number of vertical levels for each profile, and $\text{IQR}_{\text{profile}}$ is its interquartile
 805 range: $y_{z,25th}$ corresponds with the 25th percentile, and $y_{z,75th}$ corresponds with the 75th
 806 percentile.

807 Open Research

808 The data for all the analyses in the main text are available at [https://doi.org/10](https://doi.org/10.5281/zenodo.10019987)
 809 [.5281/zenodo.10019987](https://doi.org/10.5281/zenodo.10019987). The emulator code is available at [https://github.com/yqsun91/](https://github.com/yqsun91/WACCM-Emulation)
 810 [WACCM-Emulation](https://github.com/yqsun91/WACCM-Emulation). All the raw WACCM output data are available on request from authors.

811 Acknowledgments

812 We thank Andre Souza for insightful discussions. This work was supported by grants from
 813 the NSF OAC CSSI program (#2005123 , #2004512, #2004492, #2004572), and by the

814 generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program
 815 to PH, MJA, EG, and AS. PH is also supported by the Office of Naval Research (ONR)
 816 Young Investigator Award N00014-20-1-2722. SL is supported by the Office of Science, U.S.
 817 Department of Energy Biological and Environmental Research as part of the Regional and
 818 Global Climate Model Analysis program area. Computational resources were provided by
 819 NSF XSEDE (allocation ATM170020) and NCAR’s CISL (allocation URIC0009).

820 References

- 821 Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., ...
 822 Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques,
 823 applications and challenges. *Information Fusion*, *76*, 243-297. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S1566253521001081)
 824 www.sciencedirect.com/science/article/pii/S1566253521001081 doi: [https://doi](https://doi.org/10.1016/j.inffus.2021.05.008)
 825 [.org/10.1016/j.inffus.2021.05.008](https://doi.org/10.1016/j.inffus.2021.05.008)
- 826 Achatz, U. (2022). Gravity waves and their impact on the atmospheric flow. In *Atmo-*
 827 *spheric dynamics* (pp. 407–505). Berlin, Heidelberg: Springer Berlin Heidelberg. Re-
 828 trieved from https://doi.org/10.1007/978-3-662-63941-2_10 doi: 10.1007/978-3-
 829 -662-63941-2_10
- 830 Alexander, M. J., Geller, M., McLandress, C., Polavarapu, S., Preusse, P., Sassi, F., ...
 831 Watanabe, S. (2010). Recent developments in gravity-wave effects in climatemodels
 832 and the global distribution of gravity-wavemomentum flux from observations and models.
 833 *Quarterly Journal of the Royal Meteorological Society*, *136*. doi: 10.1002/qj.637
- 834 Amiramjadi, M., Plougonven, R., Mohebalhojeh, A. R., & Mirzaei, M. (2022). Using
 835 machine learning to estimate non-orographic gravity wave characteristics at source levels.
 836 *Journal of the Atmospheric Sciences*. doi: 10.1175/JAS-D-22-0021.1
- 837 Ando, S., & Huang, C. Y. (2017). Deep over-sampling framework for classifying imbalanced
 838 data. In *Machine learning and knowledge discovery in databases: European conference,*
 839 *ecml pkdd 2017, skopje, macedonia, september 18–22, 2017, proceedings, part i 10* (pp.
 840 770–785).
- 841 Bacmeister, J. T., Newman, P. A., Gary, B. L., & Chan, K. R. (1994). An algorithm for
 842 forecasting mountain wave-related turbulence in the stratosphere. *Weather Forecasting*,
 843 *9*. doi: 10.1175/1520-0434(1994)009<0241:AAFFMW>2.0.CO;2
- 844 Balaji, V. (2021). Climbing down charney’s ladder: machine learning and the post-dennard
 845 era of computational climate science. *Philosophical Transactions of the Royal Society A*,
 846 *379*(2194), 20200085.
- 847 Baldi, P., Sadowski, P., & Whiteson, D. (2014). Searching for exotic particles in high-energy
 848 physics with deep learning. *Nature communications*, *5*(1), 4308.
- 849 Ballnus, B., Hug, S., Hatz, K., Görlitz, L., Hasenauer, J., & Theis, F. J. (2017). Com-
 850 prehensive benchmarking of markov chain monte carlo methods for dynamical systems.
 851 *BMC Systems Biology*, *11*(1), 1–18.
- 852 Barnes, E. A., Barnes, R. J., & DeMaria, M. (2023). Sinh-arcsinh-normal distributions
 853 to add uncertainty to neural network regression tasks: Applications to tropical cyclone

