

1 **Machine Learning Emulation of Parameterized Gravity**
2 **Wave Momentum Fluxes in an Atmospheric Global**
3 **Climate Model**

4 **Zachary I. Espinosa¹, Aditi Sheshadri¹, Gerald R. Cain², Edwin P. Gerber³,**
5 **Kevin J. DallaSanta^{4,5}**

6 ¹Department of Earth System Science, Stanford University, Stanford, CA, USA

7 ²Department of Computer Science, Stanford University, Stanford, CA, USA

8 ³Courant Institute of Mathematical Sciences, New York University, New York, NY, USA

9 ⁴NASA Goddard Institute for Space Studies, New York, NY, USA

10 ⁵Universities Space Research Association, Columbia, MD, USA

11 **Key Points:**

- 12 • We train a machine learning model to emulate a parameterization of gravity wave
13 momentum transport
- 14 • The model reproduces the key features of the physics-based gravity wave param-
15 eterization
- 16 • The horizontal winds are the primary input features used by the model to gen-
17 erate gravity wave drag

Corresponding author: Zachary Espinosa, zespinos@stanford.edu

18 **Abstract**

19 We present a novel, single-column gravity wave parameterization (GWP) that uses ma-
 20 chine learning to emulate a physics-based GWP. An artificial neural network (ANN) is
 21 trained with output from an idealized atmospheric model and tested in an offline envi-
 22 ronment, illustrating that an ANN can learn the salient features of gravity wave momen-
 23 tum transport directly from resolved flow variables. We demonstrate that when trained
 24 on the westward phase of the Quasi-Biennial Oscillation, the ANN can skillfully gener-
 25 ate the momentum fluxes associated with the eastward phase. We also show that the merid-
 26 ional and zonal wind components are the only flow variables necessary to predict hor-
 27 izontal momentum fluxes with a globally and temporally averaged R^2 value over 0.8. State-
 28 of-the-art GWPs are severely limited by computational constraints and a scarcity of ob-
 29 servations for validation. This work constitutes a significant step towards obtaining ob-
 30 servationally validated, computationally efficient GWPs in global climate models.

31 **Plain Language Summary**

32 Atmospheric gravity waves (GWs) or “buoyancy waves” are generated by pertur-
 33 bations in a stably-stratified environment. They mediate momentum transport between
 34 the lower and middle atmospheres and play a leading-order role in driving middle atmo-
 35 spheric circulation. Due to computational constraints and a lack of observations, global
 36 climate models “parameterize” or crudely estimate the effect of GWs on the large-scale
 37 flow. Current climate predictions are sensitive to uncertainties in these representations,
 38 particularly at the regional scale. Here, we present a novel approach to parameterizing
 39 GWs by training a neural network to emulate an existing gravity wave parameterization
 40 in a global climate model. This approach represents an appealing technique to build data-
 41 driven gravity wave schemes that can reduce existing uncertainties.

42 **1 Introduction**

43 Atmospheric gravity waves (GWs) play a leading-order role in driving middle at-
 44 mospheric circulation, structure, and variability (Fritts & Alexander, 2003). By trans-
 45 porting momentum, GWs impact the mean climatology (Bretherton, 1969; Sato & Hi-
 46 rano, 2019). They are also critical for variability, impacting the jet stream and storm tracks
 47 (Fritts & Nastrom, 1992), stratospheric dynamics (Antonita et al., 2007; Kang et al., 2018;
 48 Limpasuvan et al., 2012), and processes such as stratospheric cloud formation (Hoffmann
 49 et al., 2017; S. Alexander et al., 2013). Since much of the wave spectrum (10^1 to 10^5 kilo-
 50 meters) is too fine to be captured at current model resolutions, models typically mimic
 51 its effects on the resolved circulation via explicit gravity wave parameterizations (GWPs)
 52 (Fritts & Alexander, 2003; Richter et al., 2010).

53 However, realistic representation of these effects in numerical models remains chal-
54 lenging for several reasons: i) The absolute magnitude of GW momentum flux is not par-
55 ticularly well-constrained by observational or intermodel studies (Geller et al., 2013). ii)
56 GWs are generated by a variety of sources, including orography, convection, and fron-
57 togenesis (Fritts & Alexander, 2003), but the representation of their sources is not uni-
58 form among models. iii) For a given source, the details of the GWP can vary greatly be-
59 tween models (Butchart et al., 2018), and even small changes within the same model can
60 lead to diverging regional climate projections (Schirber, 2015). iv) The horizontal prop-
61 agation of GWs is usually neglected in parameterizations, which is nonphysical and has
62 an impact on the middle atmosphere (Xu et al., 2017). v) There tends to be a compen-
63 sation between resolved and unresolved waves (Cohen et al., 2013), complicating obser-
64 vational and intermodel comparisons.

65 Despite these limitations, GWPs can be optimized towards maximizing global fore-
66 cast skill scores (Alexander et al., 2019) or “tuned” to reduce climatological biases (Garcia
67 et al., 2017). However, these limitations become apparent in simulations of future cli-
68 mate. Projections of the tropospheric and stratospheric circulations’ response to anthro-
69 pogenic forcing are sensitive to uncertainties in GWPs (Sigmond & Scinocca, 2010; Polichtchouk
70 et al., 2018). Such limitations indicate the need for the continued development of GWPs,
71 preferably with observational validation, computational efficiency, and minimal bias stem-
72 ming from underlying physical assumptions.

73 An alternative approach to physics-based parameterization has emerged in the form
74 of machine learning. For atmospheric sciences, machine learning has been employed to
75 parameterize processes such as convection (Rasp et al., 2018; Gentine et al., 2018) and
76 radiation (Brenowitz & Bretherton, 2018; Roh & Song, 2020), among other examples.
77 These applications are particularly valuable when algorithms derived from first princi-
78 ples (e.g., advection on the sphere) cannot be defined with great success. GWPs are thus
79 ripe for investigation with machine learning techniques, and relatively little work has been
80 done in this area. Matsuoka et al. (2020) made the important demonstration that a con-
81 volutional neural network can estimate the GW structure over Hokkaido, Japan when
82 trained on high-resolution reanalysis data.

83 Here, we perform a novel investigation into the efficacy of machine learning as a
84 GWP. We demonstrate that the drag due to breaking GWs can be faithfully represented
85 using an artificial neural network (ANN) that is trained using data from a global atmo-
86 spheric model embedded with an existing GWP. The machine learning model, which we
87 call WaveNet, differs from the Matsuoka et al. (2020) approach in that it is trained us-
88 ing output from an existing GWP, generates GW momentum tendencies rather than wind
89 anomalies, and provides global coverage. The value of a data-driven GWP lies in its com-
90 putational efficiency after training and its ability to be trained and evaluated using an
91 arbitrary amount of input data, which may include observations, reanalysis, and targeted

92 integrations (e.g., at ultra-high resolutions). The results presented here are a natural first
 93 step towards developing a computationally efficient, three-dimensional, observationally
 94 validated GWP.

95 **2 Data**

96 We train the ANN to emulate the M. Alexander and Dunkerton (1999) GWP as
 97 incorporated into an atmospheric model of intermediate complexity, the Model of an Ide-
 98 alized Moist Atmosphere (MiMA; Jucker and Gerber (2017)). This choice of GWP and
 99 model were based on two factors. First, it is desirable to use a model and GWP that pro-
 100 duce somewhat realistic GW behavior. MiMA captures key dynamical features of the
 101 stratosphere-troposphere system that depend critically on GWs at a resolution compa-
 102 rable to state-of-the-art stratosphere resolving global climate models (GCMs). It also
 103 produces a realistic representation of stratospheric variability (i.e., the frequency and in-
 104 tensity of sudden stratospheric warmings and a self-generated Quasi-Biennial Oscilla-
 105 tion (QBO)). Second, it is advantageous to limit additional degrees of freedom that could
 106 cause underfitting in the ANN. MiMA’s key simplification relative to a comprehensive
 107 atmospheric model lies in its idealized treatment of the hydrological cycle, its lower bound-
 108 ary (a purely thermodynamic, or slab ocean), and the absence of cloud and aerosol pro-
 109 cesses.

