

1 **Reconstructing balloon-observed gravity wave**
2 **momentum fluxes using machine learning and input**
3 **from ERA5**

4 **Sothea Has¹, Riwal Plougonven², Aurélie Fischer¹, Raj Rani³**
5 **Francois Lott³, Albert Hertzog⁴, Aurélien Podglajen³, Milena Corcos⁵**

6 ¹CNRS/Laboratoire de Probabilités, Statistique et Modélisation (LPSM), Université Paris Cité, France

7 ²Laboratoire de Météorologie Dynamique, Ecole Normale Supérieure, IPSL, Paris, France

8 ³Laboratoire de Météorologie Dynamique (LMD)/IPSL, PSL Research Institute, Paris, France

9 ⁴LMD/IPSL, Sorbonne Université, Paris, France

10 ⁵NorthWest Research Associates, Boulder, Colorado

11 **Key Points:**

- 12 • Eight superpressure balloons from the Strateole 2 mission provide observations
13 for accurate gravity wave momentum flux estimation
- 14 • Three machine learning methods are employed to probe the relationship between
15 the gravity wave momentum fluxes and ERA5's large-scale flow
- 16 • The most informative large-scale inputs are provided, along with a discussion of
17 the successes and challenges of machine learning methods

Corresponding author: Sothea Has, sothea.has@lpsm.paris

Abstract

Global atmospheric models rely on parameterizations to capture the effects of gravity waves (GWs) on middle atmosphere circulation. As they propagate upwards from the troposphere, the momentum fluxes associated with these waves represent a crucial yet insufficiently constrained component. The present study employs three tree-based ensemble machine learning (ML) techniques to probe the relationship between large-scale flow and small-scale GWs within the tropical lower stratosphere. The measurements collected by eight superpressure balloons from the Strateole 2 campaign, comprising a cumulative observation period of 680 days, provide valuable estimates of the gravity wave momentum fluxes (GWMFs). Multiple explanatory variables, including total precipitation, wind, and temperature, were interpolated from the ERA5 reanalysis at each balloon's location. The ML methods are trained on data from seven balloons and subsequently utilized to estimate reference GWMFs of the remaining balloon. We observed that parts of the GW signal are successfully reconstructed, with correlations typically around 0.54 and exceeding 0.70 for certain balloons. The models show significantly different performances from one balloon to another, whereas they show rather comparable performances for any given balloon. In other words, limitations from training data are a stronger constraint than the choice of the ML method. The most informative inputs generally include precipitation and winds near the balloons' level. However, different models highlight different informative variables, making physical interpretation uncertain. This study also discusses potential limitations, including the intermittent nature of GWMFs and data scarcity, providing insights into the challenges and opportunities for advancing our understanding of these atmospheric phenomena.

Plain Language Summary

Part of the atmosphere's large-scale circulation results from motions that are not resolved, or partly resolved, by weather or climate models. These include internal gravity waves, with horizontal scales from a few to hundreds of kilometers. The main sources occur in the troposphere, such as flow over mountains and cloud development. Their three-dimensional propagation induces major aggregated impacts in the stratosphere and mesosphere, forcing key aspects of the circulation. This forcing is accounted for in climate models by 'parameterizations', that mimics the effect of the unresolved waves based on the large-scale, resolved flow. These parameterizations necessarily retain crude approximations and introduce significant uncertainty in the models. For GWs, sources are a major uncertainty. This study makes use of the high-altitude balloon campaign Strateole 2 (Oct. 2019-Feb. 2020). Eight balloons circled Earth at heights around 18 to 20 km, providing unique observations of the GWs. These are used as targets for machine learning (ML) methods that take as inputs the information from outputs of a numerical weather prediction model describing the large-scale flow. The successes and difficulties of ML provide insights which can guide improvements of parameterizations, such as the most informative large-scale variables for estimating the unresolved waves.

1 Introduction

Climate models and Numerical Weather Prediction models resolve a widening range of atmospheric processes as computing power increases, enabling finer spatial resolution. Subgrid-scale processes persist nonetheless, and efforts to improve and constrain them better are essential. Internal gravity waves constitute one of these subgrid-scale processes, with important implications for the circulation and variability of the middle atmosphere (Fritts & Alexander, 2003). Motivations for improved modeling of the stratosphere includes climate (e.g. Solomon et al. (2010); Kremser et al. (2016)) but also predictability on shorter time scales (F. Vitart and A.W. Robertson, 2018; Butchart, 2022).

67 Gravity waves occur on scales ranging from a few to several hundreds of kilome-
68 ters. An important effect stems from their vertical propagation: gravity waves are re-
69 sponsible for vertical transfers of momentum from lower layers (troposphere: denser and
70 with more gravity wave sources) to upper layers (stratosphere and beyond), where they
71 constitute an essential driver of the overall circulation (Fritts & Alexander, 2003). A sig-
72 nificant part of the spectrum of gravity waves has been and remains unresolved in global
73 models, requiring these effects to be represented by parameterizations (Kim et al., 2003).
74 Models display sensitivity to these, calling for coordinated efforts to better constrain these
75 parameterizations from both observations and high-resolution modeling (Alexander et
76 al., 2010).

77 A global comparison of observed, resolved and parameterized gravity wave momen-
78 tum fluxes was carried out by Geller et al. (2013), highlighting significant discrepancies.
79 Although GWs parameterizations are now used routinely in climate models, their val-
80 idation against in situ observations remains a challenge. There exist global observations
81 derived from satellite observations (e.g. Ern et al. (2018)), but there are limitations on
82 the wavelengths that can be observed, and significant assumptions are needed to indi-
83 rectly deduce important quantities like the momentum fluxes from temperature fluctu-
84 ations, using polarization relations (Alexander et al., 2010; Ern et al., 2014). For these
85 reasons superpressure balloons have been highlighted as a valuable and accurate source
86 of information on gravity wave momentum fluxes (Geller et al., 2013). A downside of su-
87 perpressure balloon observations is their very sparse sampling of the lower stratosphere:
88 despite a broad coverage of the Southern Ocean (Jewtoukoff et al., 2015) and of the equa-
89 torial belt (Corcos et al., 2021), each balloon flight provides only local information: one
90 time series along its trajectory.

91 There are fundamental difficulties in validating parameterizations of gravity waves:
92 the purpose of a parameterization is to provide the forcing to the large-scale which is miss-
93 ing because of unresolved processes. Ideally, one would wish to *know* what this forcing
94 should be and validate this outcome of parameterizations. Unfortunately, this forcing
95 cannot be directly observed. Validating parameterizations by the realism of the clima-
96 tology and variability of the atmospheric circulation in global models constitutes a first
97 step, but is not a severe test and allows for compensating errors between parameterized
98 processes (Plougonven et al., 2020). More stringent tests involve comparisons to obser-
99 vations (de la Camara et al., 2014; Trinh et al., 2016). Recently, direct comparisons be-
100 tween observed and parameterized gravity waves have been carried out on the scale of
101 daily variations rather than at the level of general statistical characteristics (Lott et al.,
102 2023). The large-scale environment was described using the ERA5 reanalyses (Hersbach
103 et al., 2020), providing the background fields necessary to emulate the parameterization
104 of convectively generated waves of Lott & Guez (2013), which is the parameterization
105 used in the climate model of IPSL (Institut Pierre Simon Laplace, Boucher et al. (2020)).
106 The comparison was quite encouraging, with the gravity wave momentum fluxes hav-
107 ing the right order of magnitude, and an appropriate intermittency.

108 An essential aspect, and fundamental issue, to keep in mind when comparing ob-
109 served and modeled gravity wave momentum fluxes is their strong intermittency: in time
110 series of GWMF, one commonly finds short, intense peaks corresponding to a strong grav-
111 ity wave event, surrounded by considerably weaker values. This has been highlighted in
112 the long ‘tail’ of the Probability Density Function (PDF) of the GWMF (Alexander et
113 al., 2010; Hertzog et al., 2012), and quantified in simulations and observations (Plougonven
114 et al., 2013; Wright et al., 2013; Ern et al., 2022). This intermittency further contributes
115 to making the parameterization of gravity waves a challenging task.

116 For the improvement of parameterizations in general (not only those of gravity waves),
117 machine learning methods provide an array of possibilities. These have been explored
118 in different directions:

- 119 • Machine learning can enable the emulation of parameterizations, leading to sig-
120 nificant computational time savings (Chantry et al., 2021; de Burgh-Day & Leeuwen-
121 burg, 2023).
- 122 • Machine learning can help to capture the relationship between large-scale fields
123 and the unresolved process, as illustrated in the case of convection by Gentine et
124 al. (2018). For exploration, the dataset used as the truth came from a higher-resolution
125 simulation, not from observations; obtaining observationally based knowledge of
126 the effects to be parameterized remains a major challenge.
- 127 • Machine learning can be used to explore the relationship between the large-scale
128 flow and the resulting small-scale waves, as has been done for orographic waves
129 over Northern Japan (Matsuoka et al., 2020). Again, both the target and the in-
130 puts are modelled fields, but at different resolutions.
- 131 • As a precursor to a data-driven parameterization that would have learned from
132 observations, a machine learning-based emulator of a parameterization for grav-
133 ity waves has been used in a climate model, including under climate change con-
134 ditions (Espinosa et al., 2022).

135 The purpose and scope of the present study is to probe the relationship between
136 the large-scale flow and gravity waves in the Tropics, using machine learning approaches
137 to address fundamental issues: what fraction of the GWMF can be determined from knowl-
138 edge of the large-scale flow, and what fraction remains as *stochastic*? Which large-scale
139 variables are most informative, and do they match with our common understanding of
140 underlying gravity wave parameterizations? The present study belongs to the third cat-
141 egory outlined above for the uses of machine learning (the purpose is *not* to produce a
142 new parameterization, nor to emulate an existing one). With similar goals, Amiramjadi
143 et al. (2023) used machine learning methods to probe the relationship between the large-
144 scale flow and gravity waves, for non-orographic waves in the mid-latitudes and using
145 waves resolved in a reanalysis as a target. In contrast, the present study aims at *observed*
146 momentum fluxes in the Tropics, where the Strateole 2 campaigns provide a wealth of
147 new observations (Haase et al., 2018; Corcos et al., 2021).

148 The paper is organized as follows: Section 2.1 provides an overview of the data and
149 ML algorithms used in this study. Section 3 presents the performances of ML methods
150 in reconstructing the reference GMWFs. Section 4 discusses the factors that influence
151 the performances and addresses the limitations of ML methods. Finally, Section 5 con-
152 cludes the study with key takeaways and future directions.