- 854 intensity forecasts. *Environmental Data Science*, 2, e15.
- 855 Beck, A., Flad, D., & Munz, C.-D. (2019). Deep neural networks for data-driven LES
856 closure models. *Journal of Computational Physics*, 398, 108910.
- 857 Beljaars, A. C., Brown, A. R., & Wood, N. (2004). A new parametrization of turbulent
858 orographic form drag. *Quarterly Journal of the Royal Meteorological Society*, 130. doi:
859 10.1256/qj.03.73
- 860 Belochitski, A., & Krasnopolsky, V. (2021). Robustness of neural network emulations of ra-
861 diative transfer parameterizations in a state-of-the-art general circulation model. *Geosci-
862 entific Model Development*, 14(12), 7425–7437. Retrieved from [https://gmd.copernicus](https://gmd.copernicus.org/articles/14/7425/2021/)
863 [.org/articles/14/7425/2021/](https://gmd.copernicus.org/articles/14/7425/2021/) doi: 10.5194/gmd-14-7425-2021
- 864 Beucler, T., Ebert-Uphoff, I., Rasp, S., Pritchard, M., & Gentine, P. (2021, 05). *Machine
865 learning for clouds and climate (invited chapter for the agu geophysical monograph series
866 "clouds and climate")*. doi: 10.1002/essoar.10506925.1
- 867 Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015). *Weight uncertainty
868 in neural networks*.
- 869 Bolton, T., & Zanna, L. (2019). Applications of deep learning to ocean data inference
870 and subgrid parameterization. *Journal of Advances in Modeling Earth Systems*, 11(1),
871 376–399.
- 872 Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and
873 stabilizing machine-learning parametrizations of convection. *Journal of the Atmospheric
874 Sciences*, 77(12), 4357 - 4375. Retrieved from [https://journals.ametsoc.org/view/
875 journals/atms/77/12/jas-d-20-0082.1.xml](https://journals.ametsoc.org/view/journals/atms/77/12/jas-d-20-0082.1.xml) doi: [https://doi.org/10.1175/JAS-D-20-
876 -0082.1](https://doi.org/10.1175/JAS-D-20-0082.1)
- 877 Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network
878 parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Sys-
879 tems*, 11(8), 2728-2744. Retrieved from [https://agupubs.onlinelibrary.wiley.com/
880 doi/abs/10.1029/2019MS001711](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001711) doi: <https://doi.org/10.1029/2019MS001711>
- 881 Buda, M., Maki, A., & Mazurowski, M. A. (2018). A systematic study of the class imbalance
882 problem in convolutional neural networks. *Neural Networks*, 106, 249-259. Retrieved from
883 <https://www.sciencedirect.com/science/article/pii/S0893608018302107> doi:
884 <https://doi.org/10.1016/j.neunet.2018.07.011>
- 885 Bölöni, G., Kim, Y. H., Borchert, S., & Achatz, U. (2021). Toward transient subgrid-scale
886 gravity wave representation in atmospheric models. part i: Propagation model including
887 nondissipative wave mean-flow interactions. *Journal of the Atmospheric Sciences*, 78.
888 doi: 10.1175/JAS-D-20-0065.1
- 889 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Ma-
890 chine learning emulation of gravity wave drag in numerical weather forecasting. *Jour-
891 nal of Advances in Modeling Earth Systems*, 13(7), e2021MS002477. Retrieved
892 from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002477>
893 (e2021MS002477 2021MS002477) doi: <https://doi.org/10.1029/2021MS002477>
- 894 Chattopadhyay, A., Nabizadeh, E., & Hassanzadeh, P. (2020). Analog forecasting of
895 extreme-causing weather patterns using deep learning. *Journal of Advances in Model-*