110 MiMA is integrated with T42 spectral resolution (triangular truncation at wavenum-
 111 ber 42, roughly equivalent to a 2.8-degree grid) in the horizontal and 40 vertical levels,
 112 with a model lid at 0.18 hPa and 23 levels above 100 hPa. Following Garfinkel et al. (2020),
 113 its simple thermodynamic ocean includes a crude parameterization of oceanic heat trans-
 114 port specified by steady heat flux within the oceanic layer, often referred to as a “Q-flux”.

115 The implementation of the M. Alexander and Dunkerton (1999) GWP, hereafter
 116 referred to as the AD99 scheme, is based on its formulation within the GFDL Flexible
 117 Modeling System. It is the same scheme employed by GFDL Atmospheric Model 3 (Donner
 118 et al., 2011), except modified by Cohen et al. (2013) to ensure that no momentum flux
 119 escapes the model lid. Here, all momentum that reaches 0.85 hPa is uniformly distributed
 120 to levels layers above the stratopause. The scheme assumes a Gaussian spectrum of GWs,
 121 discretized as a function of phase speed and launched from the upper troposphere (315
 122 hPa). A broad spectrum of phase speeds is chosen, with a half width of 35 m/s, and cen-
 123 tered around the speed of the zonal wind at the launch level. The total amplitude of the
 124 momentum stress is 0.0043 Pa. The width and momentum stress were optimized to sim-
 125 ulate the QBO, but still provide a reasonable representation of waves in the extratrop-
 126 ics. The broad spectrum ensures that there is a rich source of GWs at MiMA’s lower bound-
 127 ary, and the amplitude and variability of the extratropical polar vortices are well cap-
 128 tured in the simulation. The momentum associated with each wave is deposited at its
 129 linear breaking level. The intermittency (i.e., highly skewed distribution) of observed GWs

130 is taken into account via a scaling parameter: the breaking level is based on the behav-
 131 ior of a large wave, indicative of the median amplitude of the distribution, but the mo-
 132 mentum flux is rescaled to provide the momentum deposition associated with an aver-
 133 age wave. This better captures the true breaking level of GWs, but smooths out the mo-
 134 mentum deposition in time.

135 The chief idealization of the scheme is in our choice of source spectrum. We assume
 136 a uniform spectrum that varies only with respect to the zonal winds at the source level.
 137 This crudely accounts for filtering of the wave spectrum by the troposphere and was crit-
 138 ical for the simulation of the QBO. Such fixed source schemes are widely employed in
 139 atmospheric models (Butchart et al., 2018), but are highly simplified relative to real GW
 140 sources.

141 We utilized five years of one integration of MiMA, yielding around 12 million train-
 142 ing samples per year, with one year of six-hourly data representing approximately 20 Gb.
 143 The second year of output is used for training, while all others are reserved for testing.
 144 Not all runs utilize the full year of training data, and discrepancies are specified per ex-
 145 periment. All training data are standardized by removing the mean and scaling to unit
 146 variance, calculated for each variable across each pressure level. Test data was similarly
 147 standardized using the mean and variance calculated for training data. Output variables
 148 were inverse transformed before presentation.

149 **3 Neural Network Architecture**

150 An ANN is a computing system of interconnected layers of computational nodes.
 151 Each node is comprised of a linear component, which adds a bias parameter to the in-
 152 ner product of a feature vector and a learnable weight vector. A nonlinear component
 153 then maps the linear output to an activation function. The resulting scalar from each
 154 node in a layer is passed as input to each node in the subsequent layer. With incremen-
 155 tal weight and bias adjustments, an ANN attempts to approximate nonlinear relation-
 156 ships between input and output features. The first layer, or input layer, of WaveNet ac-
 157 cepts a stacked vector representing a single vertical column of resolved flow variable out-
 158 put from MiMA. The last layer, or output layer, produces a stacked vector of gravity wave
 159 drag (GWD) generated by the ANN for a single vertical column (Table S1). The lay-
 160 ers between the first and last layers are called hidden layers. We trained two ANNs, one
 161 to generate zonal drag and another to generate meridional drag. Both ANNs contain four
 162 hidden layers, each with 256 nodes. The fourth hidden layer splits into 33 branches, one
 163 for each nontrivial vertical level (the first seven layers are below the source level). We
 164 allow the network to use resolved flow variables from below the source level since this
 165 may relate to the drag above. Each branch contains four pressure-level specific hidden
 166 layers containing 256, 128, 64, and 32 nodes. The final pressure-level specific layer feeds
 167 into an output layer. The result of each node in the output layer does not pass through

168 an activation function and produces a GWD value for the vertical column (Figure S1).
 169 For all other nodes, we use the Rectified Linear Unit (ReLU). In total, this ANN archi-
 170 tecture contains 3,848,481 trainable parameters when using the full set of resolved flow
 171 variables and shifts slightly, with a lower limit of 3,806,753, when training with a sub-
 172 set of resolved flow variables. We did not perform an analysis of performance sensitiv-
 173 ity to the number of learnable parameters. Rather, we followed standard literature sug-
 174 gesting that with an abundance of data available, deeper neural networks generally pro-
 175 duce better scores than shallow networks (Liang & Srikant, 2016).

176 During training, the ANN attempts to minimize the loss by incrementally nudg-
 177 ing each trainable parameter by a scaled version of the gradient of the loss with respect
 178 to that parameter. The loss function - in our case, the logcosh error - is computed for
 179 a minibatch of 1,024 training samples that are drawn from a pseudo-shuffled training dataset.
 180 The logcosh error is defined as

$$L(y, y^p) = \sum_{i=1}^n \log(\cosh(y_i^p - y_i))$$

181 where n is the size of the dataset and y_i and y_i^p are the i^{th} truth and prediction, respec-
 182 tively. Parameter updates are performed according to Adam optimization (Kingma &
 183 Ba, 2014). We started each training session with a learning rate of 10^{-3} and reduced it
 184 when improvement plateaued for more than 5 epochs, with an epoch defined as a sin-
 185 gle pass through the training data. We stopped training when performance plateaued
 186 for more than 10 epochs, which occurred at 200 epochs on average. We did not perform
 187 any regularization. All weights are initialized using Xavier initialization (Glorot & Ben-
 188 gio, 2010); however, initialization between runs proved to have no impact on the final
 189 results between training events.

190 4 Results

191 4.1 Evaluation of ANN Predictions

192 We start by training WaveNet on one year of MiMA output and testing it on the
 193 three years proceeding and one year preceding the training period. All subsequent anal-
 194 yses are completed using our best performing networks and the full vertical column of
 195 resolved flow variables, unless otherwise noted. Figure 1 shows strong similarities on a
 196 global scale between the zonal and meridional GWD generated by WaveNet and AD99
 197 at 10 hPa (panels a through f) and 100 hPa (panels g through l) for a single time step.
 198 Similar similarities are seen at all vertical layers and across time (Movie S1 and S2). All
 199 reported tests are conducted “offline”, such that WaveNet is not directly coupled to MiMA
 200 (i.e. WaveNet’s output is not used by MiMA to generate data at subsequent time steps).
 201