153 2 Data and methodology

154 2.1 Data

155 We use in situ observations collected from eight constant-level balloon flights (al-
156 titude between 18.5 and 20km) during the Strateole-2 mission from November 2019 to
157 February 2020 (Corcos et al., 2021). As in Corcos et al. (2021), momentum fluxes (MFs)
158 were computed from raw balloon measurements following the procedure described in Vin-
159 cent and Hertzog (2014). Essentially, the pressure and horizontal wind time series are
160 first projected in the time-frequency domain thanks to a continuous wavelet transform
161 (Torrence and Compo, 1998). The pressure observations inform on the vertical displace-
162 ments of the balloon, which are related to those of air parcels, assuming that the bal-
163 loon behaves as a perfect isopycnic tracer. The time-frequency MF decomposition is then
164 derived from the wavelet cross-spectrum of the horizontal winds and air-parcel vertical
165 displacements. Segments polluted by non-geophysical artifacts (e.g. depressurization events)
166 are discarded.

167 For our analysis, and following Corcos et al. (2021), we considered gravity wave
168 MFs integrated over two frequency bands: a high-frequency (HF) band (i.e. short pe-

169 periods, ranging from 15 minutes to 1 hour) and wide-frequency (WF) band (i.e., long pe-
 170 riods, ranging from 15 minutes to 1 day). For the sake of readability, in all that follows
 171 we focus on the HF band, unless explicitly stated. Additionally, we also differentiate be-
 172 tween eastward-propagating waves that yield positive MF in the zonal direction (east-
 173 ward) and westward-propagating waves that produce negative MF (westward). We use
 174 these MFs as a reference for the true target MFs. Then, we pair them with large-scale
 175 flow input information from ERA5, such as wind velocity (u and v), temperature (temp),
 176 total precipitation (tp) and logarithm of surface pressure (lnsp). These fields are retrieved
 177 for each balloon, from fields at a resolution of $1^\circ \times 1^\circ$, at the grid point closest to the
 178 balloon position. Additionally, the same input variables have been retrieved in the vicini-
 179 ty of a 5 by 5 horizontal square centered on the grid point closest to the balloon; in the
 180 present study, only total precipitation in this extended area around the balloon will be
 181 used. In the vertical, the ECMWF model comprises a total of 137 levels. Four levels are
 182 retained in the present study, to succinctly describe the vertical wind profile from the
 183 surface to balloon flight level (see Table 1).
 184 The inputs and the targets are interpolated and averaged into 1-hour time resolution.
 185 The three ML models are trained using 3-hour time averaging data, and their perfor-
 186 mance will be evaluated based on daily averaging time resolution, as presented in Lott
 187 et al. (2023). Table 1 presents the finalized large-scale flow variables utilized for train-
 188 ing ML models.

189 2.2 Methodology

190 In this study, three tree-based ensemble ML methods are considered: random for-
 191 est (RF) introduced in Breiman (2001), extremely randomized trees also known as extra-
 192 trees (ET) by Geurts et al. (2006), and Adaptive Boosting or Adaboost regressors by
 193 Freund & Schapire (1997). These algorithms construct multiple decision trees, and the
 194 final prediction is determined by aggregating the individual decision tree predictions.

195 It should be noted that other methods, such as deep neural networks, as well as
 196 other types of networks including convolutional and recurrent neural networks, have also
 197 been implemented. However, the performances of these methods are not comparable to
 198 the presented tree-based algorithms, as these models typically require a large number
 199 of observations to achieve comparable results. The limitations and concerns regarding
 200 the models, the large-scale input variables, the target observations, and the nature of the
 201 relation between the large-scale and small-scale flow will be discussed later in Section 4.3.

202 2.2.1 Decision tree

The decision tree algorithm (Breiman et al., 1984) is the foundational building block
 of the primary ML methods used for our predictions. They are widely used for nonlin-
 ear prediction problems due to their efficiency and interpretability. To construct a de-
 cision tree, the training data is recursively partitioned into small hyperrectangular re-
 gions of the forms $R_1 = \{X \leq \alpha\}$ and $R_2 = \{X > \alpha\}$ for some ERA5 input variable
 X (wind velocity or precipitation, for instance) and threshold α . At each step, we re-
 cursively split the input space into hyperrectangular regions that are as pure as possi-
 ble. Purity refers to the homogeneity of the training target y (GWMF) within each re-
 gion, and Total Within Sum of Squares (TWSS) is utilized as the impurity measure in
 this study. Specifically, a split is performed at any input variable X at threshold α if it
 minimizes the following TWSS criterion:

$$\sum_{y \text{ of } R_1} (y - \mu_1)^2 + \sum_{y \text{ of } R_2} (y - \mu_2)^2,$$

203 where

- 204 • R_1 and R_2 are the left and right regions of the split

205 • μ_1 and μ_2 are the average targets within region R_1 and R_2 respectively.

206 Any new observation must belong to one of these regions, and its prediction is determined
 207 by averaging the target values of all the neighboring observations within that block. Con-
 208 structing an optimal tree is generally challenging, and the tree’s structure, such as its
 209 depth and the minimum size of regions allowed to split, are hyperparameters that need
 210 to be optimized. Figure 1 below provides an example of a simple decision tree trained
 211 on 100 observations of precipitation and zonal wind velocity to predict absolute GWMF.

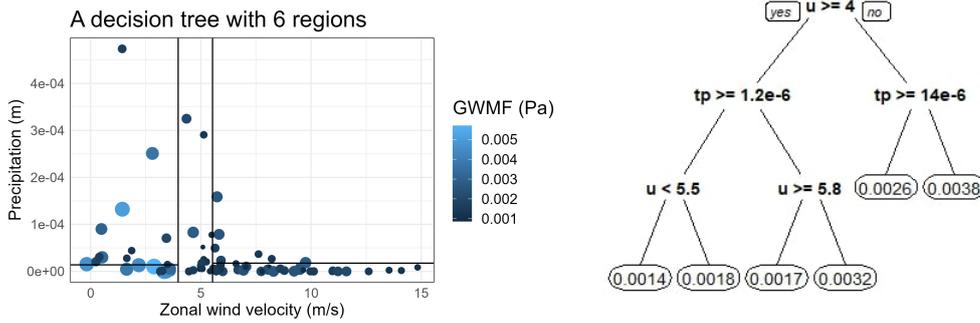


Figure 1. An example of a simple decision tree built using precipitation and wind velocity to predict absolute GWMF. The left side is the partition cell representation of the tree on the right side. The data points are colored according to the value of their target GWMF.

212

213 **2.2.2 Random forest**

214 Random forest (RF) is a powerful ensemble learning method that aims at minimiz-
 215 ing variance across a collection of decision trees by averaging their predictions (Breiman,
 216 2001). The term ‘random’ signifies the deliberate characteristic of constructing individ-
 217 ual trees using different bootstrap samples (sampling observations with replacement) and
 218 exploring only a small, randomly selected, subset of the complete input features. This
 219 approach effectively decorrelates the individual trees, resulting in a reduction of predic-
 220 tion variance. Additionally, the construction of each individual tree using only a small
 221 subset of input features enables random forest to handle high-dimensional data effectively.
 222 The key hyperparameters in a random forest are the number of trees, tree complexity,
 223 and the number of randomly selected features used in building the individual trees. Fine-
 224 tuning these hyperparameters is essential to optimize the performance of the method.

225 **2.2.3 Extremely randomized trees**

226 Extremely randomized trees or Extra-trees (ET) operates similarly to RF approach,
 227 with the distinction that each tree is constructed using the complete training data, and
 228 each split is performed at *random values* using a *random subset* of input features (Geurts
 229 et al., 2006). This results in a high degree of independence among the trees and can oc-
 230 casionally yield remarkable results compared to the random forest method.

231 **2.2.4 Adaptive boosting**

232 Adaptive boosting (Adaboost) combines weak learners to create a strong predic-
 233 tive model (Freund & Schapire, 1997). Weak learners refer to predictive models that per-

234 form slightly better than random guesses, and simple decision trees with only a few splits
 235 (stumps) are used as weak learners in this study. During each iteration, Adaboost com-
 236 bines an individual stump by using a weighted sum, where the weight assigned to the
 237 current stump is determined based on its overall performance in predicting the target
 238 variable. Additionally, the weights associated with the individual training data points
 239 are adjusted manually based on their prediction accuracy, giving more attention or weight
 240 to points with poor predictions in the next iteration. Adaboost is well known for its abil-
 241 ity to mitigate overfitting (Rätsch et al., 2001) and has achieved significant success in
 242 various prediction challenges (see, for example, Benjamin Bossan (2015) and ZEWEICHU
 243 (2019)).

244 **2.2.5 *K*-fold cross validation**

245 *K*-fold cross-validation is the most commonly used model selection technique in ma-
 246 chine learning. It involves dividing the training data into K parts or folds, namely F_1, \dots, F_K ,
 247 then a model is trained on $K-1$ folds, and it is tested on the remaining one. This pro-
 248 cess is repeated K times and the final performance is the average performance over all
 249 the K different testing folds. In this study, *K*-fold cross-validation is used to prevent over-
 250 fitting and to select the best possible hyperparameters of each ensemble method. More
 251 precisely, if f_θ is the considered method (random forest, for example) with a hyperpa-
 252 rameter $\theta \in \Theta$, then the optimal hyperparameter θ^* is defined by,

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{K} \sum_{k=1}^K \sum_{(x_i, y_i) \in F_k} (f_\theta(x_i) - y_i)^2. \quad (1)$$

253 In our study, θ consists of the depth of the decision trees (maximum number of splits
 254 performed from the root node to the leaves), the size of random subsets of the ERA5
 255 input features to be considered when building individual trees, and the number of deci-
 256 sion trees used in each ensemble learning method. All these keys are tuned using 10-fold
 257 cross-validation.

258 **2.3 Training**

We first train ML models with an extensive set of ERA5 inputs. Subsequently, we
 refine these inputs to a more manageable subset (see Table 1 below) using importance
 feature scores, which will be described in Section 3. Moreover, in order to reduce the in-
 fluence of extreme values in the target and increase its normality, the Box-Cox trans-
 formation (Box & Cox, 1964) is performed on the GWMF y to obtain the transformed
 target \tilde{y} :

$$\tilde{y} = \frac{y^\lambda - 1}{\lambda}.$$

In the experiment, the exponent $\lambda = 0.6$ is chosen based on the performance of mod-
 els trained on the corresponding transformed target data. The predictions given by ML
 models are then reverted using the inverse transformation:

$$y = (1 + \lambda\tilde{y})^{1/\lambda}.$$

259 Moreover, to predict any GWMFs (absolute, eastward, or westward GWMFs of HF
 260 or WF case) of any given balloon, the ML models are trained using data from the seven
 261 other balloons. The models are fine-tuned using a 10-fold cross-validation method to op-
 262 timize their performances.