- 896 *ing Earth Systems*, 12(2), e2019MS001958.
- 897 Chattopadhyay, A., Subel, A., & Hassanzadeh, P. (2020). Data-driven super-
898 parameterization using deep learning: Experimentation with multiscale Lorenz 96 sys-
899 tems and transfer learning. *Journal of Advances in Modeling Earth Systems*, 12(11),
900 e2020MS002084.
- 901 Chawla, N. V., Japkowicz, N., & Kotcz, A. (2004, jun). Editorial: Special issue on learning
902 from imbalanced data sets. *SIGKDD Explor. Newsl.*, 6(1), 1–6. Retrieved from <https://doi.org/10.1145/1007730.1007733> doi: 10.1145/1007730.1007733
- 903
- 904 Chen, N., & Majda, A. J. (2019). A new efficient parameter estimation algorithm for high-
905 dimensional complex nonlinear turbulent dynamical systems with partial observations.
906 *Journal of Computational Physics*, 397, 108836.
- 907 Clare, M. C., Sonnewald, M., Lguensat, R., Deshayes, J., & Balaji, V. (2022). Explainable
908 artificial intelligence for bayesian neural networks: Towards trustworthy predictions of
909 ocean dynamics. *arXiv preprint arXiv:2205.00202*.
- 910 Delle Monache, L., Eckel, F. A., Rife, D. L., Nagarajan, B., & Searight, K. (2013). Probabilistic
911 weather prediction with an analog ensemble. *Monthly Weather Review*, 141(10),
912 3498–3516.
- 913 Dong, W., Fritts, D. C., Liu, A. Z., Lund, T. S., Liu, H.-L., & Snively, J.
914 (2023). Accelerating atmospheric gravity wave simulations using machine learning:
915 Kelvin-helmholtz instability and mountain wave sources driving gravity wave break-
916 ing and secondary gravity wave generation. *Geophysical Research Letters*, 50(15),
917 e2023GL104668. Retrieved from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023GL104668> (e2023GL104668 2023GL104668) doi: <https://doi.org/10.1029/2023GL104668>
- 918
- 919
- 920 Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J. (2022).
921 Machine Learning Gravity Wave Parameterization Generalizes to Capture the QBO and
922 Response to Increased CO₂. , 49(8), e98174. doi: 10.1029/2022GL098174
- 923 Finkel, J., Gerber, E. P., Abbot, D. S., & Weare, J. (2023). Revealing the statis-
924 tics of extreme events hidden in short weather forecast data. *AGU Advances*, 4(2),
925 e2023AV000881. Retrieved from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023AV000881> (e2023AV000881 2023AV000881) doi: <https://doi.org/10.1029/2023AV000881>
- 926
- 927
- 928 Foster, D., Gagne, D. J., & Whitt, D. B. (2021, 12). Probabilistic machine learning esti-
929 mation of ocean mixed layer depth from dense satellite and sparse in situ observations.
930 *Journal of Advances in Modeling Earth Systems*, 13. doi: 10.1029/2021MS002474
- 931 Frezat, H., Sommer, J. L., Fablet, R., Balarac, G., & Lguensat, R. (2022). A posteriori learn-
932 ing for quasi-geostrophic turbulence parametrization. *arXiv preprint arXiv:2204.03911*.
933 doi: 10.1029/2022MS003124
- 934 Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine
935 learning for stochastic parameterization: Generative adversarial networks in the lorenz'96
936 model. *Journal of Advances in Modeling Earth Systems*, 12(3), e2019MS001896.
- 937 Gal, Y., & Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing