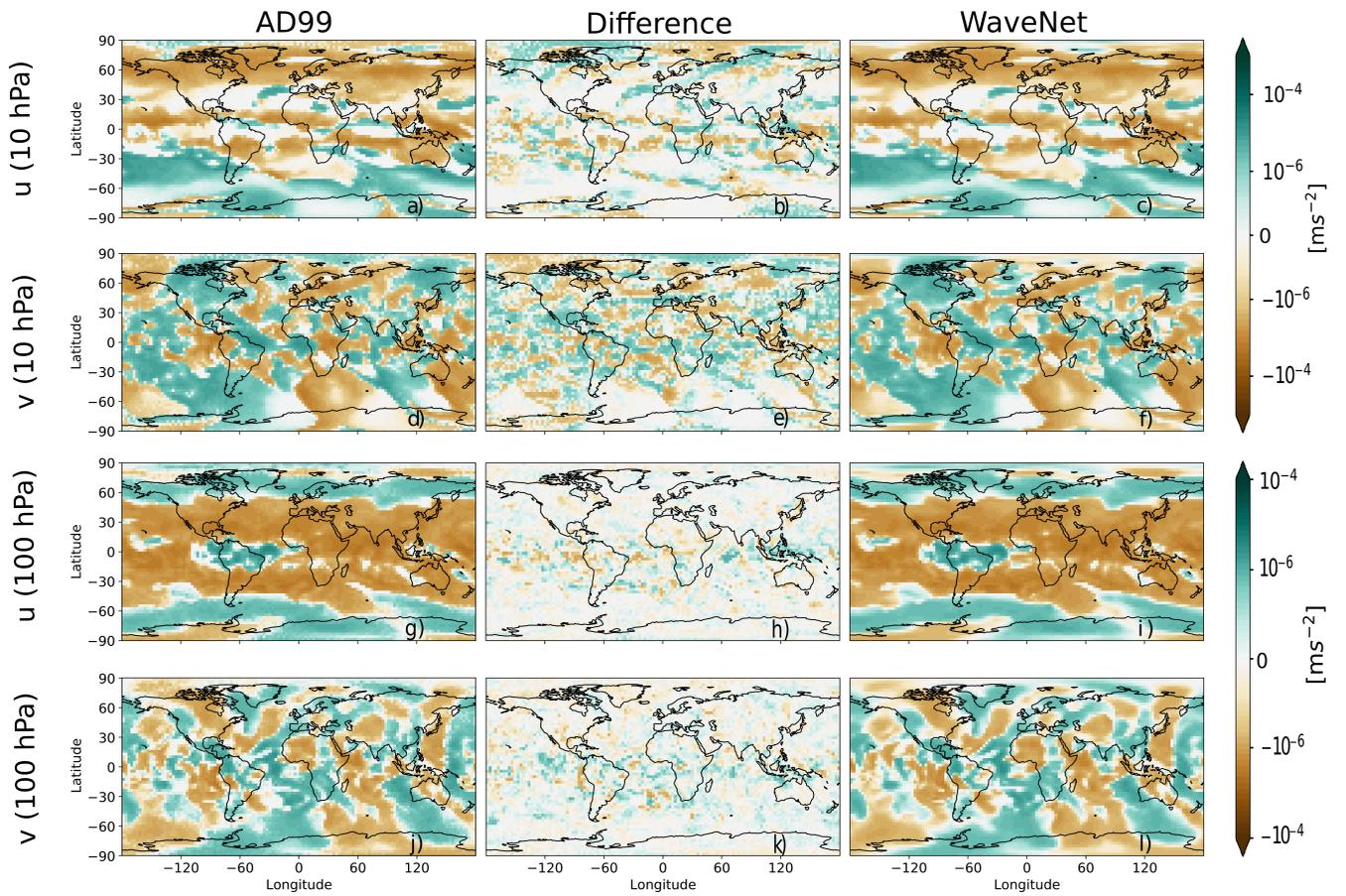


Figure 1. Latitude-longitude snapshot of zonal and meridional AD99 generated drag (panels a, d, g, and j), the corresponding ANN predictions in an offline test (panels c, f, i, and l), and their difference (panels b, e, h, and k) at model level 13 (10 hPa) and 23 (100 hPa) for one time step in the test set.

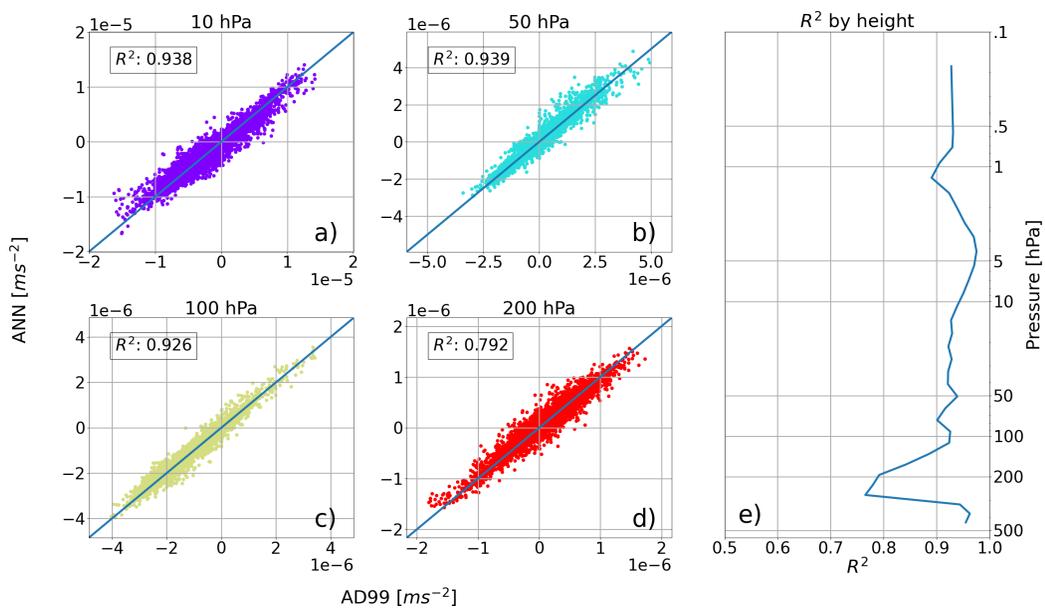


Figure 2. Panels a) through d) show ANN zonal predictions versus AD99 truth at 10 hPa, 50 hPa, 100 hPa and 200 hPa, respectively, for 10k samples in the test set. Panel e) shows pressure versus horizontal and time averaged R^2 values generated from one year of test data for zonal predictions.

202 To evaluate the quality of predictions, we calculate the R^2 coefficient of determi-
 203 nation averaged over time and horizontal dimensions for each vertical level (Figure 2).
 204 R^2 is defined as one minus the proportion of the sum of squares of residuals to the to-
 205 tal sum of squares (this is also equal to the square of the correlation coefficient):

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - \bar{y})^2}$$

206 where y_i and f_i are the AD99 and WaveNet generated drag for the i^{th} sample, respec-
 207 tively, and \bar{y} is the globally and temporally averaged AD99 generated drag. Figure 2 shows
 208 that the ANN can skillfully generate GWD similar to that generated by AD99. WaveNet
 209 produces R^2 values for its zonal and meridional predictions averaged across space and
 210 time of .92 and .85, respectively. For all metrics and across experiments, WaveNet per-
 211 forms better on zonal tendencies than meridional tendencies (Figure S2). The variance
 212 of meridional GWD generated by AD99 is smaller in absolute magnitude than the vari-
 213 ance of zonal GWD but greater relative to its mean at all pressure levels. As a result,
 214 meridional drag is likely more difficult to learn.