263 Finally, the resolutions used for the data (see Section 2.1) reflect the phenomena
 264 we aim to estimate. From large-scale information as described from reanalyses at a res-
 265 olution of $1^\circ \times 1^\circ$ and hourly in time, it is only reasonable to estimate GWMFs aver-
 266 aged over a comparable timescale (one hour). As the balloons drift at velocities typically

267 ten to twenty m.s^{-1} , this corresponds to sampling over a spatial area of several tens of
 268 kilometers. The final choice for the specific setting used has been also guided by the mo-
 269 tivation to make comparison with the results of Lott et al. (2023) possible.

270 The targeted gravity waves, as observed by the balloons, cover the whole range of
 271 intrinsic frequencies. The high frequency band (HF, see Section 2.1) may a priori be more
 272 difficult to predict from ML because it is expected to be more intermittent (Corcos et
 273 al., 2021), so that sampling will be a more severe issue than for the WF band. On the
 274 other hand, higher frequency waves propagate more vertically and are shorter-lived, both
 275 factors contributing to a stronger causal relationship between local conditions below the
 276 balloons and observed gravity at balloon level. As it has turned out that this second fac-
 277 tor is more important, we focus hereafter on HF waves as the target, while the WF cases
 278 are detailed in the supplementary document.

Name	Notation	Description
Zonal, meridional wind velocity (m.s^{-1}) & temperature (K)	u_j, v_j & temp_j	with vertical level $j \in \{0, 2, 9, 19\}$ (km), where 0 is the surface and 19 is the balloon's level.
Total precipitation (m)	tp	at center of horizontal grid points.
Mean & standard deviation of precipitation (m)	tp_{mean} & tp_{sd}	over horizontal grid points.
Solar zenith angle ($^\circ$)	sza^1	at the location of the balloon.
Log surface pressure ($\log(\text{hPa})$)	lnsp	at the surface level.

Table 1. Large-scale input data for training ML models.

279 2.4 Evaluation metric

280 An important aspect in any comparison of models to observations is the choice of
 281 a metric to evaluate the performance of the models. We explain here why, in line with
 282 Lott et al. (2023), we use correlation between modelled and observed values as our met-
 283 ric. The current study is in line with studies that have compared parameterized and ob-
 284 served gravity waves (eg Geller et al. (2013)). In such comparisons, the first aim is nat-
 285 urally to compare *mean* momentum fluxes, yet over the past decade the importance of
 286 having a realistic variability has been emphasized (Alexander et al., 2010)). This has high-
 287 lighted the notion of intermittency (Hertzog et al., 2012) and quantification of the dis-
 288 tribution of momentum fluxes when comparing parameterizations to observations (de la
 289 Camara et al., 2014; Bushell et al., 2015). These comparisons, however, concern the over-
 290 all statistics, not a direct comparison of observed and parameterized variations on a case-
 291 to-case basis. Obtaining an appropriate observational dataset and gathering the corre-
 292 sponding large-scale variables for such a case-to-case comparison has required significant
 293 work and has been achieved for the comparison of Lott et al. (2023). These datasets pro-
 294 vide a unique opportunity to investigate the co-variability of observed GWMF and es-
 295 timations from the large-scale flow, whether based on parameterizations (Lott et al., 2023)

¹ Solar zenith angle is the only input obtained from the balloons, not from the ERA5. It is a periodic function that provides an estimation of time of the day and the balloon's location.

296 or on machine learning techniques presented in this study. This is why we here focus on
 297 this co-variability, quantified by the correlation. It is expected that the averaging effect
 298 of tree-based algorithms may lead to underestimation of the target, especially when deal-
 299 ing with rare extreme values such as GWMFs. Obtaining appropriate intermittency of
 300 the reconstructed gravity wave momentum fluxes will require further efforts, and direc-
 301 tions for these efforts are discussed in the perspectives (Section 5).

302 **3 Results**

303 This section reports the correlations of ML methods in reconstructing various types
 304 of observed GWMFs. The numerical study is carried out using `sklearn.ensemble` mod-
 305 ule in Python (Pedregosa et al., 2011). In general, the three ML models exhibit very com-
 306 parable performances on any given balloon. In contrast, the performance of the ML mod-
 307 els varies significantly from one balloon to another. At their best, ML models can achieve
 308 an encouraging level of correlation larger than 0.7. The average performance over all bal-
 309 loons and data exceeds 0.5. The worst performances is found for westward GWMF for
 310 a specific balloon, with correlation down to 0.2. Overall, the performances of ML mod-
 311 els are sensitive to the choice of balloons and the types of GWs being considered (east-
 312 ward, westward or absolute GWMFs). The numerical results for HF waves are presented
 313 in the following subsections, while the WF cases are presented in supplementary docu-
 314 ment.

315 **3.1 Overall performances**

316 Three examples of observed and predicted GWMFs of the HF case are presented
 317 in Figure 2 below. Each subplot displays the eastward component of the GWMFs in the
 318 positive part and the westward ones in the negative part. It can be observed that the
 319 models effectively capture the fluctuations of the observed momentum fluxes, particu-
 320 larly on balloon 2. However, the models struggle to fully estimate the amplitudes of high-
 321 peak events, especially for balloons 3 and 7. Overall, the performances of all ML mod-
 322 els are quite similar; however, there are cases where one outperforms the others. For ex-
 323 ample, Adaboost appears to do a slightly better job on balloon 2 than the other two mod-
 324 els in capturing the amplitudes of the high-peak events. It is worth noting that balloon
 325 2 presents overall the best performance for the ML models, balloon 7 illustrates a typ-
 326 ical average case, and balloon 3 is the most challenging one to predict: this is suggested
 327 visually in Figure 2, and is confirmed quantitatively in Table 2.

328 A feature of the reconstructed GWMF is that the peak values are generally under-
 329 estimated, as can be seen even for balloon 2 in Figure 2. This is partly expected given
 330 that tree-based models involve averaging from numerous decision trees, some of which
 331 are insufficiently informed to capture extreme occurrences of GWMFs. To document the
 332 relationship between the reconstructed and observed GWMFs, scatterplots are displayed
 333 in Figure 3. These illustrate how the reconstruction captures well the variations of GWMFs,
 334 especially for rather weak variations. In contrast, for occurrences of larger MFs, the ob-
 335 served values cover a range of values that are not captured by the ML approaches. The
 336 scatterplots illustrate that those occurrences are rare, and the training data certainly con-
 337 stitutes a limiting factor. It is not clear that it may be possible to capture, in a deter-
 338 ministic way, these extremes. It is worth noting that the ML approaches do generally
 339 capture when the GWMF is at the high end of the range of reconstructed values.

340 Figure 4 presents boxplots of Pearson’s correlation coefficients between predicted
 341 and true GWMFs of the HF case. Firstly, choosing the best model is challenging due
 342 to the variability in the boxplot positions, which depends on the choices of balloons and
 343 GWMF types. For instance, on balloon 2, the correlation boxplot of Adaboost is higher
 344 than the other two methods for the absolute and westward cases but lower than Ran-
 345 dom Forest for the eastward case. However, these differences are generally insignificant

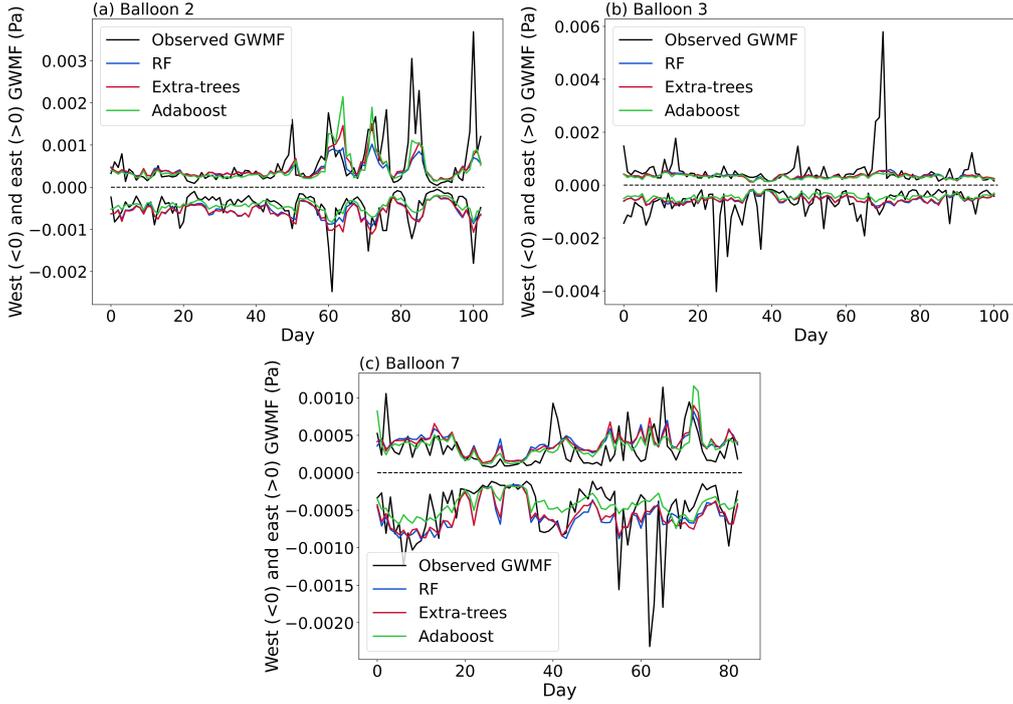


Figure 2. Observed and predicted time series of high-frequency east and westward GWMFs of the best, worst and medium cases: balloon 2, 3 and 7, respectively. The x-axis label "Day" indicates the number of days since the individual balloon was launched, with 0 corresponding to the moment of launch.

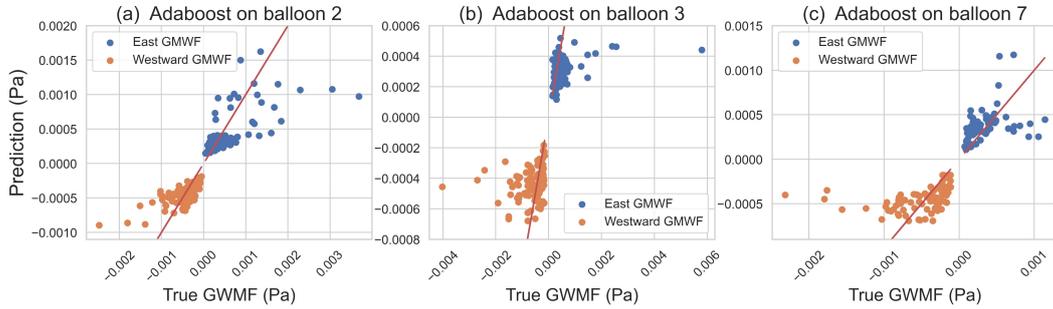


Figure 3. Scatterplots of predictions against observed (true) GWMF corresponding to the time series of Figure 2. Only the predictions of Adaboost are presented for balloon 2, 3 and 7 (from left to right). The lower groups represent the westward fluxes, while the upper groups denote the eastward ones. The red line serves as the reference 1:1 line.