- 938 model uncertainty in deep learning. In *international conference on machine learning* (pp.
939 1050–1059).
- 940 Garcia, R. R., Smith, A. K., Kinnison, D. E., Álvaro de la Cámara, & Murphy, D. J. (2017).
941 Modification of the gravity wave parameterization in the whole atmosphere community
942 climate model: Motivation and results. *Journal of the Atmospheric Sciences*, *74*(1),
943 275 - 291. Retrieved from [https://journals.ametsoc.org/view/journals/atasc/74/
944 1/jas-d-16-0104.1.xml](https://journals.ametsoc.org/view/journals/atasc/74/1/jas-d-16-0104.1.xml) doi: <https://doi.org/10.1175/JAS-D-16-0104.1>
- 945 Geller, M. A., Alexander, M. J., Love, P. T., Bacmeister, J., Ern, M., Hertzog, A., ...
946 others (2013). A comparison between gravity wave momentum fluxes in observations and
947 climate models. *Journal of Climate*, *26*(17), 6383–6405.
- 948 Gettelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J., Morrison, H.,
949 & Gantos, G. (2021). Machine learning the warm rain process. *Journal of Advances in
950 Modeling Earth Systems*, *13*. doi: 10.1029/2020MS002268
- 951 Gettelman, A., Mills, M. J., Kinnison, D. E., Garcia, R. R., Smith, A. K., Marsh, D. R.,
952 ... Randel, W. J. (2019, 12). The whole atmosphere community climate model version
953 6 (waccm6). *Journal of Geophysical Research: Atmospheres*, *124*, 12380-12403. doi:
954 10.1029/2019JD030943
- 955 Gordon, E. M., & Barnes, E. A. (2022, 8). Incorporating uncertainty into a regression
956 neural network enables identification of decadal state-dependent predictability in cesm2.
957 *Geophysical Research Letters*, *49*. doi: 10.1029/2022GL098635
- 958 Guan, Y., Chattopadhyay, A., Subel, A., & Hassanzadeh, P. (2022). Stable a posteriori
959 LES of 2D turbulence using convolutional neural networks: Backscattering analysis and
960 generalization to higher Re via transfer learning. *Journal of Computational Physics*, *458*,
961 111090.
- 962 Guan, Y., Subel, A., Chattopadhyay, A., & Hassanzadeh, P. (2023). Learning physics-
963 constrained subgrid-scale closures in the small-data regime for stable and accurate les.
964 *Physica D: Nonlinear Phenomena*, *443*, 133568. doi: [https://doi.org/10.1016/j.physd
965 .2022.133568](https://doi.org/10.1016/j.physd.2022.133568)
- 966 Guillaumin, A. P., & Zanna, L. (2021). Stochastic-deep learning parameterization of
967 ocean momentum forcing. *Journal of Advances in Modeling Earth Systems*, *13*(9),
968 e2021MS002534.
- 969 Hardiman, S. C., Scaife, A. A., Niekerk, A. v., Prudden, R., Owen, A., Adams, S. V., ...
970 Madge, S. (2023). Machine learning for non-orographic gravity waves in a climate model.
971 *Artificial Intelligence for the Earth Systems*.
- 972 Haynes, K., Lagerquist, R., McGraw, M., Musgrave, K., & Ebert-Uphoff, I. (2023). Creat-
973 ing and evaluating uncertainty estimates with neural networks for environmental-science
974 applications. *Artificial Intelligence for the Earth Systems*, 1–58.
- 975 Hertzog, A., Alexander, M. J., & Plougonven, R. (2012). On the intermittency of gravity
976 wave momentum flux in the stratosphere. *Journal of the Atmospheric Sciences*, *69*(11),
977 3433–3448.
- 978 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ... others
979 (2017). The art and science of climate model tuning. *Bulletin of the American Meteorolo-*

