

1 Using Machine Learning to Model Uncertainty for Water-Vapor 2 Atmospheric Motion Vectors

3 Joaquim V. Teixeira¹, Hai Nguyen¹, Derek J. Posselt¹, Hui Su¹, Longtao Wu¹

4 ¹Jet Propulsion Laboratory, California Institute of Technology

5 **Abstract.** Wind-tracking algorithms produce Atmospheric Motion Vectors (AMVs) by tracking clouds or water vapor
6 across spatial-temporal fields. Thorough error characterization of wind-tracking algorithms is critical in properly
7 assimilating AMVs into weather forecast models and climate reanalysis datasets. Uncertainty modelling should yield
8 estimates of two key quantities of interest: bias, the systematic difference between a measurement and the true value,
9 and standard error, a measure of variability of the measurement. The current process of specification of the errors in
10 inverse modelling is often cursory and commonly consists of a mixture of model fidelity, expert knowledge, and need
11 for expediency. The method presented in this paper supplements existing approaches to error specification by
12 providing an error-characterization module that is purely data-driven. Our proposed error-characterization method
13 combines the flexibility of machine learning (random forest) with the robust error estimates of unsupervised
14 parametric clustering (using a Gaussian Mixture Model). Traditional techniques for uncertainty modeling through
15 machine learning have focused on characterizing bias, but often struggle when estimating standard error. In contrast,
16 model-based approaches such as k-means or Gaussian mixture modelling can provide reasonable estimates of both
17 bias and standard error, but they are often limited in complexity due to reliance on linear or Gaussian assumptions. In
18 this paper, a methodology is developed and applied to characterize error in tracked-wind using a high-resolution global
19 model simulation, and it is shown to provide accurate and useful error features of the tracked wind.

20 1. Introduction

21 Reliable estimates of global winds are critical to science and application areas, including global chemical transport
22 modeling and numerical weather prediction. One source of wind measurements consists of feature-tracking based
23 Atmospheric Motion Vectors (AMVs), produced by tracking time sequences of satellite-based measurements of
24 clouds or spatially distributed water vapor fields (Mueller et al., 2017; Posselt et al., 2019). The importance of global
25 measurements of 3-dimensional winds was highlighted as an urgent need in the NASA Weather Research Community
26 Workshop Report (Zeng et al., 2016) and was identified as a priority in the 2007 National Academy of Sciences Earth
27 Science and Applications from Space (ESAS 2007) Decadal Survey and again in ESAS 2017. For instance, wind is
28 used in the study of global CO₂ transport (Kawa et al., 2004), numerical weather prediction (NWP; Cassola and
29 Burlando, 2012), as inputs into weather and climate reanalysis studies (Swail and Cox, 2000), and for estimating
30 current and future wind-power outputs (Staffell and Pfenninger, 2016).

31 Thorough error characterization of wind-track algorithms is critical in properly assimilating AMVs into forecast
32 models. Prior literature has explored the impact of ‘poor’ error-characterization in Bayesian-based approaches to
33 remote sensing applications. Nguyen et al. (2019) proved analytically that when the input bias is incorrect in Bayesian

34 methods (specifically, optimal estimation retrievals), then the posterior estimates would also be biased. Moreover,
35 they proved that when the input standard error is ‘correct’ (that is, it is as close to the unknown truth as possible), then
36 the resulting Bayesian estimate is ‘efficient’; that is, it has the smallest error among all possible choices of prior
37 standard error. Additionally, multiple active and passive technologies are being developed to measure 3D winds, such
38 as Doppler wind lidar (DWL), radar, and infrared/microwave sensors that derive AMVs using feature-tracking of
39 consecutive images. Therefore, an accurate and robust methodology for modeling uncertainty will allow for more
40 accurate assessments of mission impacts, and the eventual propagation of data uncertainties for these instruments.

41 Velden and Bedka (2009) and Salonen et al. (2015) have shown that height assignment contributes a large component
42 of uncertainty in AMVs tracked from cloud movement and from sequences of infrared satellite radiance images.
43 However, with AMVs obtained from water vapor profiling instruments (e.g., infrared and microwave sounders),
44 height assignment error cannot be directly assessed purely through analysis of the AMV extraction algorithm. Height
45 assignment is instead an uncertainty in the water vapor profile itself. Unfortunately, without the quantified
46 uncertainties on the water vapor profile necessary to pursue such a study, that is well beyond the scope of this paper.
47 As such, this study will focus on errors in the AMV estimates at a given height. Previous work has demonstrated
48 several different approaches for characterizing AMV vector error. One common approach is to employ quality
49 indicator thresholds, as described by Holmund et al (2001), which compare changes in AMV estimates between
50 sequential timesteps and neighboring pixels, as well as differences from model predictions, to produce a quality
51 indicator to which a discrete uncertainty is assigned. The Expected Error approach, developed by Le Marshal et al.
52 (2004), builds a statistical model using linear regression against AMV-radiosonde values to estimate the statistical
53 characteristics of AMV observation error.

54 In this study, we outline a data-driven approach for building an AMV uncertainty model using observing system
55 simulation experiment (OSSE) data. We build on the work by Posselt et al. (2019) in which a water vapor feature-
56 tracking AMV algorithm was applied to a high-resolution numerical simulation, thus providing a global set of AMV
57 estimates which can be compared to the reference winds produced by the simulation. In this case, a synthetic “true”
58 state is available with which AMVs can be compared and errors are quantified, and it is shown that errors in AMV
59 estimates are state dependent. Our approach will use a conjunction of machine learning (random forest) and
60 unsupervised parametric clustering (Gaussian mixture models) to build a model for the uncertainty structures found
61 by Posselt et al. (2019). The realism and robustness of the resulting uncertainty estimates depend on the realism and
62 representativeness of the reference dataset. This work builds upon the work of Bormann et al. (2014) and Hernandez-
63 Carrascal and Bormann (2014), who showed that wind tracking could be divided into distinct geophysical regimes by
64 clustering based on cloud conditions. This study supplements that approach with the addition of machine learning,
65 which, compared with traditional linear modeling approaches, should allow the model to capture more complex non-
66 linear processes in the error function.