346 compared to the variations observed between different balloons. Secondly, ML models
 347 demonstrate strong performance on balloons 2, 6, and 8 across all types of momentum
 348 fluxes, and they also excel in predicting the eastward momentum flux of balloon 1. Nev-
 349 ertheless, balloons 3, 4, 5, and 7 pose greater challenges, with the most difficult being
 350 the westward component of GWMF on balloon 3. Finally, the ML models generally out-
 351 perform the gravity wave drag scheme of the IPSL model (Lott et al., 2023), except for
 352 balloon 3 (east and westward) and balloon 4. Moreover, Table 2 provides the statisti-
 353 cal significance of the correlations presented in Figure 4.

Flight	Alt	Start	End	Duration/ DOF	Absolute			Eastward			Westward		
					RF	ET	AB	RF	ET	AB	RF	ET	AB
01_STR1	20.7	12/11/19	28/02/20	107/53	0.56	0.57	<u>0.58</u>	0.67	<u>0.69</u>	0.67	0.38	0.37	<u>0.43</u>
02_STR2	20.2	11/11/19	23/02/20	103/51	0.70	0.67	<u>0.74</u>	0.67	0.62	0.65	0.60	0.63	<u>0.70</u>
03_TTL3	19.0	18/11/19	28/02/20	101/33	0.45	0.48	<u>0.49</u>	0.41	<u>0.49</u>	0.43	0.21	<u>0.23</u>	0.18
04_TTL1	18.8	27/11/19	02/02/20	67/22	0.44	0.43	<u>0.47</u>	0.47	<u>0.48</u>	0.44	0.35	0.33	<u>0.37</u>
05_TTL2	18.9	05/12/19	23/02/20	79/19	0.51	<u>0.56</u>	0.55	0.39	<u>0.48</u>	0.35	0.35	<u>0.40</u>	0.50
06_STR1	20.5	06/12/19	01/02/20	57/10	0.72	0.74	<u>0.75</u>	0.64	0.65	0.70	0.68	<u>0.72</u>	0.57
07_STR2	20.2	06/12/19	28/02/20	83/16	0.51	<u>0.53</u>	0.48	0.46	<u>0.49</u>	0.42	0.44	<u>0.45</u>	0.32
08_STR2	20.2	07/12/19	22/02/20	77/12	0.74	<u>0.76</u>	0.72	0.71	<u>0.71</u>	0.68	0.66	<u>0.66</u>	0.64

Table 2. Average correlation coefficients between predicted and observed high-frequency GWMFs in 24h time resolution. In each case, by using decorrelated time as the degree of freedom (DOF), t-test statistics can provide the significance of each correlation with the convention: *italic boldface* = 99%, **boldface** = 95%, *italic* = 90%, and normal font = below 90% significant. For any given type of GWMF, the underlined correlations indicate the best performance of ML method on that target.

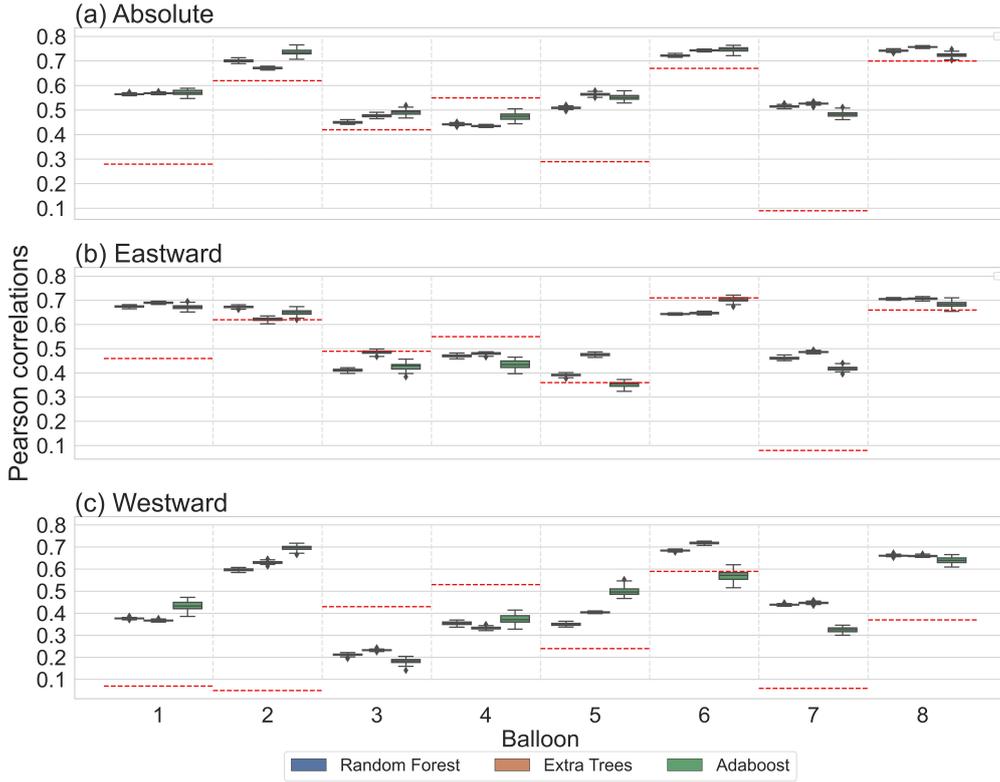


Figure 4. The boxplots display the correlations between predicted and observed high-frequency GWMFs obtained from 50 runs of ML methods as shown in Table 2. For each balloon, moving from left to right, the three boxplots correspond to the Random Forest, Extra Trees, and AdaBoost methods, respectively. The dashed horizontal red lines indicate the performance of the parameterization of the IPSL model (Lott et al., 2023).

354

3.2 Which large-scale inputs are informative for ML models?

355

356

357

358

359

360

361

The tree-based ensemble ML models employed in this study are not only proficient in predicting GWMFs but also offer valuable insights into the importance of large-scale input information during their training process. Each method exploits the feature importance (decrement of impurity measure at each split) of its individual decision trees for determining the overall feature importance, resulting in a ranking of input features from most to least important. Figure 5 showcases the ranking of the top 5 input features for all ML methods and GWMF types of the HF case.

362

363

364

365

366

367

368

369

370

371

372

Generally, the high-ranking inputs consist of variables that describe precipitation and wind velocity at and below the balloon’s level. It is important to note that different models may not rank input features in the same way for a given target (as seen along the rows), due to the variations in the way individual trees are grown. However, the three models concur on the strongly impactful input features; for example, wind velocity at the balloon’s level (u_{19}) ranked first in the eastward case (second row) for all models. This suggests that the wind velocity surrounding the balloons is the most informative large-scale variable for predicting eastward gravity wave momentum fluxes (GWMFs). Furthermore, the few most significant inputs show a similar preference in both absolute and eastward GWMFs within the same model, as demonstrated in the columns of the first and second rows. For instance, standard deviation and average total precipitation

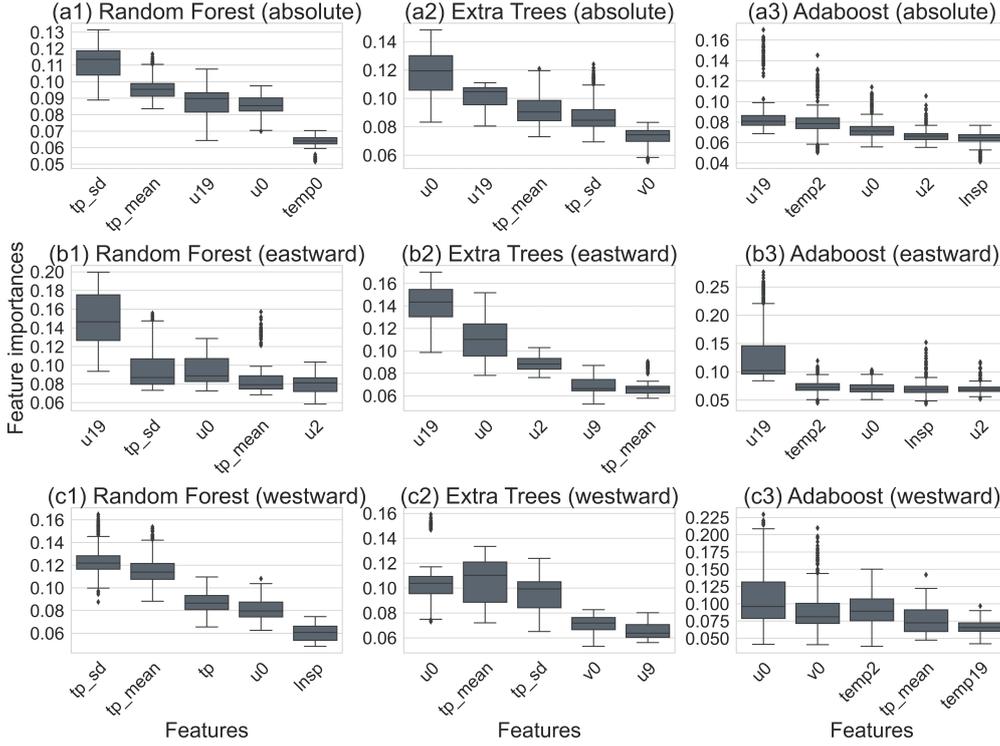


Figure 5. The boxplots show the 5 most important features given by different ML models (by column) on different types of targets (by row). Each boxplot is obtained from the same 50 simulations as displayed in Figure 4.

373 (tp_sd and tp_mean) are identified as impactful inputs in random forests, while surface
 374 zonal wind velocity (u0) is deemed the most important one in extra trees.

375 4 Discussion

376 While the results of the machine-learning models are generally encouraging, defi-
 377 ciencies and cases with poor performances were also found. The main motivation for this
 378 study being to probe the relationship between the large-scale and the unresolved process,
 379 these somewhat negative results are also of interest and can provide useful insights.
 380 Possible explanations for the main difficulties encountered are discussed below.

381 4.1 Why are westward GWMFs more challenging?