- 980 *logical Society*, 98(3), 589–602.
- 981 Huang, C., Li, Y., Loy, C. C., & Tang, X. (2016). Learning deep representation for imbal-
 982 anced classification. In *Proceedings of the IEEE conference on computer vision and pattern*
 983 *recognition* (pp. 5375–5384).
- 984 Iglesias-Suarez, F., Gentine, P., Solino-Fernandez, B., Beucler, T., Pritchard, M., Runge,
 985 J., & Eyring, V. (2023). *Causally-informed deep learning to improve climate models and*
 986 *projections*.
- 987 Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study.
 988 *Intelligent Data Analysis*, 6, 429–449. Retrieved from [https://doi.org/10.3233/IDA-](https://doi.org/10.3233/IDA-2002-6504)
 989 [-2002-6504](https://doi.org/10.3233/IDA-2002-6504) (5) doi: 10.3233/IDA-2002-6504
- 990 Johnson, J. M., & Khoshgoftaar, T. M. (2019, Mar 19). Survey on deep learning with class
 991 imbalance. *Journal of Big Data*, 6(1), 27. Retrieved from [https://doi.org/10.1186/s40537-](https://doi.org/10.1186/s40537-019-0192-5)
 992 [s40537-019-0192-5](https://doi.org/10.1186/s40537-019-0192-5) doi: 10.1186/s40537-019-0192-5
- 993 Kim, Y., Bölöni, G., Borchert, S., Chun, H. Y., & Achatz, U. (2021). Toward transient
 994 subgrid-scale gravity wave representation in atmospheric models. part ii: Wave intermit-
 995 tency simulated with convective sources. *Journal of the Atmospheric Sciences*, 78. doi:
 996 10.1175/JAS-D-20-0066.1
- 997 Kim, Y., Eckermann, S. D., & Chun, H. (2003). An overview of the past, present and
 998 future of gravity-wave drag parametrization for numerical climate and weather prediction
 999 models. *Atmosphere-Ocean*, 41, 65–98. Retrieved from [https://doi.org/10.3137/ao-](https://doi.org/10.3137/ao.410105)
 1000 [.410105](https://doi.org/10.3137/ao.410105) doi: 10.3137/ao.410105
- 1001 Kingma, D. P., & Welling, M. (2014). Auto-encoding variational bayes. *2nd International*
 1002 *Conference on Learning Representations, ICLR 2014 - Conference Track Proceedings*.
- 1003 Köhler, L., Green, B., & Stephan, C. C. (2023). Comparing loon superpressure balloon
 1004 observations of gravity waves in the tropics with global storm-resolving models. *Journal*
 1005 *of Geophysical Research: Atmospheres*, 128(15), e2023JD038549.
- 1006 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). New approach to
 1007 calculation of atmospheric model physics: Accurate and fast neural network emulation
 1008 of longwave radiation in a climate model. *Monthly Weather Review*, 133(5), 1370 -
 1009 1383. Retrieved from [https://journals.ametsoc.org/view/journals/mwre/133/5/](https://journals.ametsoc.org/view/journals/mwre/133/5/mwr2923.1.xml)
 1010 [mwr2923.1.xml](https://journals.ametsoc.org/view/journals/mwre/133/5/mwr2923.1.xml) doi: <https://doi.org/10.1175/MWR2923.1>
- 1011 Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., ... Courville,
 1012 A. (2021). Out-of-distribution generalization via risk extrapolation (rex). In *International*
 1013 *conference on machine learning* (pp. 5815–5826).
- 1014 Kruse, C. G., Alexander, M. J., Hoffmann, L., Niekerk, A. V., Polichtchouk, I., Bacmeister,
 1015 J. T., ... Stein, O. (2022). Observed and modeled mountain waves from the surface to
 1016 the mesosphere near the drake passage. *Journal of the Atmospheric Sciences*, 79. doi:
 1017 10.1175/JAS-D-21-0252.1
- 1018 Li, X., Dai, Y., Ge, Y., Liu, J., Shan, Y., & Duan, L.-Y. (2022). Uncertainty modeling for
 1019 out-of-distribution generalization. *arXiv preprint arXiv:2202.03958*.
- 1020 Ling, J., & Templeton, J. (2015, 08). Evaluation of machine learning algorithms for
 1021 prediction of regions of high Reynolds averaged Navier Stokes uncertainty. *Physics of*

- 1022 *Fluids*, 27(8). Retrieved from <https://doi.org/10.1063/1.4927765> (085103) doi:
1023 10.1063/1.4927765
- 1024 Liu, Y., Racah, E., Prabhat, M., Correa, J., Khosrowshahi, A., Lavers, D., ... Collins,
1025 W. (2016, 05). Application of deep convolutional neural networks for detecting extreme
1026 weather in climate datasets.
- 1027 Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. (2022). Global extreme heat
1028 forecasting using neural weather models. *Artificial Intelligence for the Earth Systems*, 1
1029 - 41. Retrieved from [https://journals.ametsoc.org/view/journals/aies/aop/AIES-](https://journals.ametsoc.org/view/journals/aies/aop/AIES-D-22-0035.1/AIES-D-22-0035.1.xml)
1030 [D-22-0035.1/AIES-D-22-0035.1.xml](https://journals.ametsoc.org/view/journals/aies/aop/AIES-D-22-0035.1/AIES-D-22-0035.1.xml) doi: 10.1175/AIES-D-22-0035.1
- 1031 Lu, L., Jin, P., Pang, G., Zhang, Z., & Karniadakis, G. E. (2021). Learning nonlinear
1032 operators via deepnet based on the universal approximation theorem of operators. *Nature*
1033 *machine intelligence*, 3(3), 218–229.
- 1034 Maalouf, M., & Siddiqi, M. (2014). Weighted logistic regression for large-scale imbalanced
1035 and rare events data. *Knowledge-Based Systems*, 59, 142–148.
- 1036 Maalouf, M., & Trafalis, T. B. (2011). Robust weighted kernel logistic regression in imbal-
1037 anced and rare events data. *Computational Statistics & Data Analysis*, 55(1), 168–183.
- 1038 Maddox, W. J., Izmailov, P., Garipov, T., Vetrov, D. P., & Wilson, A. G. (2019). A
1039 simple baseline for bayesian uncertainty in deep learning. *Advances in Neural Information*
1040 *Processing Systems*, 32.
- 1041 Mamalakis, A., Barnes, E. A., & Ebert-Uphoff, I. (2022). Investigating the fidelity of
1042 explainable artificial intelligence methods for applications of convolutional neural networks
1043 in geoscience. *Artificial Intelligence for the Earth Systems*, 1(4), e220012.
- 1044 Matsuoka, D., Watanabe, S., Sato, K., Kawazoe, S., Yu, W., & Easterbrook, S. (2020).
1045 Application of Deep Learning to Estimate Atmospheric Gravity Wave Parameters in Re-
1046 analysis Data Sets. , 47(19), e89436. doi: 10.1029/2020GL089436
- 1047 Maulik, R., San, O., Rasheed, A., & Vedula, P. (2019). Subgrid modelling for two-
1048 dimensional turbulence using neural networks. *Journal of Fluid Mechanics*, 858, 122–
1049 144.
- 1050 McGovern, A., Lagerquist, R., Gagne, D. J., Jergensen, G. E., Elmore, K. L., Homeyer,
1051 C. R., & Smith, T. (2019). Making the black box more transparent: Understanding
1052 the physical implications of machine learning. *Bulletin of the American Meteorological*
1053 *Society*, 100. doi: 10.1175/BAMS-D-18-0195.1
- 1054 Miller, J. P., Taori, R., Raghunathan, A., Sagawa, S., Koh, P. W., Shankar, V., ... Schmidt,
1055 L. (2021). Accuracy on the line: on the strong correlation between out-of-distribution
1056 and in-distribution generalization. In *International conference on machine learning* (pp.
1057 7721–7735).
- 1058 Miloshevich, G., Cozian, B., Abry, P., Borgnat, P., & Bouchet, F. (2023, Apr). Probabilistic
1059 forecasts of extreme heatwaves using convolutional neural networks in a regime of lack
1060 of data. *Phys. Rev. Fluids*, 8, 040501. Retrieved from [https://link.aps.org/doi/](https://link.aps.org/doi/10.1103/PhysRevFluids.8.040501)
1061 [10.1103/PhysRevFluids.8.040501](https://link.aps.org/doi/10.1103/PhysRevFluids.8.040501) doi: 10.1103/PhysRevFluids.8.040501
- 1062 Nagarajan, V., Andreassen, A., & Neyshabur, B. (2020). Understanding the failure modes
1063 of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*.