67 Traditional techniques for modeling uncertainty through machine learning have focused on characterizing bias but
68 often struggle when estimating standard error. By pairing a random forest algorithm with unsupervised parametric

69 clustering, we propose a data-driven, cluster-based approach for quantifying both bias and standard error from
70 experimental data. According to the theory developed by Nguyen et al. (2019), these improved error characterizations
71 should then lead to improved error characteristics (e.g., lower bias, more accurate uncertainties) in subsequent analyses
72 such as flux inversion or data assimilation.

73 This paper does not purport that the specific algorithm detailed here should supplant error characterization approaches
74 for all AMVs; indeed, most commonly assimilated AMVs are based on tracking cloud features, not water vapor
75 profiles. In addition, this algorithm is trained and developed for a specific set of AMVs extracted from a water vapor
76 field associated with a particular range of flow features. As such, application of our algorithm to modeled or observed
77 AMVs will be most appropriate in situations with similar dynamics to our training set. However, we intend in this
78 paper to demonstrate that the methodology is successful in characterizing errors for this set of water vapor AMVs and
79 suggest that this approach— that is, capturing state-dependent uncertainties in feature-tracking algorithms through a
80 combination of clustering and random forest— could be implemented in other feature-tracking AMV extraction
81 methods and situations.

82 The rest of the paper is organized as follows: In Section 2, we give an overview of the simulation which provides the
83 training data for our machine learning approach. We then motivate and define the specific uncertainties this study
84 aims to characterize. In Section 3, we describe the error characterization approach with the specifics of our error
85 characterization model, including both the implementation of and motivations for employing the random forest and
86 Gaussian mixture model. In Section 4, we provide a validation of our methods, attempting to assess the bias of our
87 predictions. In Section 5, we discuss the implications of our error characterization approach, both on AMV estimation
88 and data assimilation more broadly.

89 **2. Experimental Set-up**

90 **2.1 Simulation and Feature-Tracking Algorithm**

91 We trained our model on the simulated data used by Posselt et al. (2019), which applied an AMV algorithm to outputs
92 from the NASA Goddard Space Flight Center (GSFC) Global Modeling and Assimilation Office (GMAO) GEOS-5
93 Nature Run (G5NR; Putman et al. 2014). The Nature Run is a global dataset with ~7 km horizontal grid spacing that
94 includes, among other quantities, three-dimensional fields of wind, water vapor concentration, clouds, and
95 temperature. Note that throughout the text we will use the term ‘Nature Run wind’ to refer to reference winds in the
96 simulation dataset used to train the uncertainty model. The AMV algorithm is applied on four pressure levels (300hPa,
97 500hPa, 700hPa, and 850hPa) at 6-hourly intervals, using three consecutive global water vapor fields spaced one hour
98 apart, and for a 60-day period from 07/01/2006 to 08/30/2006. The water-vapor fields from GEOS5 were input to a
99 local-area pattern matching algorithm that approximates wind speed and direction from movement of the matched
100 patterns. The algorithm searches a pre-set number of nearby pixels to minimize the sum-of-absolute-differences
101 between aggregated water vapor values across the pixels. Posselt et al. (2019) describes the sensitivity of the tracking

102 algorithm and the dependency of the tracked winds on atmospheric states in detail. The coordinates of the data are on
103 a 5758 x 2879 x 240 spatio-temporal grid for the longitude, latitude, and time dimension, respectively.

104 It is important to note that the AMV algorithm tracks water vapor on fixed pressure levels. In practice, these would be
105 provided by satellite measurements, whereas in this paper we use simulated water vapor from the GEOS-5 Nature
106 Run. In this simulation height assignment of the AMVs is assumed to be perfectly known. This assumption is far from
107 guaranteed in real world applications but, as previously discussed, its implications are not pursued in this paper. As
108 such, we focus solely on observational AMV error and not on height assignment error. We note that in practice, one
109 approach to understanding the behavior and accuracy of the wind-tracking algorithm is to apply it to modeled data
110 (e.g., Posselt et al., 2019). Our approach seeks to complement this approach by providing a machine-
111 learning/clustering hybrid approach that can further divide comparison domains into ‘regimes’ which may provide
112 further insights into the behavior of the errors and/or feedback into the wind-tracking algorithm.

113 A snapshot of the dataset at 700hPa is given in Figure 1, where we display the water vapor from Nature Run (top left
114 panel), the wind speed from Nature Run (top right panel), the tracked wind from the AMV-tracking algorithm (bottom
115 right panel), and the difference between the Nature Run and tracked wind (bottom left panel). Note that the wind-
116 tracking algorithm tends to have trouble in region where the Nature Run water vapor content is close to zero. It is clear
117 that while the wind-tracking algorithm tends to perform well in most regions (we can classify these regions as areas
118 where the algorithm is skilled), in some regions the algorithm is unable to reliably make a reasonable estimate of the
119 wind speed (unskilled). We will examine these skilled and unskilled regimes (and their corresponding contributing
120 factors) in section 3.