215 In order to assess how well WaveNet captures GWD associated with large-scale cir-
 216 culation, we analyze the vertical profile of equatorial drag tendencies (Figure 3). WaveNet
 217 is trained on one year of global data (months 12-24), containing the westward phase of
 218 the QBO. The ANN captures the changes in GW driving associated with the eastward
 219 phases preceding and proceeding the training period. While this is an offline test, this

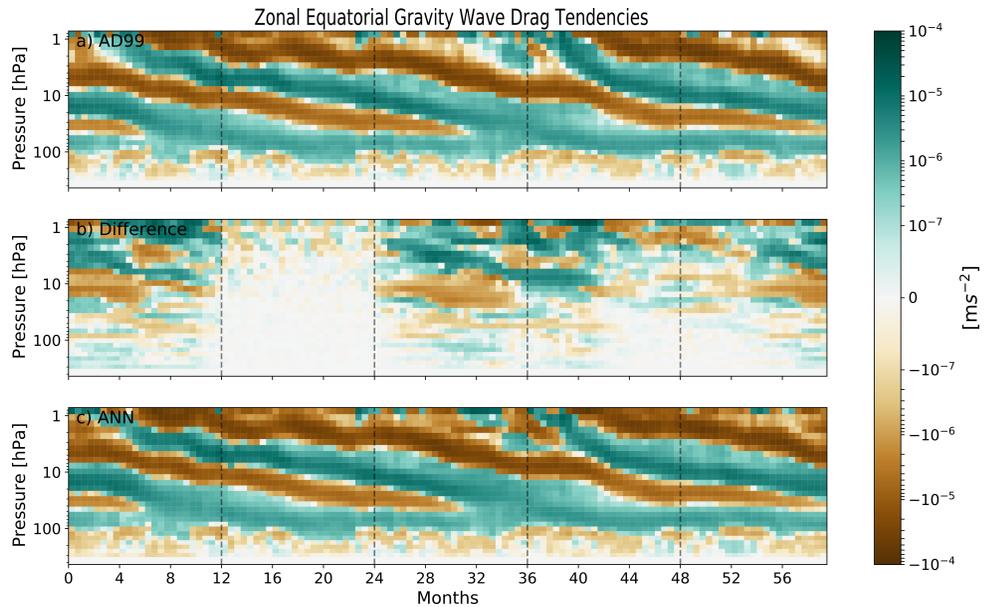


Figure 3. Pressure-Time profile of zonal mean drag 15-day tendencies at 0° latitude for zonal AD99 generated GW drag (a), the corresponding ANN predictions (c), and their difference (b). Vertical dashed lines separate years. The ANN is trained on months 12-24 and tested on all other months. The westward (brown) and eastward (green) bands correspond to drag associated with breaking GWs in opposite phases of the QBO.

220 result suggests that WaveNet can generalize outside of its training sample. That is, it
 221 can transfer learning from other regions (horizontally or vertically) containing eastward
 222 wind samples to a region where it has not experienced similar samples during the train-
 223 ing period. While comparison between the top and bottom panels of Figure 3 shows that
 224 WaveNet can capture the gross features of the “out-of-sample” periods, differences be-
 225 tween AD99 and WaveNet are larger at these times, as seen in the middle panel. No-
 226 tably, the errors decrease around month 48, when the winds are more similar in struc-
 227 ture to the training period. This suggests more varied training data may improve WaveNet’s
 228 performance. Nonetheless, this result is reassuring with respect to WaveNet’s potential
 229 to globally generalize from regional observations or datasets generated from high-resolution
 230 simulations, a subject of future studies. Furthermore, it suggests that the ANN does not
 231 trivially depend on the source function’s artificial behavior. A similar analysis performed
 232 at 60°N reveals that WaveNet captures the seasonal cycle of GWD associated with the
 233 polar vortex, further supporting the claim that WaveNet can capture large-scale circu-
 234 lation patterns (Figure S3).

235 4.2 Interpretability of ANN

236 Although ANNs have achieved great success in a range of applications, their lack
 237 of interpretability has become a significant obstacle to their widespread acceptance and
 238 made them generally unsuitable for conceptual model building. Here, WaveNet is trained
 239 using ten distinct subsets of the full resolved flow variable set (Figure 4ab) to determine

240 which features are most critical. Figure 4a and 4b show horizontally and temporally av-
 241 eraged R^2 values by height per feature set for zonal and meridional drag predictions. We
 242 conclude that the zonal and meridional wind components are the only flow variables nec-
 243 essary to predict horizontal drag with a globally and temporally averaged R^2 value over
 244 0.8. Additional training features mildly improve performance, with varied effects at dif-
 245 ferent pressure levels. Note that with the exclusion of latitude and longitude features,
 246 the general performance of WaveNet does not significantly drop. These results are in agree-
 247 ment with those in subsection 4.1, in that the ANN is learning to emulate GW dissipa-
 248 tion rather than climatological GWD properties or a function of latitude. For these ex-
 249 periments, the total number of trainable parameters varies by roughly 1.0%, with fewer
 250 input features corresponding to fewer parameters. This variance did not impact the rel-
 251 ative performance of each experiment.

252 As a preliminary analysis of the role of horizontal wind components in predictions,
 253 we calculate Shapley Additive Explanations (SHAP; Lundberg and Lee (2017)). SHAP
 254 values represent the effect a feature has on the model’s prediction if that feature is in-
 255 cluded in the input. To compute SHAP values, the model is retrained on all feature sub-
 256 sets $S \subseteq F$, where F is the set of all features. The SHAP value for a feature is then the
 257 weighted sum of the conditional expectation of the marginal contribution of including
 258 that feature in the prediction. The SHAP values for this study are calculated using Deep
 259 SHAP, an approximation technique that combines DeepLIFT with Shapley values from
 260 collinear cooperative game theory (Shrikumar et al., 2017; Lundberg & Lee, 2017). This
 261 approximation technique avoids retraining the network N times, where N is the power
 262 set of F . Unsurprisingly, the results show that on average the horizontal wind compo-
 263 nents on levels directly above and below the level of prediction have the largest contri-
 264 bution to the model’s prediction (Figure S4). This is consistent with our physical un-
 265 derstanding of GWs, where dissipation is linked to critical levels (i.e., where the phase
 266 speed of a wave is equal to the speed of the background flow). This spatially local de-
 267 pendence suggests that our approach may generalize well to observational datasets that
 268 contain measurements at a small range of pressure levels, e.g., observations from super-
 269 pressure balloons (Podglajen et al., 2016; Lindgren et al., 2020).

270 4.3 Sensitivity to Amount of Training Data

271 To understand how performance degrades as less training data is made available
 272 to the ANN, we incrementally decrease the number of training samples from one year
 273 to one day (Figure 4cd). To account for the effects of seasonality, for each test, we sam-
 274 ple uniformly in space and time from one year of training data to generate a subset with
 275 coverage of the entire seasonal cycle. While more data generally leads to better perfor-
 276 mance, we find that one fourth of one year of data (equivalent to three months or roughly
 277 2.9 million training samples) is sufficient to learn most of the salient features of AD99.