- 1064 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameterize moist
 1065 convection: Potential for modeling of climate, climate change, and extreme events. *Journal*
 1066 *of Advances in Modeling Earth Systems*, *10*(10), 2548-2563. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351)
 1067 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351 doi: [https://](https://doi.org/10.1029/2018MS001351)
 1068 doi.org/10.1029/2018MS001351
- 1069 Oh, S. M., Rehg, J. M., Balch, T., & Dellaert, F. (2005). Data-driven mcmc for learning and
 1070 inference in switching linear dynamic systems. In *Proceedings of the national conference*
 1071 *on artificial intelligence* (Vol. 20, p. 944).
- 1072 Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., ... Snoek, J. (2019).
 1073 Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset
 1074 shift. *Advances in neural information processing systems*, *32*.
- 1075 Pahlavan, H. A., Hassanzadeh, P., & Alexander, M. J. (2023). Explainable offline-online
 1076 training of neural networks for parameterizations: A 1d gravity wave-qbo testbed in the
 1077 small-data regime. *arXiv preprint arXiv:2309.09024*.
- 1078 Palmer, T. (2019). Stochastic weather and climate models. *Nature Reviews Physics*, *1*(7),
 1079 463–471.
- 1080 Polichtchouk, I., Van Niekerk, A., & Wedi, N. (2023). Resolved gravity waves in the
 1081 extratropical stratosphere: Effect of horizontal resolution increase from o (10) to o (1)
 1082 km. *Journal of the Atmospheric Sciences*, *80*(2), 473–486.
- 1083 Prein, A. F., Langhans, W., Fossier, G., Ferrone, A., Ban, N., Goergen, K., ... others
 1084 (2015). A review on regional convection-permitting climate modeling: Demonstrations,
 1085 prospects, and challenges. *Reviews of geophysics*, *53*(2), 323–361.
- 1086 Psaros, A. F., Meng, X., Zou, Z., Guo, L., & Karniadakis, G. E. (2023). Uncertainty
 1087 quantification in scientific machine learning: Methods, metrics, and comparisons. *Journal*
 1088 *of Computational Physics*, *477*, 111902.
- 1089 Qi, D., & Majda, A. J. (2020). Using machine learning to predict extreme events in complex
 1090 systems. *Proceedings of the National Academy of Sciences*, *117*(1), 52-59. Retrieved
 1091 from <https://www.pnas.org/doi/abs/10.1073/pnas.1917285117> doi: 10.1073/pnas
 1092 .1917285117
- 1093 Ramachandran, P., Zoph, B., & Le, Q. V. (2017). Searching for activation functions. *arXiv*
 1094 *preprint arXiv:1710.05941*.
- 1095 Rasp, S. (2020). Coupled online learning as a way to tackle instabilities and biases in neural
 1096 network parameterizations: general algorithms and lorenz 96 case study (v1.0). *Geosci-*
 1097 *entific Model Development*, *13*(5), 2185–2196. Retrieved from [https://gmd.copernicus](https://gmd.copernicus.org/articles/13/2185/2020/)
 1098 [.org/articles/13/2185/2020/](https://gmd.copernicus.org/articles/13/2185/2020/) doi: 10.5194/gmd-13-2185-2020
- 1099 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid
 1100 processes in climate models. *Proceedings of the National Academy of Sciences*, *115*(39),
 1101 9684–9689.
- 1102 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., &
 1103 Prabhat. (2019, Feb 01). Deep learning and process understanding for data-driven earth
 1104 system science. *Nature*, *566*(7743), 195-204. Retrieved from [https://doi.org/10.1038/](https://doi.org/10.1038/s41586-019-0912-1)
 1105 [s41586-019-0912-1](https://doi.org/10.1038/s41586-019-0912-1) doi: 10.1038/s41586-019-0912-1