121 2.2 Importance of Uncertainty Representation in Data Assimilation

122 Proper error characterization for any measurement, including AMVs, is important in data assimilation. Data
123 assimilation often uses a regularized matrix inverse method based on Bayes’ theorem, which, when all probability
124 distributions in Bayes’ relationship are assumed to be Gaussian, reduces to minimizing a least-squares (quadratic) cost
125 function Eq (1):

$$126 \quad \mathbf{J} = (\mathbf{x} - \mathbf{x}_b)\mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b) + ((\hat{\mathbf{y}} - \mathbf{a}) - \mathbf{H}[\mathbf{x}])^T \mathbf{R}^{-1}((\hat{\mathbf{y}} - \mathbf{a}) - \mathbf{H}[\mathbf{x}]) \quad (1)$$

127 where \mathbf{x} represents the analysis value, \mathbf{x}_b represents the background field (first guess), \mathbf{B} represents the background
128 error covariance, \mathbf{y} represents the observation, and \mathbf{H} represents the forward operator that translates model space into
129 observation space. This translation may consist of spatial and/or temporal interpolation if \mathbf{x} and \mathbf{y} are the same variable
130 (e.g., if the observation of temperature comes from a radiosonde), or may be far more complicated (e.g., a radiative
131 transfer model in the case of satellite observations). \mathbf{R} represents the observation error covariance, and \mathbf{a} represents
132 the accuracy, or bias, in the observations. The right hand side of Eq. (1) can be interpreted as a sum of the contribution
133 of information from the data ($\mathbf{y} - \mathbf{H}[\mathbf{x}] - \mathbf{a}$) and the contribution from the prior ($\mathbf{x} - \mathbf{x}_b$), which are weighted by their

134 respective covariance matrices. In our analysis, the AMVs obtained from the wind-tracking algorithm is used as ‘data’
135 in subsequent analysis. That is, the tracked wind data $\hat{\mathbf{y}}$ is a biased and noisy estimator of the true wind \mathbf{y} , and might
136 be assumed to follow the model Eq. (2):

$$137 \quad \hat{\mathbf{y}} = \mathbf{y} + \epsilon \quad (2)$$

138 where ϵ is an error term, commonly assumed to be Gaussian with mean \mathbf{a} and covariance matrix \mathbf{R} (i.e., $\epsilon \sim N(\mathbf{a}, \mathbf{R})$),
139 which are the same two terms that appear in Equation (1). As such, for data assimilation to function, it is essential to
140 correctly specify the bias vector \mathbf{a} and the standard error matrix \mathbf{R} . Incorrect characterizations of either of these
141 components could have adverse consequences on the resulting data assimilation analyses with respect to bias and/or
142 the standard error (Nguyen et al., 2019).

143 **3 Methodology**

144 **3.1 Generalized Error Characterization Model**

145 An overview of our approach is outlined in Figure 2. Given a set of training predictors X , training responses \hat{Y} , and
146 simulated true response Y , our approach begins with two independent steps. In one step, a Gaussian mixture model is
147 trained on the set of X , \hat{Y} , and Y . This clustering algorithm identifies geophysical regimes where the nonlinear
148 relationships between the three variables differ. In the other step, a random forest is used to model Y based on X and
149 \hat{Y} . This step produces an estimate of the true response (we call this \tilde{Y}) using only the training predictors and response.
150 We then employ the Gaussian mixture model to estimate the clusters which the set of X , \hat{Y} , and \tilde{Y} pertain to.
151 Subsequently, we compute the error characteristics of each cluster of X , \hat{Y} , and \tilde{Y} in the training dataset. Thereafter,
152 given a new point consisting solely of X and \hat{Y} , we can assign it to a specific cluster and ascribe to it a set of error
153 characteristics.

154 In this paper, we are primarily interested in the distribution of a retrieved quantity versus the truth. That is, given a
155 retrieved value \hat{Y}_i , we are interested in the first and second moments (i.e., $E(\hat{Y}_i - Y)$ and $\text{var}(\hat{Y}_i - Y)$), respectively.
156 We note that there is a large body of existing work on uncertainty modeling in the machine learning literature (e.g.,
157 Coulston et al., 2016; Tripathy et al., 2018; Tran et al., 2019; Kwon et al., 2020), although these approaches primarily
158 define the uncertainty of a prediction as $\text{var}(\hat{Y}_i)$, or quantify how sensitive that prediction is to tiny changes in the
159 models/inputs. Our approach, on the other hand, characterizes the error as $\text{var}(\hat{Y}_i - Y)$, which describes how accurate
160 a prediction is relative to the *true value*. For this reason, our methodology is more stringent in that it requires
161 knowledge of the true field (which comes naturally within OSSE framework) or some proxies such as independent
162 validation data or reanalysis data. In return, the error estimates from our methodology fit naturally within the data
163 assimilation framework (that is, it constitutes the parameter \mathbf{R} in Eq. (1)).

164 What follows in this paper is an implementation of the error characterization model obtained for a subsample of the
165 GEOS-5 Nature Run at a fixed height of 700hPa. In particular, we trained the error characterization on a random

166 subsample from the first 1.5 months of the Nature Run, and show the results obtained when applying it to a test
167 subsample drawn from the subsequent 0.5 months of the Nature Run.

168 3.2 Error Regime

169 When examining the relationship between AMVs and Nature Run winds in Figure 3, it is clear that there are two
170 distinct ‘error-regimes’ present in the dataset. The majority of AMV estimates can be categorized as ‘skilled’, wherein
171 their estimate lies clearly along a one-to-one line with the Nature Run wind. However, there is also clearly an
172 ‘unskilled’ regime, for which the AMV estimate is very close to zero when there are actually moderate or large Nature
173 Run wind values present. Our goal is to provide unique error characterizations for each error regime, because the error
174 dynamics are different within each regime. Furthermore, when we analyze this error and its relationship to water
175 vapor, we see that ‘unskilled’ regime correlates highly with areas of low water vapor in Figure 4. This matches the
176 error patterns discussed in Posselt et al. (2019). We note that the division between skilled and unskilled regimes does
177 not need to be binary. For instance, in some regions the wind-tracking algorithm might be unbiased with high-
178 correlation with the true winds, and in other regions the algorithm might still be unbiased relative to the true winds,
179 but with higher errors. The second situation is clearly less skilled than the first, although it might still be considered
180 ‘skilled’, and this separation of the wind-tracking estimates into various ‘grades’ of skill forms the basis of our error
181 model.