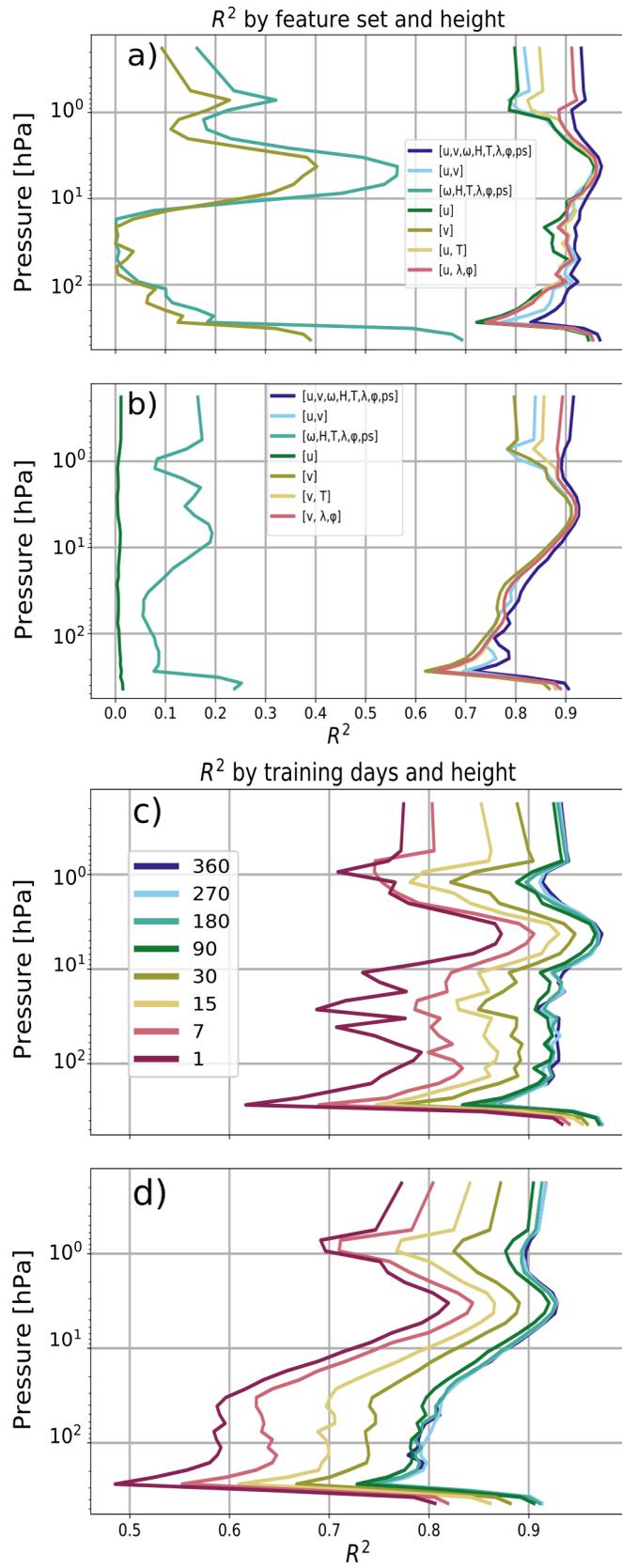


Figure 4. Panels a) through d) show horizontally and temporally averaged R^2 values computed for each model pressure level for either each feature experiment (a and b) or for each data availability experiment (c and d).

278 This suggests the promise of a data-driven GWP to be trained and validated on obser-
 279 vational datasets of similar resolution and size, as well as to generalize to simulations that
 280 more realistically model the physics of GW generation, propagation, and dissipation. Ad-
 281 ditionally, given the deep architecture of WaveNet, that no effort was made to improve
 282 performance at each duration, and that it is known that deep networks require large train-
 283 ing sets, three months can be regarded as a threshold beyond which optimal performance
 284 is likely for an ANN similar to WaveNet.

285 5 Discussion and Conclusion

286 We have demonstrated that an ANN can skillfully learn the salient features of GW
 287 momentum transport directly from resolved flow variables. The concept was demonstrated
 288 in an idealized setting using the M. Alexander and Dunkerton (1999) GWP in a simpli-
 289 fied atmospheric model, MiMA. We have shown that the most important input features
 290 for WaveNet’s predictions are the horizontal wind components local to the vertical level
 291 of prediction. Moreover, WaveNet is skillfully able to reproduce drag associated with the
 292 eastward phase of the QBO after being trained on data representing the westward phase.
 293 In doing so, we have demonstrated that WaveNet can spatially and temporally gener-
 294 alize. The success seen in this context implies that an approach like WaveNet may open
 295 a new avenue by which the advantages of high-resolution GW simulations (Remmler et
 296 al., 2015) or observational datasets (Lindgren et al., 2020) can be incorporated into cur-
 297 rent GCMs.

298 There are, however, a number of challenges that may emerge before the advantages
 299 of an approach like WaveNet can be fully realized in a GCM. First, many studies (e.g.,
 300 Brenowitz and Bretherton (2018)) have shown that machine learning schemes which per-
 301 form very well offline, i.e., reproducing the correct tendencies, given the correct model
 302 state, do not work as well (or at all) when the scheme is coupled with a GCM in an “on-
 303 line” integration. Second, ANNs do not inherently conserve energy or momentum, and
 304 additional assumptions may be made to conserve these quantities: for example, artifi-
 305 cially scaling the positive and negative fluxes, or depositing the remaining momentum
 306 at pre-determined levels. Third, the lack of interpretability of ANNs may serve as a sub-
 307 stantial barrier to their widespread adoption. Additional effort is necessary to consider
 308 how WaveNet’s behavior may relate to the GW dispersion relations. Fourth, WaveNet
 309 is an extremely large network. In order to make coupling WaveNet with a GCM com-
 310 putationally feasible, a cost-performance analysis should be performed to reduce WaveNet’s
 311 complexity. Finally, a next test is to examine how WaveNet generalizes when trained on
 312 regional datasets and orographic and nonorographic GWPs.

313 Nevertheless, our results suggest that machine learning may represent a powerful
 314 alternative to existing GWPs. An approach like WaveNet is naturally suited for data as-
 315 simulation, and WaveNet may be completely or partially trained and validated using ob-

316 servational datasets. Moreover, existing GWP ignore horizontal GW propagation due
 317 to computational limitations. A machine learning approach such as WaveNet may af-
 318 ford the computational efficiency needed to develop a three-dimensional GWP. Projec-
 319 tions of the large-scale climate response to anthropogenic warming are sensitive to un-
 320 certainties in existing GWPs. Developing an observationally validated, three-dimensional
 321 GWP may more accurately capture the physics of GWs. The approach presented here
 322 constitutes a first step toward obtaining such GWPs for global climate prediction.

323 **Acknowledgments**

324 We thank Joan Alexander and Pedram Hassanzadeh for useful discussions. ZIE and AS
 325 acknowledge support from the NSF through grant OAC-2004492. EPG acknowledges sup-
 326 port from the NSF through grant OAC-2004572. KD was funded by the NASA Post-
 327 doctoral Program at the Goddard Institute for Space Studies.

328 **Data Availability**

329 The ANN is built using Keras, a deep learning framework that wraps Tensorflow.
 330 All source code is available at <https://doi.org/10.5281/zenodo.4428931>. Training
 331 took on the order of 48 hr on a graphical processing unit and varied according to data
 332 size and trainable parameters. The implementation of SHAP values is available at [https://](https://github.com/slundberg/shap)
 333 github.com/slundberg/shap and is not maintained or owned by this project group. MiMA
 334 is documented by Jucker and Gerber (2017) and Garfinkel et al. (2020). The source code,
 335 run parameters, and modified configuration for MiMA are available at: [https://doi.org/](https://doi.org/10.5281/zenodo.1401407)
 336 [10.5281/zenodo.1401407](https://doi.org/10.5281/zenodo.1401407)

337 **References**

- 338 Alexander, Bacmeister, J., Ern, M., Gisinger, S., Hoffmann, L., Holt, L., . . . Wright,
 339 C. (2019). *Seeking new quantitative constraints on orographic gravity*
 340 *wave stress and drag to satisfy emerging needs in seasonal-to-subseasonal*
 341 *and climate prediction*. Retrieved from [http://www.issibern.ch/teams/](http://www.issibern.ch/teams/consonorogravity/#)
 342 [consonorogravity/#](http://www.issibern.ch/teams/consonorogravity/#)
- 343 Alexander, M., & Dunkerton, T. (1999). A spectral parameterization of mean-flow
 344 forcing due to breaking gravity waves. *Journal of the Atmospheric Sciences*,
 345 *56*(24), 4167–4182.
- 346 Alexander, S., Klekociuk, A., McDonald, A., & Pitts, M. (2013). Quantifying the
 347 role of orographic gravity waves on polar stratospheric cloud occurrence in
 348 the antarctic and the arctic. *Journal of Geophysical Research: Atmospheres*,
 349 *118*(20), 11–493.
- 350 Antonita, T. M., Ramkumar, G., Kumar, K. K., Appu, K., & Nambhoodiri, K.
 351 (2007). A quantitative study on the role of gravity waves in driving the trop-