- 1106 Richter, J. H., Butchart, N., Kawatani, Y., Bushell, A. C., Holt, L., Serva, F., . . . Yukimoto,
 1107 S. (2022, 4). Response of the quasi-biennial oscillation to a warming climate in global
 1108 climate models. *Quarterly Journal of the Royal Meteorological Society*, *148*, 1490-1518.
 1109 doi: 10.1002/qj.3749
- 1110 Richter, J. H., Sassi, F., & Garcia, R. R. (2010). Toward a physically based gravity
 1111 wave source parameterization in a general circulation model. *Journal of the Atmospheric
 1112 Sciences*, *67*. doi: 10.1175/2009JAS3112.1
- 1113 Ross, A. S., Li, Z., Perezhugin, P., Fernandez-Granda, C., & Zanna, L. (2022). Benchmark-
 1114 ing of machine learning ocean subgrid parameterizations in an idealized model.
- 1115 Sato, K., Watanabe, S., Kawatani, Y., Tomikawa, Y., Miyazaki, K., & Takahashi, M. (2009).
 1116 On the origins of mesospheric gravity waves. *Geophysical Research Letters*, *36*. doi:
 1117 10.1029/2009GL039908
- 1118 Schneider, T., Jeevanjee, N., & Socolow, R. (2021). Accelerating progress in climate science.
 1119 *Physics Today*, *74*(6), 44–51.
- 1120 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling 2.0: A
 1121 blueprint for models that learn from observations and targeted high-resolution simula-
 1122 tions. *Geophysical Research Letters*, *44*. doi: 10.1002/2017GL076101
- 1123 Scinocca, J. F., & McFarlane, N. A. (2000). The parametrization of drag induced by
 1124 stratified flow over anisotropic orography. *Quarterly Journal of the Royal Meteorological
 1125 Society*, *126*. doi: 10.1002/qj.49712656802
- 1126 Seifert, A., & Rasp, S. (2020). Potential and limitations of machine learning for modeling
 1127 warm-rain cloud microphysical processes. *Journal of Advances in Modeling Earth Sys-
 1128 tems*, *12*(12), e2020MS002301. Retrieved from [https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020MS002301)
 1129 [.com/doi/abs/10.1029/2020MS002301](https://doi.org/10.1029/2020MS002301) (e2020MS002301 10.1029/2020MS002301) doi:
 1130 <https://doi.org/10.1029/2020MS002301>
- 1131 Shamekh, S., Lamb, K. D., Huang, Y., & Gentine, P. (2023). Implicit learning of convective
 1132 organization explains precipitation stochasticity. *Proceedings of the National Academy
 1133 of Sciences*, *120*(20), e2216158120. Retrieved from [https://www.pnas.org/doi/abs/
 1134 10.1073/pnas.2216158120](https://www.pnas.org/doi/abs/10.1073/pnas.2216158120) doi: 10.1073/pnas.2216158120
- 1135 Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., & Cui, P. (2021). Towards out-of-
 1136 distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*.
- 1137 Song, H.-J., & Roh, S. (2021). Improved weather forecasting using neural network emulation
 1138 for radiation parameterization. *Journal of Advances in Modeling Earth Systems*, *13*(10),
 1139 e2021MS002609. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/
 1140 abs/10.1029/2021MS002609](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002609) (e2021MS002609 2021MS002609) doi: [https://doi.org/
 1141 10.1029/2021MS002609](https://doi.org/10.1029/2021MS002609)
- 1142 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014).
 1143 Dropout: a simple way to prevent neural networks from overfitting. *The journal of
 1144 machine learning research*, *15*(1), 1929–1958.
- 1145 Stensrud, D. J. (2007). *Parameterization schemes: Keys to understanding numerical
 1146 weather prediction models*. Cambridge University Press. Retrieved from [https://
 1147 www.cambridge.org/core/product/identifier/9780511812590/type/book](https://www.cambridge.org/core/product/identifier/9780511812590/type/book) doi: 10