182 3.3 Gaussian Mixture Model

183 These distinct regimes present an opportunity to employ machine learning. Bormann et al. (2014) and Hernandez-
184 Carrascal and Bormann (2014) demonstrated that cluster (also called regime) analysis is a successful approach for
185 wind-tracking error characterization, and so we aim to train a clustering algorithm that will cluster a given individual
186 AMV estimate to various ‘grades’ of skill. In particular, we use a clustering algorithm that can take advantage of the
187 underlying geophysical dynamics. To this end, we employ a Gaussian mixture model, an unsupervised clustering
188 algorithm based on estimating a training set as a mixture of multiple Gaussian distributions. A mathematical overview
189 follows:

- 190 1. Define each location containing Nature Run winds, water vapor, and AMV estimates as a random variable
191 x_i
- 192 2. Define θ as the population that consists of all x_i in the training dataset
- 193 3. Model the distribution of the population $P(\theta)$ as:

$$194 \quad P(\theta) = \sum_j^K \pi_j N(\mu_j, \Sigma_j) \quad (3)$$

195 Where $N(\mu_j, \Sigma_j)$ is the normal distribution with mean μ_j and covariance Σ_j of the j -th cluster,

196 K is the number of clusters, and π_j is the mixture proportion.

- 197 4. Determine π_j, μ_j, Σ_j for K clusters using the Expectation–Maximization Algorithm
- 198 5. From 3. and 4., estimate the probability of a given x_i belonging to the j-th cluster as $P(x_i \in k_j) = p_{ij}$
- 199 6. Assign point x_i to the cluster with the maximum probability p_{ij}

200 The mixture model clustering is based on the R package ‘Mclust’ developed by Fraley et al. (2012), which builds upon
201 the theoretical work of Fraley and Raftery (2002) for model-based clustering and density estimation. The process uses
202 an Expectation-Maximization algorithm to cluster the dataset, estimating a variable number of distinct multivariate
203 Gaussian distributions from a sample dataset. Training the Gaussian mixture model on this dataset provides a
204 clustering function which outputs a unique cluster for any data point with the same number of variables.

205 In one dimension, a Gaussian mixture model looks like the distributions depicted in Figure 5: instead of modeling a
206 population as a single distribution (Gaussian or otherwise), the GMM algorithm fits multiple Gaussian distributions
207 to a population. One key aspect of this algorithm is the capability of assigning a new point to the most likely
208 distribution. For example, in the 1-D figure, a normalized AMV estimate with a value of 10 would be more likely to
209 originate from the broad cluster ‘2’ than the narrow cluster ‘4’. In this case, we model the population as a Gaussian
210 mixture model in five-dimensional space, which consists of two Nature Run wind vector components (u and v), two
211 AMV estimates of these wind components (\hat{u} and \hat{v}), and the simulated water vapor values, all of which have been
212 standardized to have mean 0 and standard deviation of 1. Each cluster has a 5-dimensional mean vector for the center
213 and a 5x5 covariance matrix defining their multivariate Gaussian shape. The estimation of a covariance matrix allows
214 for the characterization of the relationships between the different dimensions within each cluster, and as such the
215 gaussian mixture model approach provides greater potential for understanding the geophysical basis of error regimes
216 than other unsupervised clustering approaches.

217 We note that the choice of inputs to the clustering methodology is limited, and that a more successful clustering may
218 be achieved by including additional meteorological or geographic information. However, the intention of this paper
219 is to study the ability of a purely data-driven approach, where no additional information or assumptions are passed to
220 the machine learning model outside of the inputs and outputs to the AMV algorithm itself. Posselt et al. (2019) showed
221 that state dependent uncertainties are a major source of error in water vapor AMVs; introducing further information
222 may cloud our ability to discern these specific uncertainties. While scaling this methodology to other applications may
223 incentivize tailoring to specific conditions, this paper aims to demonstrate that modifications are encouraged for
224 improvement, but not necessary for success.

225 Having trained the Gaussian mixture model on the 1.5 month training dataset, we applied the clustering algorithm to
226 a testing dataset sampled from the subsequent 0.5 months of the nature run. By re-analyzing the AMV estimate in
227 relation to the Nature Run winds within each cluster (**Error! Reference source not found.**), we find that the clustering
228 approach successfully separates the AMV estimates according to their ‘skillfulness’. Essentially, we repeat Figure 3

229 but divide the AMV estimates by cluster. We see that, for example, clusters 4, 5, and 7 clearly represent cases in which
230 the feature-tracking algorithm provides an accurate estimate of the Nature Run winds, with very low variance around
231 the one-to-one line (i.e., low estimation errors). Clusters 1, 2, 3, and 9 are somewhat noisier than the low-variance
232 clusters, with error characteristics similar to those of the entirety of the dataset. That is, they are considered less skilled,
233 but their estimates still lie on a one-to-one line with respect to the true wind. Clusters 6 and 8, on the other hand, are
234 clearly unskilled in different ways. Cluster 6 is a noisy regime, which captures much of the more extreme differences
235 between the AMV estimates and the Nature Run winds. Cluster 8, on the other hand, represents the low AMV estimate,
236 high Nature Run wind regime. This cluster is returning AMVs with values of zero where the Nature Run wind is
237 clearly non-zero because of the very low water vapor present. We further see the stratification of the regimes when
238 analyzing the absolute AMV error in relation to the water vapor content (Figure 7). We see that clusters that have
239 similar behaviors in the error pattern (such as 1, 2, and 3) represent different regimes of water vapor content.