382 Figure 4 displays the performances of the ML models and those of the parameteriza-
 383 tion used in the IPSL climate model. Balloon 4 constitutes an exception, for which
 384 the parameterization systematically performs better than the ML methods. Leaving bal-
 385 loon 4 aside, ML approaches unambiguously outperform the parameterization for the ab-
 386 solute momentum fluxes. For the eastward momentum fluxes, ML approaches generally
 387 perform better or are similar to the parameterization. In contrast, both ML approaches
 388 and the parameterization have poorer performances for westward MF, and with greater
 389 variability for both: for five balloons, ML outperforms clearly the parameterization, whereas
 390 for two balloons (including balloon 4) the parameterization clearly outperforms the ML.
 391 The present section discusses possible reasons for this difficulty in reproducing the west-
 392 ward momentum fluxes.

393 Figure 6 displays the Probability Density Function of winds for three balloons as
 394 blue curves: balloon 2 has flown in winds that include a majority of westward, strong
 395 winds. Like balloon 1, it traveled near 10°S in easterly flow for a significant portion of
 396 its flight. In contrast, balloons 3 and 7 have flown in weaker winds, with a mild dom-
 397 inance of westerly winds. Also plotted in Figure 6 are conditional PDFs of the zonal winds,
 398 conditioned on the intensity of the absolute GWMF. The purpose is to detect if strong
 399 values of GWMF were associated to specific wind conditions. For balloon 2, strong GWMF
 400 values were found mostly for moderate to strong easterly winds, and this distribution
 401 is insensitive to the quantile chosen for the GWMF (90th, 95th or 99th percentile). For
 402 balloon 7, the distribution is somewhat sensitive to the quantile chosen. Finally, for bal-
 403 loon 3, the conditional distribution of zonal wind dramatically changes when it is restricted
 404 to the 99th percentile. This detects a particularly intermittent time series, with variabil-
 405 ity dominated by one extreme event, as seen from Figure 2. These findings contribute
 406 to explaining the poor performances for balloon 3: the variability of GWMF was dom-
 407 inated there by one (or very few) extreme events, occurring in a specific condition with
 408 very weak winds (close to zero, less than 5 m.s⁻¹). In contrast, the good performances
 409 for balloon 2 occur in a case with less intermittency, for which large GWMF are found
 410 in strong (easterly) winds.

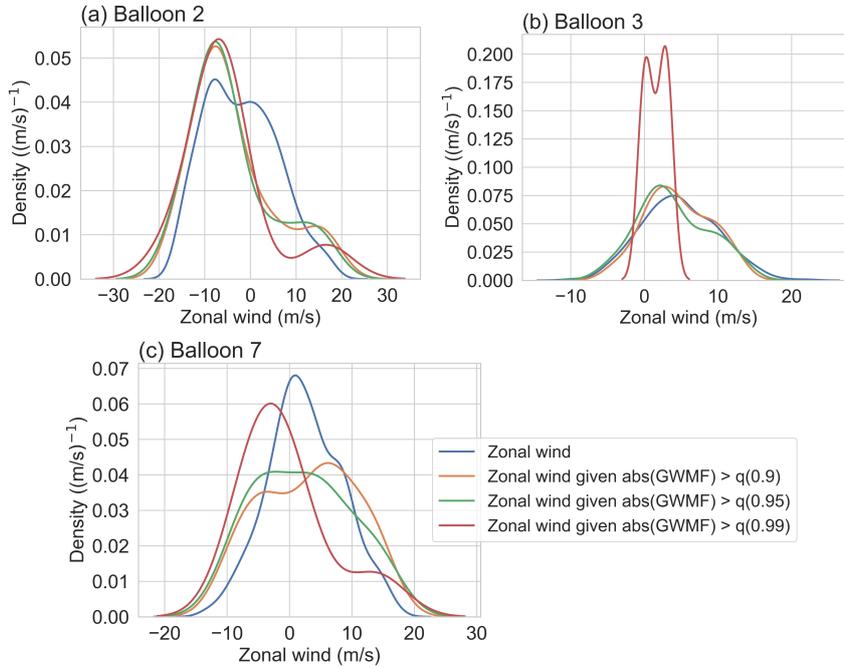


Figure 6. Conditional densities of zonal wind given different values of high-frequency westward GWMFs. Here, $q(0.9)$, $q(0.95)$ and $q(0.99)$ are the 90%, 95% and 99% quantiles of the absolute value of high-frequency westward GWMFs, respectively.

411 From Table 2, Figure 4 and the trajectories of the balloons (Corcos et al., 2021),
 412 it appears that drifting with easterly winds may constitute a favorable factor (balloon
 413 2), but neither a sufficient one (the correlation for westward momentum fluxes for bal-
 414 loon 1, which has a similar trajectory, is moderate, 0.43 at most) nor a necessary one:
 415 balloons 6 and 8 generally drift eastward, but good performances are found for the ML
 416 reconstruction of the westward MF (0.66 and 0.72 respectively).

417 Another aspect that influences the performances is the geographical location, and
 418 more specifically the latitude of the balloons. Figure 7 displays the PDF of latitude for
 419 the eight balloons, distinguishing those for which the ML reconstruction of westward MF
 420 is satisfactory (balloons 1, 2, 6 and 8, full lines) from those for which it remains chal-
 421 lenging (balloons 3, 4, 5 and 7). Here again, one does not isolate a necessary condition,
 422 but the balloons for which reconstruction remain challenging are those that remain clos-
 423 est to the equator. This is consistent with the general expectation that dynamics is more
 424 complicated near the Equator, although it is not completely clear why this should mat-
 425 ter for a small-scale process such as convectively generated gravity waves. It may be that
 426 it is not the dynamics itself that is intrinsically more difficult to capture at the Equa-
 427 tor: it may be the input variables that are poorer, less accurate, very close to the Equa-
 428 tor. It is known indeed that significant errors, in particular in the wind, are present in
 429 the reanalyses very near the Equator (Podglajen et al., 2014; Baker et al., 2014; Ern et
 430 al., 2023) and the errors are enhanced within a few degrees of the Equator (roughly be-
 431 tween 8°S and 8°N).

432 **4.2 Why are some balloons easier to predict than others?**

433 Figure 7 indicates that the predictability of the observed GWMFs is influenced by
 434 the balloons’ position, specifically, their distances from the equator. Balloons that trav-
 435 eled farther from the equator, primarily south (except for balloon 6, which also explored
 436 farther to the north), were found to be easier to predict. This tendency is observed for
 437 balloons 1, 2, 6, and 8 which are the well-predicted balloons. In contrast, the challeng-
 438 ing balloons spent most of their time flying within a few degrees of the equator, where
 the atmospheric conditions are not well described by ERA5 data.

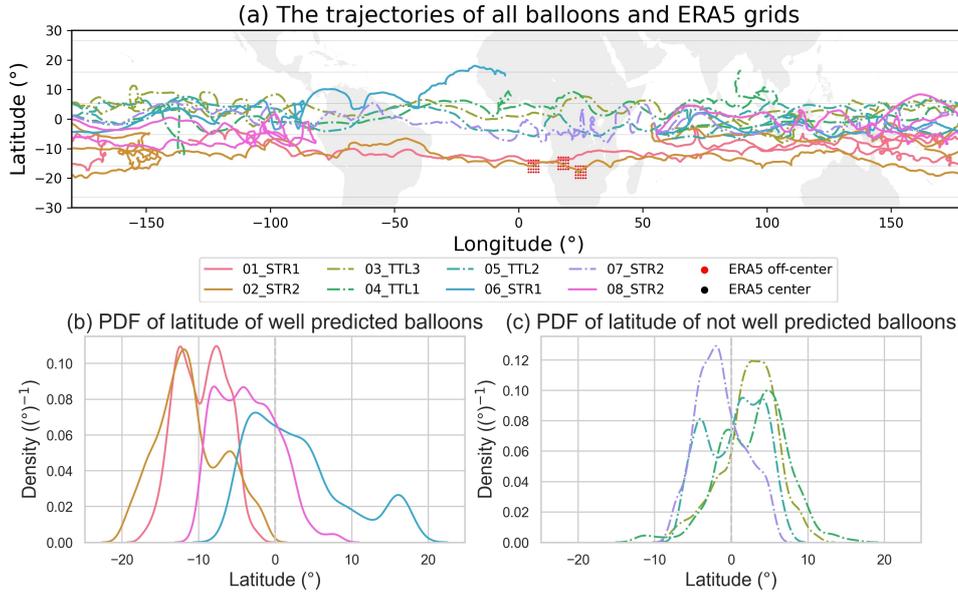


Figure 7. The trajectories of the balloons during the whole flight (a), and their latitude PDFs (b) and (c). Dashed lines correspond to balloons that pose challenges in prediction.

439

440 **4.3 Exploring potential reasons for unsatisfactory cases**

441 Several factors are expected to limit the ability to estimate the observed GWMFs
 442 from inputs describing the large-scale flow:

- 443 A. Part of the relationship between the large-scale flow and a subgrid-scale process
 444 such as gravity waves is non-deterministic, or stochastic: for given values of the
 445 large-scale fields, a range of different realizations of the subgrid-scale process is
 446 possible. It depends on the process: orographic gravity waves are likely more pre-
 447 dictable than convective processes for instance.
- 448 B1. The estimate of GWMFs from superpressure balloons is very local and samples
 449 only along its trajectory. This is only partly mitigated by the hourly averaging.
 450 The GWMFs time series certainly remain sensitive to the specific location of the
 451 balloon. At present, it is difficult to estimate this sensitivity. Investigations with
 452 virtual balloons in high-resolution simulations shall be informative on this issue.
- 453 B2. A second concern regarding the target used for the ML is the observational error
 454 present in the estimates of the GWMFs from balloon measurements. These es-
 455 timates are regarded as accurate because several variables are measured simulta-
 456 neously and because of the quasi-Lagrangian nature of the measurements (Geller
 457 et al., 2013; Vincent & Hertzog, 2014). There remains nonetheless observational
 458 error.
- 459 C1. Concerns are also present for the input variables, and in particular it is known that
 460 the description of the equatorial dynamics is challenging, with significant errors
 461 remaining present in the reanalysis especially for wind (Podglajen et al., 2014).
- 462 C2. Another concern regarding input variables is that we may have omitted variables
 463 that could have been informative.

464 In our study, we mitigated the concern of omitting informative variables (C2.) by
 465 initially training ML models on a large set of ERA5 inputs, then selectively reducing them
 466 to a reasonably small subset, as described in Section 2.1. This approach ensures that es-
 467 sential ERA5 inputs are not inadvertently omitted. Furthermore, fine-tuning the hyper-
 468 parameters of the models enhances their predictive capacity. Regarding the concern of
 469 large-scale variables (C1.), a sensitivity test to the error of ERA5’s wind is described at
 470 the end of Section 5 (key messages).