- 352 ical stratospheric semiannual oscillation. *Journal of Geophysical Research:*
353 *Atmospheres*, 112(D12).
- 354 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural net-
355 work unified physics parameterization. *Geophysical Research Letters*, 45(12),
356 6289–6298.
- 357 Bretherton, F. P. (1969). Momentum transport by gravity waves. *Quarterly Journal*
358 *of the Royal Meteorological Society*, 95(404), 213–243.
- 359 Butchart, N., Anstey, J., Hamilton, K., Osprey, S., McLandress, C., Bushell, A.,
360 ... others (2018). Overview of experiment design and comparison of models
361 participating in phase 1 of the sparc quasi-biennial oscillation initiative (qboi).
362 *Geoscientific Model Development*, 11(3).
- 363 Cohen, N. Y., Gerber, & Oliver Bühler, E. P. (2013). Compensation between
364 resolved and unresolved wave driving in the stratosphere: Implications for
365 downward control. *Journal of the atmospheric sciences*, 70(12), 3780–3798.
- 366 Donner, L. J., Wyman, B. L., Hemler, R. S., Horowitz, L. W., Ming, Y., Zhao, M.,
367 ... others (2011). The dynamical core, physical parameterizations, and basic
368 simulation characteristics of the atmospheric component am3 of the gfdl global
369 coupled model cm3. *Journal of Climate*, 24(13), 3484–3519.
- 370 Fritts, D. C., & Alexander, M. J. (2003). Gravity wave dynamics and effects in the
371 middle atmosphere. *Reviews of geophysics*, 41(1).
- 372 Fritts, D. C., & Nastrom, G. D. (1992). Sources of mesoscale variability of grav-
373 ity waves. part ii: Frontal, convective, and jet stream excitation. *Journal of the*
374 *Atmospheric Sciences*, 49(2), 111–127.
- 375 Garcia, R. R., Smith, A. K., Kinnison, D. E., Cámara, Á. d. l., & Murphy, D. J.
376 (2017). Modification of the gravity wave parameterization in the whole at-
377 mosphere community climate model: Motivation and results. *Journal of the*
378 *Atmospheric Sciences*, 74(1), 275–291.
- 379 Garfinkel, C. I., White, I., Gerber, E. P., Jucker, M., & Erez, M. (2020). The build-
380 ing blocks of northern hemisphere wintertime stationary waves. *Journal of Cli-*
381 *mate*, 33(13), 5611–5633.
- 382 Geller, M. A., Alexander, M. J., Love, P. T., Bacmeister, J., Ern, M., Hertzog, A.,
383 ... others (2013). A comparison between gravity wave momentum fluxes in
384 observations and climate models. *Journal of Climate*, 26(17), 6383–6405.
- 385 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could
386 machine learning break the convection parameterization deadlock? *Geophysical*
387 *Research Letters*, 45(11), 5742–5751.
- 388 Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feed-
389 forward neural networks. In *Proceedings of the thirteenth international confer-*
390 *ence on artificial intelligence and statistics* (pp. 249–256).
- 391 Hoffmann, L., Spang, R., Orr, A., Alexander, M. J., Holt, L. A., & Stein, O. (2017).

- 392 A decadal satellite record of gravity wave activity in the lower stratosphere
 393 to study polar stratospheric cloud formation. *Atmospheric Chemistry and*
 394 *Physics*, 17(4), 2901–2920.
- 395 Jucker, M., & Gerber, E. (2017). Untangling the annual cycle of the tropical
 396 tropopause layer with an idealized moist model. *Journal of Climate*, 30(18),
 397 7339–7358.
- 398 Kang, M.-J., Chun, H.-Y., Kim, Y.-H., Preusse, P., & Ern, M. (2018). Momentum
 399 flux of convective gravity waves derived from an offline gravity wave param-
 400 eterization. part ii: Impacts on the quasi-biennial oscillation. *Journal of the*
 401 *Atmospheric Sciences*, 75(11), 3753–3775.
- 402 Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv*
 403 *preprint arXiv:1412.6980*.
- 404 Liang, S., & Srikant, R. (2016). Why deep neural networks for function approxima-
 405 tion? *arXiv preprint arXiv:1610.04161*.
- 406 Limpasuvan, V., Richter, J. H., Orsolini, Y. J., Stordal, F., & Kvissel, O.-K. (2012).
 407 The roles of planetary and gravity waves during a major stratospheric sud-
 408 den warming as characterized in waccm. *Journal of atmospheric and solar-*
 409 *terrestrial physics*, 78, 84–98.
- 410 Lindgren, E. A., Sheshadri, A., Podglajen, A., & Carver, R. W. (2020). Sea-
 411 sonal and latitudinal variability of the gravity wave spectrum in the lower
 412 stratosphere. *Journal of Geophysical Research: Atmospheres*, 125(18),
 413 e2020JD032850.
- 414 Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model
 415 predictions. In *Advances in neural information processing systems* (pp. 4765–
 416 4774).
- 417 Matsuoka, D., Watanabe, S., Sato, K., Kawazoe, S., Yu, W., & Easterbrook, S.
 418 (2020). Application of deep learning to estimate atmospheric gravity wave
 419 parameters in reanalysis data sets. *Geophysical Research Letters*, 47(19),
 420 e2020GL089436.
- 421 Podglajen, A., Hertzog, A., Plougonven, R., & Legras, B. (2016). Lagrangian tem-
 422 perature and vertical velocity fluctuations due to gravity waves in the lower
 423 stratosphere. *Geophysical Research Letters*, 43(7), 3543–3553.
- 424 Polichtchouk, I., Shepherd, T. G., & Byrne, N. J. (2018). Impact of parametrized
 425 nonorographic gravity wave drag on stratosphere-troposphere coupling in the
 426 northern and southern hemispheres. *Geophysical Research Letters*, 45(16),
 427 8612–8618.
- 428 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid
 429 processes in climate models. *Proceedings of the National Academy of Sciences*,
 430 115(39), 9684–9689.
- 431 Remmler, S., Hickel, S., Fruman, M. D., & Achatz, U. (2015). Validation of large-

- 432 eddy simulation methods for gravity wave breaking. *Journal of the Atmo-*
433 *spheric Sciences*, *72*(9), 3537–3562.
- 434 Richter, J. H., Sassi, F., & Garcia, R. R. (2010). Toward a physically based gravity
435 wave source parameterization in a general circulation model. *Journal of the At-*
436 *mospheric Sciences*, *67*(1), 136–156.
- 437 Roh, S., & Song, H.-J. (2020). Evaluation of neural network emulations for radia-
438 tion parameterization in cloud resolving model. *Geophysical Research Letters*,
439 *47*(21), e2020GL089444.
- 440 Sato, K., & Hirano, S. (2019). The climatology of the brewer–dobson circulation and
441 the contribution of gravity waves. *Atmos. Chem. Phys*, *19*, 4517–4539.
- 442 Schirber, S. (2015). Influence of enso on the qbo: Results from an ensemble of ideal-
443 ized simulations. *Journal of Geophysical Research: Atmospheres*, *120*(3), 1109–
444 1122.
- 445 Shrikumar, A., Greenside, P., & Kundaje, A. (2017). Learning important features
446 through propagating activation differences. *arXiv preprint arXiv:1704.02685*.
- 447 Sigmond, M., & Scinocca, J. F. (2010). The influence of the basic state on the
448 northern hemisphere circulation response to climate change. *Journal of cli-*
449 *mate*, *23*(6), 1434–1446.
- 450 Xu, X., Wang, Y., Xue, M., & Zhu, K. (2017). Impacts of horizontal propagation
451 of orographic gravity waves on the wave drag in the stratosphere and lower
452 mesosphere. *Journal of Geophysical Research: Atmospheres*, *122*(21), 11–301.