240 We specified 9 individual clusters due to a combination of quantitative and qualitative reasons. Quantitatively, the
241 ‘Mclust’ package uses the Bayesian Information Criterion (BIC), a model selection criterion based on the likelihood
242 function which attempts to penalize overfitting, to select the optimal number of clusters given an input range. Using
243 an input range of one through nine, the BIC indicated the highest number of clusters would be optimal. More
244 importantly, however, the 9 clusters can be physically distinguished and interpreted. Plots of the geophysical variables
245 in the testing set associated with each of the clusters are shown in Figures 8-11. Specifically, Figure 8 plots the
246 distribution of water vapor for each cluster, while Figure 9 plots the mean wind magnitude in each direction by cluster.
247 Figure 10 plots the correlation matrix for each cluster and Figure 11 show the geographic distribution of each cluster.
248 In looking at these in combination, we see discernable and discrete clusters with unique characteristics. For example,
249 cluster 1 captures the very dry, high-wind regime in the southern hemisphere visible in Figure 2. Cluster 7
250 encompasses the tropics, while cluster 3 captures mid-latitude storm systems. Clusters 6, 8, and 9 are all characterized
251 by a much worse performance of the AMV tracking algorithm, exhibited both in Figure 7 and in Figure 8 but all
252 encompass different geographic and geophysical regimes. We see that the clustering algorithm succeeds in capturing
253 physically interpretable clusters without having any knowledge of the underlying physical dynamics. We note that in
254 other applications, the optimal number of clusters will change and the researcher will need to explore various choices
255 of this parameter in their modeling, although this tuning process should be greatly simplified by the inclusion of an
256 information criterion (e.g., BIC) in the GMM algorithm.

257 **3.5 Random Forest**

258 The clustering algorithm requires the Nature Run wind vector component values (u and v) in order to classify the
259 AMV error. When applying the algorithm in practice to tracked AMV wind from real observations, the true winds are
260 unknown. To represent the fact that we will not know the true winds in practice, we develop a proxy for the Nature
261 Run winds using only the AMV estimates and the simulated water vapor itself. This is an instance in which the
262 application of machine learning is desirable, since machine learning excels at learning high-dimensional non-linear

263 relationships from large training datasets. In this case, we specifically use random forest to create an algorithm which
264 predicts the Nature Run wind values as a function of the tracked wind values and water vapor.

265 Random forest is a machine learning regression algorithm which, as detailed by Breiman (2001), employs an ensemble
266 of decision trees to model a nonlinear relationship between a response and a set of predictors from a training dataset.
267 Here, we chose random forest specifically because it possesses certain robustness properties that are more appropriate
268 for our applications than other machine learning methods. For instance, random forest will not predict values that are
269 outside the minimum and maximum range of the input dataset, whereas other methods such as neural networks can
270 exceed the training range, sometimes considerably so. Random forest, due to the sampling procedure employed during
271 training, also tends to be robust to overtraining in addition to requiring fewer tuning parameters compared with
272 methods such as neural networks.

273 We trained a random forest with 50 trees on a separate set of tracked winds and water vapor values to predict Nature
274 Run winds using the ‘randomForest’ package in the R programming language. While the random forest estimate as a
275 whole does not perform much better than the AMV values in estimating the Nature Run wind (2.89 RMSE for random
276 forest vs 2.91 RMSE for AMVs), as shown in Figure 12, it does not display the same discrete regimentation as the
277 AMV estimates in Figure 3. As such, the random forest estimates can act as a proxy for Nature Run wind values in
278 our clustering algorithm — they remove the regimentation which is a critical distinction between the AMV estimates
279 and the Nature Run wind values.

280 **3.6 Finalized Error Characterization Model**

281 The foundation of the error characterization approach is to combine the random forest and clustering algorithm. We
282 apply the Gaussian mixture model, as trained on the Nature Run winds (in addition to the AMVs and water vapor), to
283 each point of water vapor, AMV estimate, and associated random forest estimate. This produces a set of clusters
284 which, when implemented, require no direct knowledge of the actual Nature Run state (Figure 13).

285 Naturally, the clustering algorithm performs better when applied to the dataset with the Nature Run winds, as
286 opposed to winds generated from the random forest algorithm. The former is created with direct knowledge of the
287 Nature Run winds, and any approximation will lead to increased uncertainties. In practice, the performance of the
288 cluster analysis can be improved by enhancing the performance of the random forest itself. As with any machine
289 learning algorithm, the random forest contains hyperparameters that can be optimized for specific applications. In
290 addition, performance could be improved by including additional predictor variables. Our intent is not to use the
291 random forest as a wind tracking algorithm; rather, the random forest is presented in this paper as a proof of concept.

292 Nonetheless, we see in Figure 13 and Figure 14 that the error characterization still discretizes the testing data set into
293 meaningful error regimes. The algorithm manages to separate the AMV estimates into appropriate error clusters. Once
294 again, clusters 6 and 8 manage to capture unskilled regimes, and cluster 7, and to a lesser extent clusters 4 and 5,
295 remain skillful. By taking the mean and standard deviation of the difference between AMV estimates and Nature Run

296 winds in each cluster, we develop error characteristics for each cluster (Figure 15); these quantities are precisely the
297 bias and uncertainty that we require for the cost function J in Eq (1). We see that the unskilled clusters have very high
298 standard errors and they correspond roughly to the areas of unskilled regimes in Figure 3. Similarly, skilled clusters
299 5, 4 and 7 have standard errors below that of the entire dataset. Since each cluster now has associated error
300 characteristics (e.g., bias and standard deviation), it is then straightforward to assign the bias and uncertainty for any
301 new tracked wind observation by computing which regime it is likely to belong to.