471 In addition, we observe that all the balloons often flew over many convective pro-
 472 cesses, and the high-peak events often correspond to deep convective systems, as illus-
 473 trated for selected cases in Figure 8 below. On January 12th, 2020, balloon 2 was fly-
 474 ing in an area of convection (upper panels (a1) and (a2)), which is likely responsible for
 475 the highest peaks in its GWMF time series. Interestingly, for balloon 2, almost all events
 476 correspond very well with precipitation as described by ERA5 (first column of Figure 9).
 477 On the contrary, there is only one big event that happened for balloon 3 around Jan-
 478 uary 29th, 2020 (lower left panel (b)). However, the ML models failed to capture it, as
 479 it appears to be absent from the ERA5 input variables (not reflected in precipitation nor
 480 winds as shown in the second column of Figure 9). This is also true for other challeng-
 481 ing balloons, such as the 4th and 5th. Regarding balloon 7, the large-scale flow provide
 482 partial information for the high-peak events, resulting in partial success in the model’s
 483 predictions.

484 5 Conclusion and perspectives

485 5.1 Key messages

486 The relationship between the large-scale atmospheric flow and gravity waves in the
 487 lower stratosphere has been investigated using Machine Learning (ML) approaches. This
 488 relationship is accounted for in global models through *parameterizations*. ML approaches
 489 allow us to revisit these in several ways, notably investigating how much of the subgrid-
 490 scale signal may be estimated *deterministically*, and which are the key variables for that
 491 purpose.

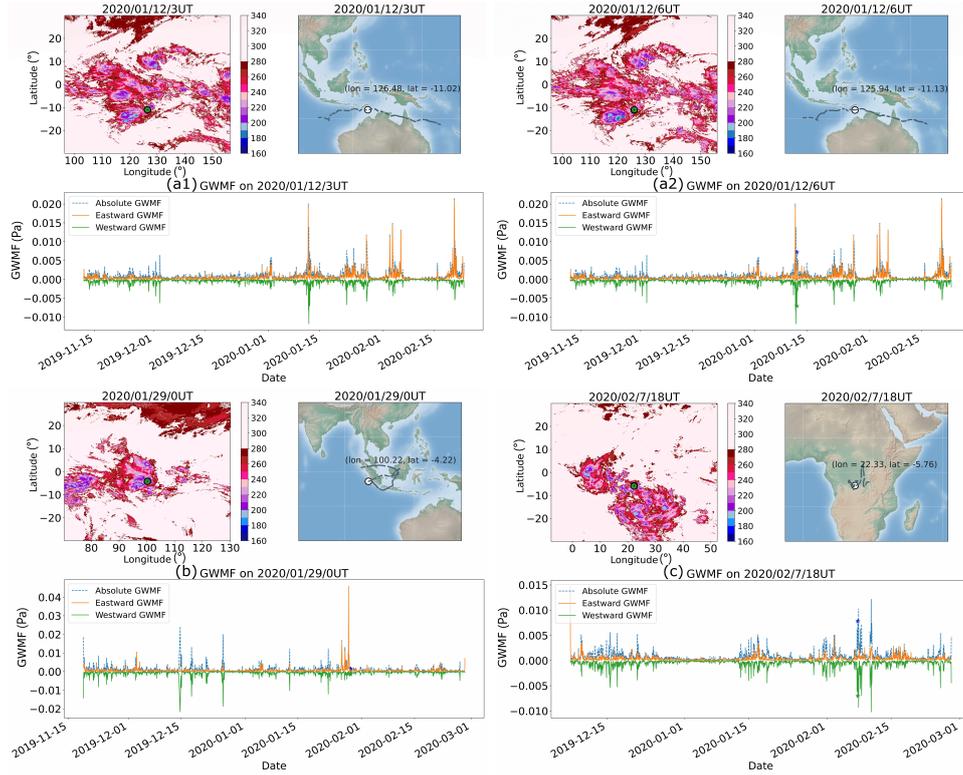


Figure 8. Brightness temperature from NOAA/NCEP GPM_MERGIR product (Janowiak, 2017), positions, and the corresponding observed GWMFs at the high-peak events of balloon 2 (top), balloon 3 (lower left) and balloon 7 (lower right).

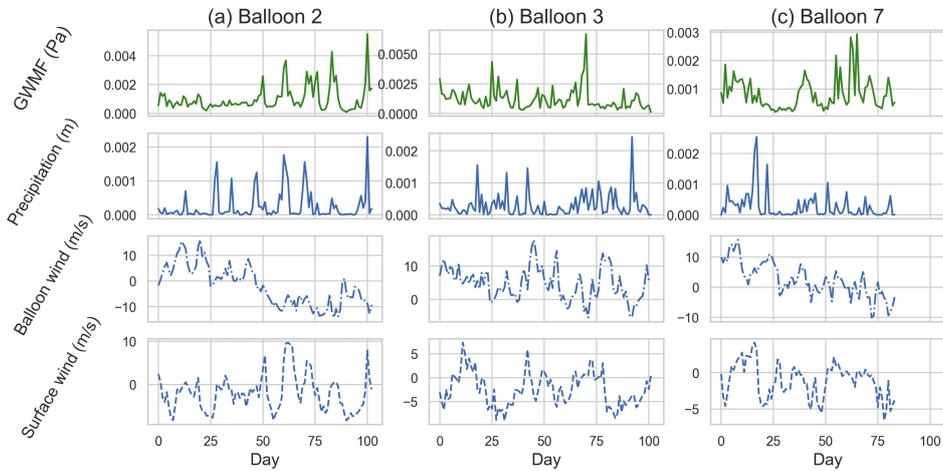


Figure 9. Time series of absolute GWMFs and the most informative ERA5 inputs in daily time resolution. The clear correspondence between precipitation and GWMF of balloon 2 can be visually observed in column (a). In contrast, this is not the case at all for balloon 3 as shown in column (b), and it partially presents in column (c) of balloon 7.

492
493

Estimates from superpressure balloon measurements were chosen as the target observations for gravity wave momentum fluxes (GWMF). The first campaign of the Stra-

teole 2 project (Haase et al., 2018) consisted of eight balloons flying an average of about 85 days each around the globe in the equatorial band. The quasi-Lagrangian nature of the balloons allows an accurate estimate of gravity wave momentum fluxes (Geller et al., 2013), the latter being a key quantity for parameterizations (Alexander et al., 2010). Analysis of the GWMF estimated from measurements in this first campaign has highlighted and confirmed convection as the main source of gravity waves in this region, especially for waves with high frequencies (periods shorter than one hour); see Corcos et al. (2021).

The description of the large-scale flow environment was provided from the ERA5 reanalysis, along with vertical profiles co-located with each balloon at each time. These variables included wind, pressure, temperature, and precipitation. The latter being a noisy and uncertain field, values of total precipitation were retrieved in a $500 \times 500 \text{ km}^2$ area around each balloon location, and was generally described by the mean and standard deviation over this area.

The ML models used are tree-based methods: random forests, extremely randomized trees, and adaptive boosting. Other methods were also investigated, as sensitivity experiments, without yielding major improvements. For each method, seven out of eight balloons were used for *training*, and the last balloon was used for *testing*.

The main results obtained from these investigations are as follows:

1. Based on the information provided by the large-scale flow data from ERA5, ML methods can reconstruct the observed GWMFs with correlations exceeding 0.7 in certain cases (balloon 2, 6, and 8), which is encouraging. Overall, the majority of the correlations are statistically significant at least at the 95% level, except for a few cases, as indicated in Table 2. The performances of ML methods, however, vary considerably from one balloon to another, with correlations down to 0.4 for some other balloons, and even down to 0.2 in one case. The overall average correlation for the HF case is 0.54, while a slightly lower average correlation of 0.49 is obtained in the WF case. In general, the correlations for WF waves are slightly weaker than those for HF waves (refer to the supplementary document for details).
2. The variations in performance are much larger between different balloons, than they are for a given balloon between ML approaches. This suggests that the performances are limited by the datasets, not by the choice of ML approach. The tree-based methods proved generally efficient, but there is not an overwhelming preference for one of them. Adaptive boosting frequently performed a bit better, but all three failed to capture the intensity of the (very intermittent) peaks in GWMF.
3. The most informative explanatory variables are those describing the precipitation and the zonal wind at and below the balloon's level. It is indeed an advantage of tree-based methods to provide information about the usefulness of the different inputs, e.g. through the Gini importance (Hastie et al., 2001). The importance of precipitation is consistent with the convective generation of the waves (Lott & Guez, 2013; Corcos et al., 2021). The importance of winds is consistent with the general understanding of the generation and propagation of waves (Kim et al., 2003); the relevance of wind at the balloon level is reminiscent of previous findings (Plougonven et al., 2017; Amiramjadi et al., 2023).
4. The ML methods were more efficient at reconstructing the part of GWMF associated with high-frequency waves (periods shorter than an hour) than the whole spectrum. This is consistent with the local character of the explanatory variables provided as inputs: high-frequency waves will be shorter-lived and propagate more vertically.
5. Different decompositions of the GWMF were used: absolute, eastward and westward GWMF. Interestingly, the performances significantly differed between these. The most difficult to reconstruct was found to be westward GWMF. Reasons for this likely include limitations of the dataset, to be further discussed below.

546 However, there are still parts where the large-scale flow are not informative enough
 547 in the estimation. There are cases where high peaks are present in the observed target,
 548 which indicates interesting events; however, large-scale flow are missed to describe them.
 549 As a result, the models failed to reconstruct such events in GWMFs (balloon 1 and 3,
 550 for example).

551 In addition, we have also implemented ML models by replacing ERA5's winds with
 552 balloon-observed winds at the balloon's level. This tests the sensitivity to errors in the
 553 input variables, for the variables for which we have direct observations, and which is known
 554 in the reanalysis to include significant error. The results suggest there is some sensitiv-
 555 ity, but it is not extensive. Overall, the performances on some challenging balloons such
 556 as balloon 3 and 5 are significantly improved when using observed winds instead of ERA5's
 557 winds. In contrast, the performance on balloon 8 drops quite a bit compared to the model
 558 with ERA5's winds. Overall, the models utilizing observed wind achieve an average cor-
 559 relation of 0.53 in the HF case and 0.47 in the WF case. These results can be found in
 560 the supplementary document.