Supporting Information for “Machine Learning Emulation of Parameterized Gravity Wave Momentum Fluxes in an Atmospheric Global Climate Model”

Zachary I. Espinosa¹, Aditi Sheshadri¹, Gerald R. Cain², Edwin P. Gerber³,
Kevin J. DallaSanta^{4,5}

¹Department of Earth System Science, Stanford University, Stanford, CA, USA

²Department of Computer Science, Stanford University, Stanford, CA, USA

³Courant Institute of Mathematical Sciences, New York University, New York, NY, USA

⁴NASA Goddard Institute for Space Studies, New York, NY, USA

⁵Universities Space Research Association, Columbia, MD, USA

Contents of this file

1. Table S1
2. Figures S1-S4

Additional Supporting Information (Files uploaded separately)

1. Captions for Movies S1 and S2

Corresponding author: Zachary I. Espinosa, (zespinos@stanford.edu)

Introduction The supporting information includes one table, four figures and two movies. Table S1 shows a list of input variables accepted by the ANN and output variables generated by the ANN. Figure S1 shows a schematic of the ANN architecture. Figure S2 shows ANN meridional predictions versus AD99 truth at four model levels and a globally and temporally averaged R^2 -Pressure plot for one year of test data. Figure S3 is a pressure-time profile of zonal mean drag 15-day tendencies at 60N for zonal ANN predictions, the corresponding AD99 truths, and their difference. Figure S4 shows SHAP bar plots of the ten most important meridional and zonal wind features used by WaveNet to generate GWD at 10 hPa and 100 hPa. Movies S1 and S2 are a time series animation of Figure 2 at 10 hPa, for zonal and meridional predictions, respectively.

Movie S1. A latitude-longitude time series of zonal ANN predictions, AD99 truth, and their difference at model level 13 (10 hPa) for half a year of test data. Panels a through c in Figure 1 are a single snapshot of this animation.

Movie S2. A latitude-longitude time series of meridional ANN predictions, AD99 truth, and their difference at model level 13 (10 hPa) for half a year of test data. Panels d through f in Figure 1 are a single snapshot of this animation.

Table S1. List of input variables accepted by the ANN and output variables generated by the ANN. The total input feature vector contains 203 elements, and the output is 33 GWD values. Two networks are trained, one for zonal drag and one for meridional drag.

List of Input and Output Variables Used for ANN			
Input Variables	Vertical Levels	Output Variables	Vertical Levels
Zonal Wind ($\frac{m}{s}$)	40	GW zonal drag ($\frac{m}{s^2}$)	33
Meridional Wind ($\frac{m}{s}$)	40	GW meridional drag ($\frac{m}{s^2}$)	33
Vertical Wind ($\frac{m}{s}$)	40		
Temperature (K)	40		
Height (m)	40		
Latitude (λ)	1		
Longitude (ϕ)	1		
Surface Pressure (hPa)	1		
Size of Stacked Array	203		33

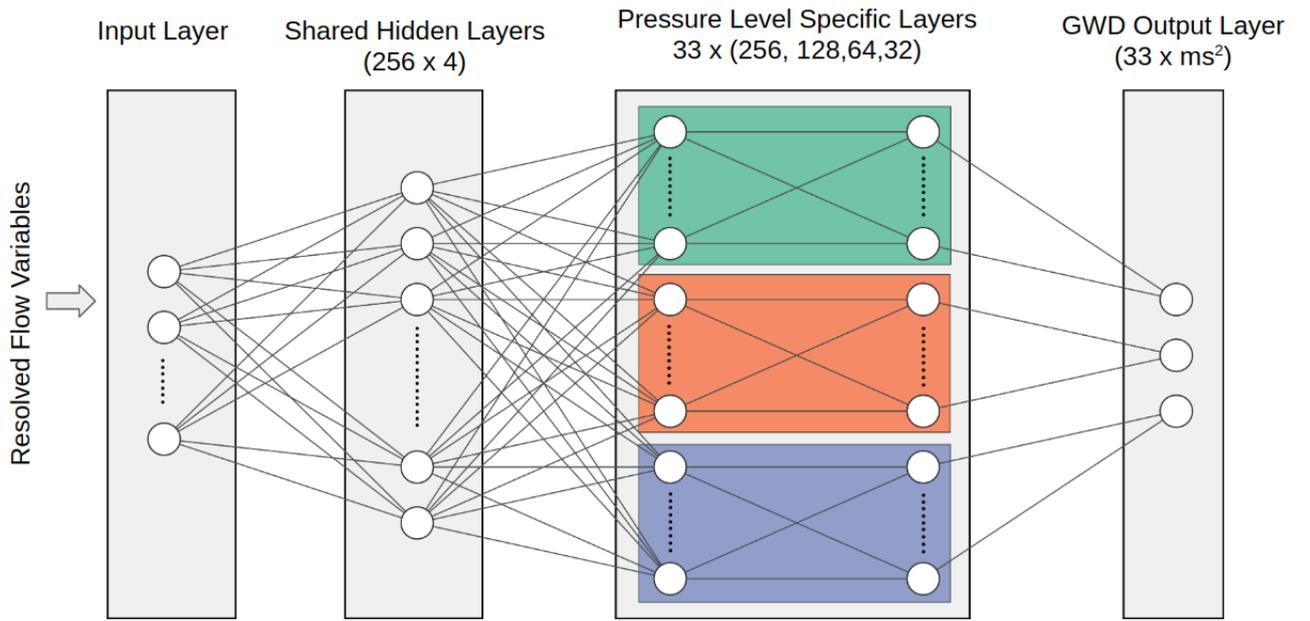


Figure S1. WaveNet contains 4 shared hidden layers, each with 256 neurons. WaveNet then splits into 33 branches (one branch per nontrivial vertical layer) each containing 4 pressure level specific layers with 256, 128, 64, and 32 neurons, respectively. Each branch then outputs a single value corresponding to the gravity wave drag at that vertical layer.

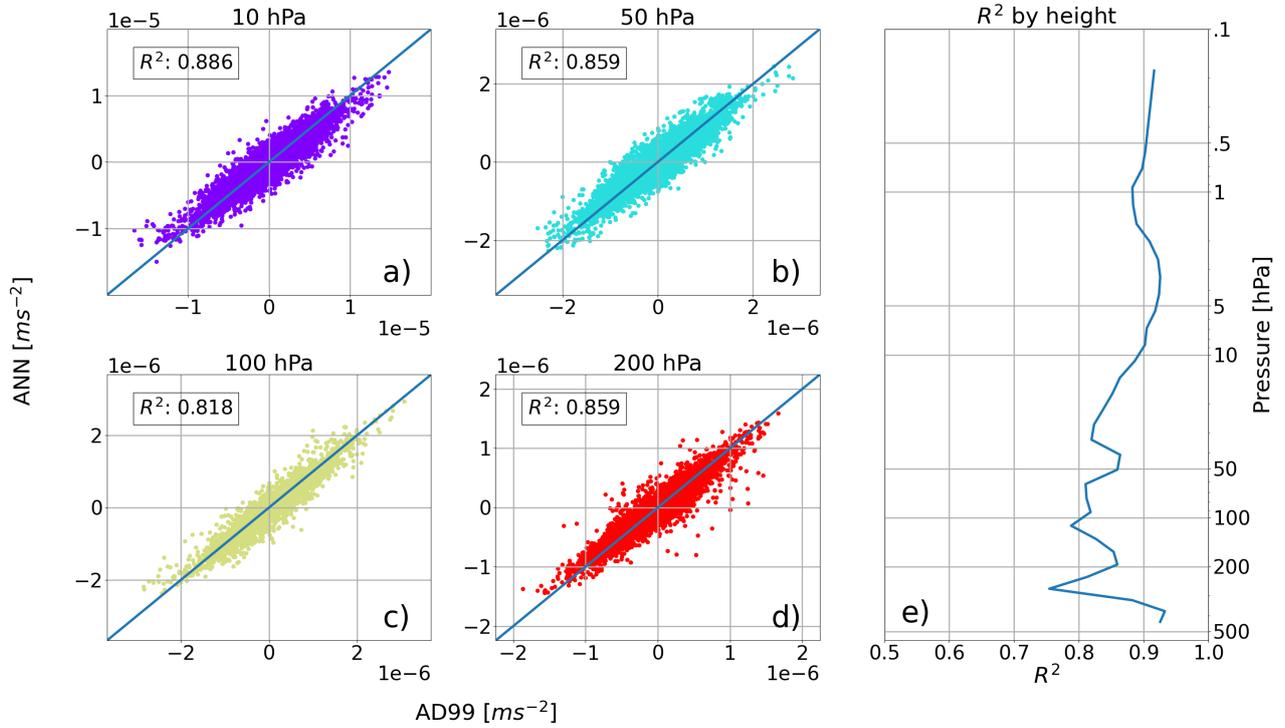


Figure S2. As in Figure 2, but for meridional GWD, panels a) through d) show ANN meridional GWD versus AD99 truth at 10 hPa, 50 hPa, 100 hPa and 200 hPa, respectively, for 10k samples in the test set. Panel e) shows pressure versus horizontally and temporally averaged R^2 values generated from one year of test data for zonal predictions.