302 **3.7 Experimental Set up**

303 In this section we will describe our experimental setup for training our model on the GEOS-5 Nature Run data and
304 testing its performance on a withheld dataset. We divide the dataset into two parts: a training set consisting of the first
305 1.5 months of the GEOS-5 Nature Run, and a testing set consisting of the last 0.5 month of the Nature Run. Our
306 training/testing procedure for the simulation data and tracked wind is as follows:

- 307 1. Divide the simulation data and tracked wind into two sets: training set of 1,000,000 points from the first 1.5
308 months of the Nature Run and a testing set of 1,000,000 points from the final 0.5 months of the Nature Run.
- 309 2. We train a Gaussian Mixture Model on a normalized random sample of observations from the training dataset
310 of Nature Run winds (u and v direction), tracked winds (u and v direction), and water vapor with $n=9$ clusters.
- 311 3. We train two separate random forests on a different random sample of 750,000 observations from the training
312 dataset. We use tracked wind (u and v direction) and water vapor to model, separately, Nature Run winds in
313 both the u and v directions.
- 314 4. We apply the random forests to the dataset used for the Gaussian Mixture Model. This provides a random
315 forest estimate for each point, which is used as a substitute for Nature Run wind values in the next step.
- 316 5. We predict the Gaussian mixture component assignment for each point of water vapor, tracked winds, and
317 random forest estimate using the GMM parameters estimated in Step 2.
- 318 6. We compute the mean and standard deviation of the difference between the tracked winds and the Nature
319 Run winds, per direction, for each Gaussian mixture model cluster assignment. This provides a set of error
320 characteristics that are specific to each cluster.
- 321 7. We can apply the random forest, and then the cluster estimation, to any set of water vapor and tracked AMV
322 estimates. Thusly, any set of tracked AMV estimates and water vapor can be mapped to a specific cluster,
323 and therefore its associated error characteristics.

324 **4 Results and Validation**

325 In this section, we compare our clustering method against a simple alternative, and we quantitatively demonstrate
326 improvements that result from our error characterization. Recall that in Section 3, we divided the wind-tracking
327 outputs into 9 regimes, which range from very skilled to unskilled. For the i -th regime, we can quantify the predicted
328 uncertainty estimate as a gaussian distribution with mean m_i and standard deviation σ_i , which has a well-defined

329 cumulative distribution function which we denote as F_i . To test the performance of our uncertainty forecast, we divide
 330 the dataset described in Section 2 into a training dataset (first 1.5 month) and a testing dataset (last 0.5 month). Having
 331 trained our model using the training dataset, we apply the methodology to the testing dataset, and we compare the
 332 performance of the predicted probability distributions against the actual wind error (tracked winds - Nature Run
 333 winds). This is a type of probabilistic forecast assessment, and we assess the quality of the prediction using a scoring
 334 rule called continuous ranked probability score (CRPS), which is defined as a function of a cumulative distribution
 335 function F and an observation x as follows:

$$336 \quad \text{CRPS}(F, x) = \int_{-\infty}^{\infty} (F(x) - \mathbb{1}(y - x))^2 dy \quad (4)$$

337 Where $\mathbb{1}(\cdot)$ is the Heaviside step function and denotes a step function along the real line that is equal to 1 if the argument
 338 is positive or zero, and it is equal zero if the argument is negative (Gneiting and Katzfuss, 2014) . The continuous rank
 339 probability score here is strictly proper, which means that the function $\text{CRPS}(F, x)$ attains the minimum if the data x
 340 is drawn from the same probability distribution as the one implied by F . That is, if the data x is drawn from the
 341 probability distribution given by F , then $\text{CRPS}(F, x) < \text{CRPS}(G, x)$ for all $G \neq F$.

342 The alternative error characterization method that we test against is a simple marginal mean and marginal standard
 343 deviation of the entire tracked subtract Nature Run wind dataset. This is essentially equivalent to an error
 344 characterization scheme that utilizes one regime, where m and σ are given as the marginal mean and the marginal
 345 standard deviation of the residuals (i.e., tracked wind minus Nature Run winds). Here, we use a negatively oriented
 346 version of the CRPS (i.e., Eq.(4) without the minus sign), which implies that lower is better. A histogram evaluating
 347 the performance of our methodology against the naive error characterization method is given in Figure 16.

348 The relative behavior of the CRPS is consistent between u and v winds. The CRPS tends to have to wider distribution
 349 when applied to the regime-based error characterization. Compared to the alternative error characterization scheme,
 350 our methodology produces a cluster of highly accurate predictions (low CRPS scores), in addition to some cluster of
 351 very uninformative predictions (high CRPS scores). These clusters correspond to the highly skilled cluster (e.g.,
 352 Cluster 3) and the unskilled clusters (Cluster 6 and 8), respectively. Overall, the mean of the CRPS is lower for our
 353 methodology than it is for the alternative method, indicating that as a whole our method produces a more accurate
 354 probabilistic forecast.

355 Thus far we have shown that our method produces more accurate error-characterization than an alternative method
 356 based on marginal means and variance. Now, we assess whether our methodology provides valid probabilistic
 357 prediction; that is, we test whether the uncertainty estimates provided are consistent with the empirical distribution of
 358 the validation data. To assess this, we construct a metric in which we normalize the difference between the Nature
 359 Run wind and the tracked wind by the predicted variance. That is, for the i -th observation, we compute the normalized
 360 values for u_i and v_i using the following equations:

361

$$z_{u,i} = \frac{u_i - \hat{u}_i}{\sigma_{u,i}}$$

362

$$z_{v,i} = \frac{v_i - \hat{v}_i}{\sigma_{v,i}} \quad (5)$$

363 Where u_i is the i -th Nature Run u wind from the Nature Run data, \hat{u}_i is the tracked-wind, and $\sigma_{u,i}$ is the error as
364 assessed by our model (recall that it is a function of the regime index to which \hat{u}_i has been assigned). The values for
365 the v-wind are defined similarly. The residuals in Eq (5) can be considered as a variant of the z-score, and it is
366 straightforward to see that if our error estimates are valid (i.e., accurate), then the normalized residuals in Eq. (5)
367 should have a standard deviation of 1. If our uncertainty estimates $\sigma_{u,i}$ and $\sigma_{v,i}$ are too large, then the standard deviation
368 of $z_{u,i}$ and $z_{v,i}$ should be less than 1; similarly, if our uncertainty estimates are too small, then the standard deviation
369 of $z_{u,i}$ and $z_{v,i}$ should be larger than 1. In *Figure 17*, we display the histogram of the normalized residuals z_u and z_v .
370 It is clear that for both types of wind, the standard deviation of $z_{u,i}$ and $z_{v,i}$ are 1.003 and 1.009, respectively, indicating
371 that our error characterization model is highly accurate when forecasting uncertainties.