561 5.2 Perspectives

562 Although the ML approaches have performed well, and nearly always better than
 563 the parameterization, there are clear limitations to the current investigation, calling for
 564 further research. The very strong sensitivity of the performances to the balloon that is
 565 left out and then used for testing is a clear indication that we lack data: the results strongly
 566 depend on the split of the data for training and testing, the performances are far from
 567 convergence. This is consistent with the strong intermittency of the GWMF (Hertzog
 568 et al., 2012; Plougonven et al., 2013) and with the illustrative time series of Figure 2:
 569 for each balloon, GWMF are dominated by a few events, such that even with 680 days
 570 of balloon measurements, only a few handfuls of GWMF peaks are described. This is
 571 too little for data-driven methods. This also explains why clear distinctions between the
 572 different methods are not found: the ML methods do their best but still lack data to clearly
 573 separate a better method for this problem, if there is one.

574 Ways forward include:

- 575 • Obtaining more observations to use as the target, keeping the same framework for
 576 the ML. Additional observations would come from the second Strateole 2 campaign
 577 (in 2021) and from Loon balloons (Schoeberl et al., 2017; Köhler et al., 2023). The
 578 additional Strateole data would enhance the data by less than a factor 2 and is
 579 therefore not expected to suffice to make a dramatic change. The Loon data would
 580 come with other difficulties as the observations were not made for research pur-
 581 poses and come with their own challenges (Green et al., 2023).
- 582 • Additional data could be provided not for the targets, but for the explanatory vari-
 583 ables. A first step could be including additional input variables from the reanal-
 584 yses. However, preliminary attempts have not suggested significant gains from the
 585 most evident additional culprits. A second step would consist of providing much
 586 more detailed and more accurate information about the background flow: this could
 587 be obtained from satellite observations, such as the observations of brightness tem-
 588 peratures from geostationary satellites shown in Figure 8. This would constitute
 589 a very interesting new study but in a profoundly new framework and with differ-
 590 ent aims: to fully use the information available from satellites would a priori re-
 591 quire providing maps (or images, or 2D fields) as input variables (more akin to
 592 Matsuoka et al. (2020), although their inputs were from models, not observations).
 593 The ML used would need to be reassessed (Matsuoka et al. (2020) used neural net-
 594 works, for instance). Such a study would be of great interest because the perfor-
 595 mance of the ML methods would much less be tainted by the uncertainty (or er-
 596 rors) present in the inputs that serve to describe the background. Additionally,

597 much more detailed information would be provided about the background flow,
 598 allowing the ML methods to tap into a greater reservoir of potentially relevant in-
 599 formation, and hence providing more precise answers regarding the relationship
 600 of the large-scale flow to the gravity wave signal. However, if the outcome of such
 601 an exercise would be of interest fundamentally, it would be more removed from
 602 the framework in which current parameterizations operate.

- 603 • A shortcoming of the present ML approaches is that they underestimate the peak
 604 values for GWMF (see Figures 2 and 3. This is expected, given the averaging in-
 605 volved in tree-based method and the limited number of strong events present in
 606 the training data. However, this implies that the distribution of reconstructed mo-
 607 mentum fluxes misses the tail of intense, rare events, which are known to matter
 608 for atmospheric gravity waves (Hertzog et al., 2012; de la Camara et al., 2016).
 609 One way to overcome this would be to aim not at a deterministic reconstruction
 610 of the momentum fluxes, but at reconstructing a probability density function of
 611 these. This change of framework, equivalent to changing from a deterministic to
 612 a stochastic parameterization, would in fact be more consistent for three reasons:
 613 first, given some large-scale conditions, there are certainly several different small-
 614 scale configurations with different resulting gravity waves that can occur. Second,
 615 for any given realization of the small-scale flow corresponding to large-scale con-
 616 ditions, our observed values depend on the specific sampling by the balloon. At
 617 present, we do not fully know how sensitive the observed gravity wave momen-
 618 tum fluxes are to this sampling. Finally, the estimate of gravity wave momentum
 619 fluxes from the observed balloon measurements involves assumptions and method-
 620 ological choices, and there is as always an observational error in the estimates for
 621 GWMF. Given that the ML methods do capture rather well the occurrence of larger
 622 values, using ML methods to reconstruct a PDF of likely fluxes, rather than a sin-
 623 gle, deterministic value, could give room to better represent the observed GWMF,
 624 although only in a probabilistic way.
- 625 • A fourth way forward consists in applying similar investigations on datasets where
 626 more data is available, albeit at the cost of more uncertainty on the realism of the
 627 data. High-resolution models such as global convection permitting simulations (Stephan
 628 et al., 2019) provide a wealth of information on the resolved gravity wave field,
 629 and many studies have repeatedly highlighted the ability of models to simulate
 630 efficiently many features of the observed gravity wave field (Plougonven & Teit-
 631 elbaum, 2003; Wu & Eckermann, 2008; Preusse et al., 2014; Stephan et al., 2019).
 632 Model output from global simulations would provide amounts of data for which
 633 the sampling limitations of the Strateole balloons would not be present. The down-
 634 side is the limitations of model data, relative to observations, and the need for strate-
 635 gies to validate which aspects of the simulations are realistic.

636 Acknowledgments

637 This work and Sothea Has are supported by the Institut des Mathématiques pour la Planète
 638 Terre (IMPT). This work has also received support from the ANR project BOOST3R
 639 (ANR-17-CE01-0016-01) and the French-American project Strateole 2 (CNES). More-
 640 over, we gratefully acknowledge the support and collaborative efforts extended by mem-
 641 bers of the DataWave consortium, a Virtual Earth System Research Institute (VESRI)
 642 Schmidt Futures project.

643 6 Open research

644 Balloon data used in this study are presented in Haase et al. (2018) of the STRA-
 645 TEOLE 2 mission and can be extracted from the following website: <https://webstr2.ipsl.polytechnique.fr>. The ERA5 input variables are described in Hersbach et al. (2020) and can be obtained from the COPERNICUS access hub using the following web-

648 site: <https://scihub.copernicus.eu/>. The machine learning algorithms implemented
 649 in our analysis are available in the `scikit-learn` Python library (Pedregosa et al., 2011)
 650 and can be downloaded from its website: [https://scikit-learn.org/stable/install](https://scikit-learn.org/stable/install.html)
 651 [.html](https://scikit-learn.org/stable/install.html). Finally, the source codes for implementing machine learning methods in our anal-
 652 ysis are made available at the following GitHub repository: [https://github.com/hassothea/](https://github.com/hassothea/Reconstruction_of_GWMF_using_ML_ERA5)
 653 `Reconstruction_of_GWMF_using_ML_ERA5`.

654 References

- 655 Alexander, M., Geller, M., McLandress, C., Polavarapu, S., Preusse, P., Sassi, F.,
 656 ... Watanabe, S. (2010, July). Recent developments in gravity-wave effects in
 657 climate models and the global distribution of gravity-wave momentum flux from
 658 observations and models. *J. Geophys. Res.*, *115*, 1103-1124.
- 659 Amiranjadi, M., Plougonven, R., Mohebalhojeh, A. R., & Mirzaei, M. (2023). Using
 660 machine learning to estimate nonorographic gravity wave characteristics at source
 661 levels. *Journal of the Atmospheric Sciences*, *80*(2), 419-440.
- 662 Baker, W. E., Atlas, R., Cardinali, C., Clement, A., Emmitt, G. D., Gentry, B. M.,
 663 ... others (2014). Lidar-Measured Wind Profiles: The Missing Link in the Global
 664 Observing System. *Bull. Am. Meteor. Soc.*, *95*(10.1175/2010JAS3455.1), 543-564.
- 665 Benjamin Bossan, W. K., Josef Feigl. (2015). *Otto group product classification chal-*
 666 *lenge*. Kaggle. Retrieved from [https://kaggle.com/competitions/otto-group](https://kaggle.com/competitions/otto-group-product-classification-challenge)
 667 [-product-classification-challenge](https://kaggle.com/competitions/otto-group-product-classification-challenge)
- 668 Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov,
 669 V., ... others (2020). Presentation and evaluation of the ipsl-cm6a-lr climate
 670 model. *Journal of Advances in Modeling Earth Systems*, *12*(7).
- 671 Box, G. E., & Cox, D. R. (1964). An analysis of transformations. *Journal of the*
 672 *Royal Statistical Society. Series B (Methodological)*, *26*(2), 211-252. doi: 10.1111/
 673 [j.2517-6161.1964.tb00553.x](https://doi.org/10.1111/j.2517-6161.1964.tb00553.x)
- 674 Breiman, L. (2001). Random forests. *Machine Learning*, *45*(1), 5-32.
- 675 Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. A. (1984). Classification and
 676 regression trees. *Wadsworth International Group*.
- 677 Bushell, A. C., Butchart, N., Derbyshire, S. H., Jackson, D. R., Shutts, G. J.,
 678 Vosper, S. B., & Webster, S. (2015). Parameterized gravity wave momentum
 679 fluxes from sources related to convection and large-scale precipitation processes
 680 in a global atmosphere model. *Journal of the Atmospheric Sciences*, *72*(11),
 681 4349-4371.
- 682 Butchart, N. (2022). The stratosphere: a review of the dynamics and variability.
 683 *Weather and Climate Dynamics*, *3*(4), 1237-1272.
- 684 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Ma-
 685 chine learning emulation of gravity wave drag in numerical weather forecasting.
 686 *Journal of Advances in Modeling Earth Systems*, *13*(7), e2021MS002477.
- 687 Corcos, M., Hertzog, A., Plougonven, R., & Podglajen, A. (2021). Observation of
 688 gravity waves at the tropical tropopause using superpressure balloons. *Journal of*
 689 *Geophysical Research: Atmospheres*, *126*(15), e2021JD035165.
- 690 de Burgh-Day, C. O., & Leeuwenburg, T. (2023). Machine learning for numer-
 691 ical weather and climate modelling: a review. *Geoscientific Model Development*,
 692 *16*(22), 6433-6477. Retrieved from [https://gmd.copernicus.org/articles/16/](https://gmd.copernicus.org/articles/16/6433/2023/)
 693 [6433/2023/](https://gmd.copernicus.org/articles/16/6433/2023/) doi: 10.5194/gmd-16-6433-2023
- 694 de la Camara, A., Lott, F., & Hertzog, A. (2014). Intermittency in a stochastic
 695 parameterization of nonorographic gravity waves. *J. Geophys. Res. Atmos.*, *119*,
 696 11,905-11,919. doi: 10.1002/2014JD022002
- 697 de la Camara, A., Lott, F., Jewtoukoff, V., Plougonven, R., & Hertzog, A. (2016).
 698 On the gravity wave forcing during the southern stratospheric final warming in
 699 LMDz. *J. Atmos. Sci.*, *73*, 3213-3226. doi: 10.1175/JAS-D-15-0377.1