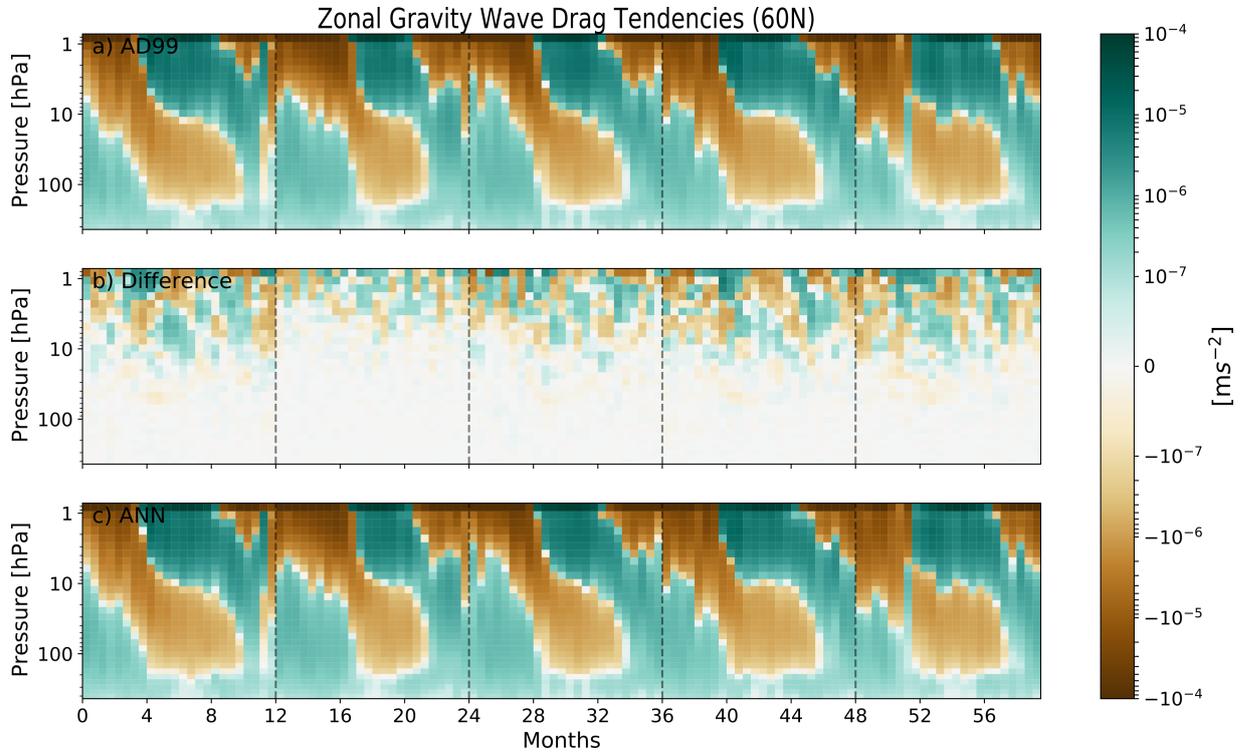


Figure S3. Pressure-Time profile of zonal mean drag 15-day tendencies at 60N latitude for zonal AD99 generated GW drag (a), the corresponding ANN predictions (c), and their difference (b). Vertical dashed lines separate years. The ANN is trained on months 12-24, and tested on all other months. The westward (brown) and eastward (green) bands correspond to drag associated with the seasonal cycle of the Polar Vortex.

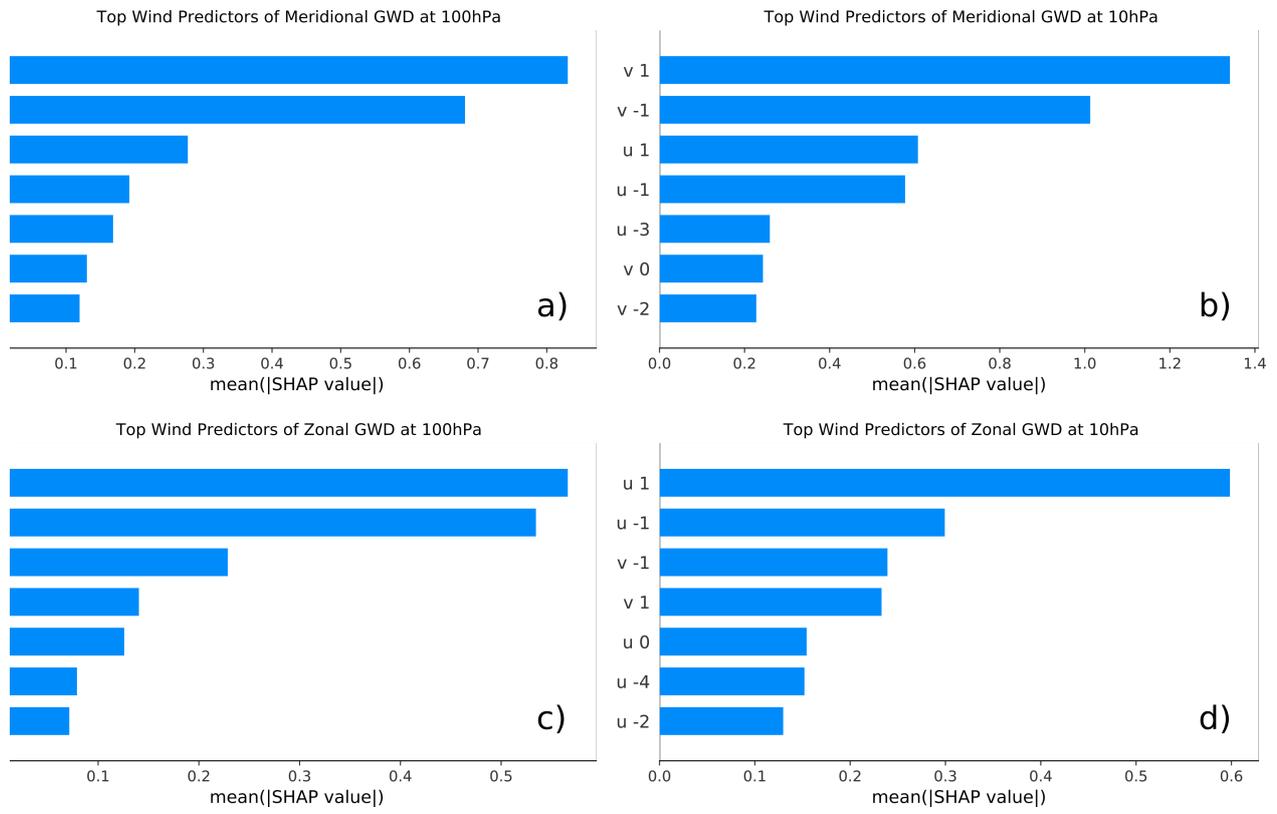


Figure S4. Panels a) through d) are SHAP bar plots of the ten most important meridional (v; panels a and b) and zonal (u; panels c and d) wind features used by WaveNet to generate meridional or zonal GWD at 10 hPa and 100 hPa. The vertical axes' values indicate displacements in vertical pressure levels, with positive and negative values being above and below the level of prediction, respectively (e.g., u -1 indicates zonal wind at the vertical level directly below the level of prediction). The results suggest that vertically local wind fields are the dominant features used by WaveNet to generate GWD.