372 A further validation of our methods encompasses an analysis of the statistical significance of the uncertainty in our
373 model. To this end, we constructed confidence intervals for the bias and standard deviation within each regime using
374 the bootstrap (Efron and Tibshirani, 1993). The procedure of our bootstrap is as follows

- 375 1. Subset the data to retain only observations with regime index j . Let's assume that we have N_j observation
376 within this data subset
- 377 2. Sample *with replacement* N_j observations from this subset. This forms a bootstrap sample
- 378 3. From 2., compute an estimate of the bias and standard deviation.
- 379 4. Repeat step 2-3 for 1000 times, giving us 1000 estimates of the bias and 1000 estimates of the standard
380 deviation within regime j .
- 381 5. Compute 95% confidence intervals from the 1000 estimates of bias and standard deviation from 4.

382 The results for the confidence intervals (in graphical form in Figure 18). We note that the figure indicates that for
383 many of the biases, they can be considered unbiased since their confidence interval includes 0 (e.g., regimes 2-8 for
384 u-wind). However, the plot also clearly indicates that two regimes are statistically different from 0 (regime 1 and 9).
385 We also note that for the standard deviation maps, the CI's indicate that they are fairly stable (small narrow range)
386 and that most of the regimes have statistically different standard deviation (denoted here visually as CI's that do not
387 overlap one another). We also note that u and v wind direction tend to have very similar patterns, indicating that our
388 regime classification is persistent across u and v. To summarize, the CI plot above indicate that the differences in
389 standard deviation between different regimes are highly statistically significant (as evidenced by the small
390 confidence intervals and their spacing). For the biases, 3 of the regimes are statistically significantly different from
391 the rest (i.e., regimes 1, 6, and 9), while the rest are likely relatively unbiased (i.e., bias = 0).

392 5 Conclusion and Discussion

393 Error characterization is an important component of data validation and scientific analysis. For wind-tracking
394 algorithms, whose outputs (tracked u and v) are often used as observations in data assimilation analyses, it is necessary
395 to accurately characterize the bias and standard error (e.g., see Section 2.2). Nguyen et al. (2019) illustrated that
396 incorrect specification of these uncertainties (\mathbf{a} and \mathbf{R} in Eq. (1)) can adversely affect the assimilation results –
397 mischaracterization of bias will systematically offset a tracked wind, while an erroneous standard error could
398 incorrectly weigh the cost function.

399 In this paper we demonstrate the application of a machine learning uncertainty modeling framework to AMVs derived
400 from water vapor profiles intended to mimic hyper-spectral sounder retrievals. The methodology, based on a
401 combination of gaussian mixture model clustering and random forest, identified distinct geophysical regimes and
402 provided uncertainties specific to each regime. This was achieved in a purely data-driven framework; nothing was
403 known to the model except the specific inputs and outputs of the AMV algorithm, deducing the relationship between
404 regime and uncertainty from the underlying multivariate distribution of water vapor, Nature Run wind, and tracked
405 wind. Our algorithm does require one major tuning parameter in the number of clusters for the GMM algorithm,
406 although the search for the ‘optimal’ number of clusters can be aided by the inclusion of an information criterion (e.g.,
407 the BIC) in the GMM model. This implementation is not intended as a ‘ready-to-go’ algorithm for general use. Instead,
408 we lay the foundation of an uncertainty modelling approach which we plan to implement at a larger scale in subsequent
409 work. Nonetheless this bare bones implementation is sufficient to produce improved error estimates of state-dependent
410 uncertainties as detailed in Posselt et al. (2019).

411 We introduce this framework in an environment that is limited and well-behaved, but which nonetheless we believe
412 provides insight into how such an approach would perform at a larger scale. Of course, there are issues when moving
413 from the controlled environment of the simulation study to large scale applications. We understand these to be: (1) the
414 existence of uncertainty on the tracked humidity values, and (2) the ability of the training dataset to adequately capture
415 both the range of conditions of water vapor and wind speed, and their inherent relationship.

416 The simulation used for introducing this framework was a ‘perfect-observation’ environment; that is, the water
417 vapor was assumed to be perfectly known to the wind tracking algorithm. In real world scenarios, this is obviously
418 not the case. However, we believe that this is mitigated by two factors. Firstly, Posselt et al (2019) also conducted a
419 study where measurement noise was added to the water vapor measurement. This did not show to have an effect on
420 the uncertainty in the AMV estimate, except where there was the presence of strong vertical wind shear, a situation
421 which can be identified a larger scale application. Secondly, given quantified uncertainties on the water vapor
422 retrievals themselves (the scope of which is decidedly outside the work of this paper), these could be assimilated
423 into the uncertainty modelling framework in a straightforward manner by adding them as a prediction variable in
424 both the regime classification and emulator. This would allow for the model to itself ascertain the relationship

425 between water vapor uncertainty and AMV estimate uncertainty, without breaking the foundational aspect of being
426 data-driven.

427 The reliability of the training dataset is the fundamental assumption of any machine learning approach. To reiterate,
428 we present a methodology which aims to characterize the uncertainty in the difference between a measurement \hat{X}
429 and its true target X (that is, $\text{var}(\hat{X} - X)$). As such, we require some proxy for the truth in the development of our
430 model (call this X^*). To expand further, we are modelling the relationship between \hat{X} and X as a function of water
431 vapor Y , with $f(Y) = \hat{X}$ and $g(Y) = X$, where f represents the AMV algorithm and g the ‘true’ relationship
432 between wind speed and water vapor. Thus, we additionally require a proxy function g^* , which is the relationship
433 implied by the training data output of water vapor and reference winds. In the implementation presented in this
434 paper, g^* is represented by the underlying physical models that model the motion of water vapor and windspeed in
435 the GEOS-5 Nature Run.