- 700 Ern, M., Diallo, M. A., Khordakova, D., Krisch, I., Preusse, P., Reitebuch, O., ...
 701 Riese, M. (2023). The quasi-biennial oscillation (qbo) and global-scale tropi-
 702 cal waves in aeolus wind observations, radiosonde data, and reanalyses. *Atmo-*
 703 *spheric Chemistry and Physics*, *23*(16), 9549–9583. Retrieved from [https://](https://acp.copernicus.org/articles/23/9549/2023/)
 704 acp.copernicus.org/articles/23/9549/2023/ doi: 10.5194/acp-23-9549-2023
- 705 Ern, M., Ploeger, F., Preusse, P., Gille, J., Gray, L. J., Kalisch, S., ... Riese,
 706 M. (2014). Interaction of gravity waves with the QBO: A satellite perspec-
 707 tive. *Journal of Geophysical Research: Atmospheres*, *119*, 2329 - 2355. doi:
 708 10.1002/2013JD020731
- 709 Ern, M., Preusse, P., & Riese, M. (2022). Intermittency of gravity wave poten-
 710 tial energies and absolute momentum fluxes derived from infrared limb sounding
 711 satellite observations. *Atmospheric Chemistry and Physics*, *22*(22), 15093–15133.
 712 Retrieved from <https://acp.copernicus.org/articles/22/15093/2022/> doi:
 713 10.5194/acp-22-15093-2022
- 714 Ern, M., Trinh, Q. T., Gille, P. P. J., Mlynczak, M., Russell, J., & Riese, M. (2018).
 715 GRACILE: a comprehensive climatology of atmospheric gravity wave parameters
 716 based on satellite limb soundings. *Earth System Science Data*, *10*, 857–892. doi:
 717 10.5194/essd-10-857-2018
- 718 Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J.
 719 (2022). Machine learning gravity wave parameterization generalizes to capture
 720 the qbo and response to increased co2. *Geophysical Research Letters*, *49*(8),
 721 e2022GL098174.
- 722 F. Vitart and A.W. Robertson (Ed.). (2018). *Sub-seasonal to seasonal prediction*. El-
 723 sevier.
- 724 Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-
 725 line learning and an application to boosting. *Journal of Computer and System*
 726 *Sciences*, *55*(1), 119–139.
- 727 Fritts, D., & Alexander, M. (2003). Gravity wave dynamics and effects in the
 728 middle atmosphere. *Reviews of Geophysics*, *41*(1), 1003. doi: doi:10.1029/
 729 2001RG000106
- 730 Geller, M., Alexander, M., Love, P., Bacmeister, J., Ern, M., Hertzog, A., ... Zhou,
 731 T. (2013). A comparison between gravity wave momentum fluxes in observations
 732 and climate models. *J. Clim.*, *26*, 6383–6405. doi: 10.1175/JCLI-D-12-00545.1
- 733 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could
 734 machine learning break the convection parameterization deadlock? *Geophys. Res.*
 735 *Lett.*, *45*, 5742–5751. doi: 10.1029/2018GL078202
- 736 Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine*
 737 *Learning*, *63*(1), 3–42.
- 738 Green, B., Sheshadri, A., Alexander, M. J., Bramberger, M., & Lott, F. (2023).
 739 Gravity wave momentum fluxes estimated from project loon balloon data. *Au-*
 740 *thorea Preprints*.
- 741 Haase, J., Alexander, M., Hertzog, A., Kalnajs, L., Deshler, T., Davis, S., ...
 742 Venel, S. (2018, March). Around the world in 84 days. *EOS*, *99*. doi:
 743 10.1029/2018EO091907
- 744 Hastie, T., Tibshirani, R., & Friedman, J. (2001). *The elements of statistical learn-*
 745 *ing*. New York, NY, USA: Springer New York Inc.
- 746 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater,
 747 J., ... Thépaut, J.-N. (2020). The ERA5 global reanalysis. *Quarterly Jour-*
 748 *nal of the Royal Meteorological Society*, *146*(730), 1999–2049. Retrieved from
 749 <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3803> doi:
 750 <https://doi.org/10.1002/qj.3803>
- 751 Hertzog, A., Alexander, M., & Plougonven, R. (2012). On the probability density
 752 functions of gravity waves momentum flux in the stratosphere. *J. Atmosph. Sci.*,
 753 *69*, 3433–3448.

- 754 Janowiak, J. B. . X. P., J. (2017). *Ncep/cpc l3 half hourly 4km global (60s - 60n)*
755 *merged ir v1*. In A. Savtchenko, & M. D. Greenbelt (Eds.), Goddard Earth sci-
756 ences data and information services center (GES DISC). (Accessed 2020).
- 757 Jewtoukoff, V., Hertzog, A., Plougonven, R., de la Camara, A., & Lott, F. (2015).
758 Gravity waves in the Southern Hemisphere derived from balloon observations and
759 ECMWF analyses. *J. Atmos. Sci.*, *72*, 3449-3468.
- 760 Kim, Y.-J., Eckermann, S., & Chun, H.-Y. (2003). An overview of the past,
761 present and future of gravity-wave drag parametrization for numerical climate
762 and weather prediction models. *Atmosphere-Ocean*, *41*, 65-98.
- 763 Köhler, L., Green, B., & Stephan, C. C. (2023). Comparing loon superpressure
764 balloon observations of gravity waves in the tropics with global storm-resolving
765 models. *Journal of Geophysical Research: Atmospheres*, *128*(15), e2023JD038549.
- 766 Kremser, S., Thomason, L. W., von Hobe, M., Hermann, M., Deshler, T., Timmreck,
767 C., ... others (2016). Stratospheric aerosol—observations, processes, and impact
768 on climate. *Reviews of Geophysics*, *54*(2), 278–335.
- 769 Lott, F., & Guez, L. (2013). A stochastic parameterization of the gravity waves due
770 to convection and its impact on the equatorial stratosphere. , *118*, 8897-8909.
- 771 Lott, F., Rani, R., Podglajen, A., Codron, F., Guez, L., Hertzog, A., & Plougonven,
772 R. (2023). Direct comparison between a non-orographic gravity wave drag scheme
773 and constant level balloons. *Journal of Geophysical Research: Atmospheres*,
774 *128*(4), e2022JD037585.
- 775 Matsuoka, D., Watanabe, S., Sato, K., Kawazoe, S., Yu, W., & Easterbrook, S.
776 (2020). Application of deep learning to estimate atmospheric gravity wave
777 parameters in reanalysis data sets. *Geophysical Research Letters*, *47*(19),
778 e2020GL089436.
- 779 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
780 ... Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of*
781 *Machine Learning Research*, *12*, 2825–2830.
- 782 Plougonven, R., de la Camara, A., Hertzog, A., & Lott, F. (2020). How does
783 knowledge of atmospheric gravity waves guide their parameterizations? *Q.J.*
784 *Roy. Meteorol. Soc.*, 1-15. doi: 10.1002/qj.3732
- 785 Plougonven, R., Hertzog, A., & Guez, L. (2013, January). Gravity waves over
786 Antarctica and the Southern Ocean: consistent momentum fluxes in mesoscale
787 simulations and stratospheric balloon observations. *Quart. J. Roy. Meteorolog.*
788 *Soc.*, *139*, 101-118.
- 789 Plougonven, R., Jewtoukoff, V., de la Camara, A., Hertzog, A., & Lott, F.
790 (2017). On the relation between gravity waves and wind speed in the lower
791 stratosphere over the Southern Ocean. *J. Atmos. Sci.*, *74*, 1075-1093. doi:
792 10.1175/JAS-D-16-0096.1
- 793 Plougonven, R., & Teitelbaum, H. (2003). Comparison of a large-scale inertia-
794 gravity wave as seen in the ECMWF and from radiosondes. *Geophys. Res. Let.*,
795 *30*(18), 1954.
- 796 Podglajen, A., Hertzog, A., Plougonven, R., & Zagar, N. (2014). Assessment of the
797 accuracy of (re)analyses in the equatorial lower stratosphere. , *119*, 11,166-11,188.
798 doi: 10.1002/2014JD021849
- 799 Preusse, P., Ern, M., Bechtold, P., Eckermann, S., Kalisch, S., Trinh, Q., & Riese,
800 M. (2014). Characteristics of gravity waves resolved by ECMWF. *Atmos. Chem.*
801 *Phys.*, *14*, 10483-10508. doi: 10.5194/acp-14-10483-2014
- 802 Rättsch, G., Onoda, T., & Müller, K.-R. (2001). Soft margins for adaboost. *Machine*
803 *Learning*, *42*.
- 804 Schoeberl, M. R., Jensen, E., Podglajen, A., Coy, L., Lodha, C., Candido, S., &
805 Carver, R. (2017). Gravity wave spectra in the lower stratosphere diagnosed from
806 project loon balloon trajectories. *Journal of Geophysical Research: Atmospheres*,
807 *122*(16), 8517–8524.

- 808 Solomon, S., Rosenlof, K., Portmann, R., Daniel, J., Davis, S., Sanford, T., &
809 Plattner, G.-K. (2010, March). Contributions of stratospheric water vapor
810 to decadal changes in the rate of global warming. *Science*, 1219-1223. doi:
811 10.1126/science.118248
- 812 Stephan, C., Strube, C., Klocke, D., Ern, M., Hoffmann, L., Preusse, P., & Schmidt,
813 H. (2019). Gravity Waves in Global High-Resolution Simulations With Ex-
814 plicit and Parameterized Convection. *J. Geophys. Res.*, 124, 4446-4459. doi:
815 10.1029/2018JD030073
- 816 Trinh, Q., Kalisch, S., Preusse, P., Ern, M., Chun, H., Eckermann, S., . . . Riese, M.
817 (2016). Tuning of a gravity wave source scheme based on HIRDLS observations.
818 *Atmos. Chem. Phys.*, 16, 7335-7356. doi: 10.5194/acp-16-7335-2016
- 819 Vincent, R., & Hertzog, A. (2014). The response of superpressure balloons to gravity
820 wave motions. *Atmospheric Measurement Techniques*, 7(4), 1043-1055.
- 821 Wright, C., Osprey, S., & Gille, J. (2013). Global observations of gravity wave inter-
822 mittency and its impact on the observed momentum flux morphology. *J. Geophys.*
823 *Res.*, 118, 10,980-10,993. doi: 10.1002/jgrd.50869
- 824 Wu, D., & Eckermann, S. (2008). Global gravity wave variances from Aura MLS:
825 characteristics and interpretation. , 65(12), 3695-3718.
- 826 ZEWEICHU. (2019). *2019 ttic 31020 hw4 spam (adaboost)*. Kaggle. Retrieved from
827 <https://kaggle.com/competitions/2019-ttic-31020-hw4-spam-adaboost>