- 1148 .1017/CBO9780511812590
- 1149 Subel, A., Chattopadhyay, A., Guan, Y., & Hassanzadeh, P. (2021). Data-driven subgrid-
1150 scale modeling of forced Burgers turbulence using deep learning with generalization to
1151 higher Reynolds numbers via transfer learning. *Physics of Fluids*, *33*(3), 031702.
- 1152 Subel, A., Guan, Y., Chattopadhyay, A., & Hassanzadeh, P. (2023). Explaining the physics
1153 of transfer learning in data-driven turbulence modeling. *PNAS nexus*, *2*(3), pgad015.
- 1154 Sun, Y., Hassanzadeh, P., Alexander, M., & Kruse, C. (2023, 12). Quantifying 3d grav-
1155 ity wave drag in a library of tropical convection-permitting simulations for data-driven
1156 parameterizations.
1157 doi: 10.1029/2022MS003585
- 1158 Sun, Y., Wong, A., & Kamel, M. S. (2009, 11). Classification of imbalanced data: a review.
1159 *International Journal of Pattern Recognition and Artificial Intelligence*, *23*, 687-719. doi:
1160 10.1142/S0218001409007326
- 1161 Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., & Liu, C. (2018). A survey on deep
1162 transfer learning. , 270–279.
- 1163 Wedi, N. P., Polichtchouk, I., Dueben, P., Anantharaj, V. G., Bauer, P., Boussetta, S., ...
1164 others (2020). A baseline for global weather and climate simulations at 1 km resolution.
1165 *Journal of Advances in Modeling Earth Systems*, *12*(11), e2020MS002192.
- 1166 Wu, D., Gao, L., Xiong, X., Chinazzi, M., Vespignani, A., Ma, Y.-A., & Yu, R. (2021).
1167 *Quantifying uncertainty in deep spatiotemporal forecasting*.
- 1168 Wu, G., & Chang, E. Y. (2003). Adaptive feature-space conformal transformation for
1169 imbalanced-data learning. In *Proceedings of the twentieth international conference on
1170 international conference on machine learning* (p. 816–823). AAAI Press.
- 1171 Ye, H., Xie, C., Cai, T., Li, R., Li, Z., & Wang, L. (2021). Towards a theoretical framework
1172 of out-of-distribution generalization. *Advances in Neural Information Processing Systems*,
1173 *34*, 23519–23531.
- 1174 Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in
1175 deep neural networks? *Advances in neural information processing systems*, *27*.
- 1176 Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid
1177 processes for climate modeling at a range of resolutions. *Nature communications*, *11*(1),
1178 1–10.
- 1179 Zhang, C., Perezhugin, P., Gultekin, C., Adcroft, A., Fernandez-Granda, C., & Zanna, L.
1180 (2023). *Implementation and evaluation of a machine learned mesoscale eddy parameteri-
1181 zation into a numerical ocean circulation model*.
- 1182 Zhang, D., Ahuja, K., Xu, Y., Wang, Y., & Courville, A. (2021). Can subnetwork structure
1183 be the key to out-of-distribution generalization? In *International conference on machine
1184 learning* (pp. 12356–12367).
- 1185 Zhu, Y., Zabarar, N., Koutsourelakis, P.-S., & Perdikaris, P. (2019). Physics-constrained
1186 deep learning for high-dimensional surrogate modeling and uncertainty quantification
1187 without labeled data. *Journal of Computational Physics*, *394*, 56–81.