436 The fidelity of our framework relies upon the assumption $X^* \sim X$ and $g^* \sim g$. In the simulation study, X^* is the first
437 1.5 months of a nature run simulation, which is used as a proxy for an X which consists of the last .5 months of a
438 nature run simulation. We have given the algorithm a training dataset with what we believe is a plausible range of
439 conditions which could occur in X . To the extent that errors may be seasonally and regionally dependent, it will be
440 more effective to train the error estimation algorithm on data that is expected to represent the specific flow regimes
441 and water vapor features valid for a particular forecast or assimilation period. A range of model data encompassing
442 enough seasonal variability should be a reasonable proxy for the possible range of true X . This would significantly
443 increase the computational demands of training the model (~ 1 day on a single processor, per pressure level to train
444 the current implementation of the algorithm and an average of 3 days per pressure level, on a non-optimized cluster
445 network to run the AMV extraction on the nature run), although such concerns could be mitigated by strategic
446 subsampling approaches.

447 On the other hand, in this implementation g^* is a perfectly known representation of g , which is the GEOS-5 model
448 that runs the simulation. This is where the simulation approach might create the largest source of uncertainty and
449 unreliability in the model. The true process g can only ever be approximated, and different attempts to do so will
450 involve different tradeoffs when implementing this framework. Users could, for example, use high quality validation
451 data such as matchups with radiosondes. In theory, this provides the best possible approximation of the true process
452 g , but could involve a sparsity of data such that the range of, X^* supplied is too narrow for a useful model (indeed,
453 the data might be so sparse as to— from a pure machine learning aspect— reduce the overall fidelity of the model
454 itself). On the other hand, model or reanalysis data can provide dense and diverse training datasets, but rely on the
455 assumption that the underlying physical models in those simulations are an adequate representation of the true
456 process. At the core of atmospheric models such as GOES-5 are the laws of fluid dynamics and thermodynamics. In
457 this context, water vapor is advected by the mean wind and as such the wind and water vapor are intrinsically related
458 in these models. This has been the case since the first atmospheric weather prediction models have been developed.
459 There are of course uncertainties associated with the discretization of the fluid dynamics equations, and sometimes

460 also with parameterizations depending on the physical constraints. But these uncertainties are likely small for the
461 water vapor structures that are selected for the wind tracking algorithm.

462 In both these cases, the model could likely be improved by the inclusion of additional variables in the clustering
463 algorithm. These could include a variety of parameters to address different potential problem areas in the model. As
464 mentioned previously, including quantified values of uncertainty in water vapor estimates would algorithmically
465 link the uncertainty in the humidity retrieval with the uncertainty in the AMV tracking. Similarly, including
466 parameters that correlate with geophysical phenomena where the AMV algorithm is known to perform poorly (such
467 as a marker for vertical wind shear or frontal features) would enable domain knowledge to inform the clustering
468 algorithm and emulator. Finally, it is likely that the several parameters used in formulating both the Quality
469 Indicator (Holmlund et al. 1998) and Expected Error (Le Marshall et al. 2004) approaches would be informative in
470 enhancing the algorithm. One critical aspect for users to consider is that these variables must be continuous
471 parameterizations, rather than discrete markers (which are often used in quality control); discrete variables cannot be
472 easily incorporated into a Gaussian mixture model, or indeed most clustering algorithms. Furthermore, we would
473 recommend that users implement parameters that are readily available at the same measurement location and time as
474 the AMV estimate itself. Part of the motivation for the purely state dependent approach in this framework is ease of
475 implementation; collocation and interpolation could add further uncertainty to the model.

476 We note that in real applications, using a proxy X^* instead of the true X will result in our algorithm estimating the
477 variability $\text{var}(\hat{X} - X^*)$ instead of $\text{var}(\hat{X} - X)$. Therefore, the degree to which $\text{var}(\hat{X} - X^*)$ approximates $\text{var}(\hat{X} - X)$
478 relies on the accuracy on the proxy data relative to the true uncertainty. Ultimately, implementing this
479 methodology at scale requires confidence in the training dataset employed by the user. As with most machine
480 learning approaches, a thorough understanding of the relative strengths and weaknesses of the training dataset is the
481 most critical consideration for users. This means not only ensuring that the training data is variable and diverse
482 enough to encapsulate the entirety of the true domain, but possessing some understanding of how and where
483 portions of the training dataset might be less representative of reality. There are a few practical ways in which users
484 could attempt to address this issue. Given adequate resources and time, users could train the uncertainty model under
485 various training datasets. While this would not necessarily give a greater understanding of the training data's
486 relationship with the truth, the differences between the produced models would provide some quantification of the
487 effect of the training data on the estimated uncertainties. Similarly, if users have some quantified understanding of
488 areas wherein the training dataset might be less useful (e.g., collocation errors), they could leverage this to inform
489 the uncertainty model. In this case, it is likely such decisions would manifest themselves in the final uncertainty
490 product. Nonetheless, as much as users should try to mitigate the potential for problems, there is always an
491 underlying leap of faith that they have chosen a training dataset that adequately represents the truth in their
492 application. Like any modeling approach, this methodology relies on a set of assumptions; this is one such
493 assumption. This is why domain knowledge is critical in developing a similar uncertainty model. Thoughtful and
494 careful implementations by users, keeping in mind the prescriptions and concepts detailed above, should mitigate the
495 training data dependent uncertainty.

496 Future users would also be wise to consider improvements in the random forest step of the framework. The
497 capability of this implementation in discerning accurate error regimes degrades substantially with the introduction of
498 the random forest wind estimates. This work focused on the ability to capture regime dependent error, and as such
499 the random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the
500 uncertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator;
501 ultimately, and even more so than the regime classification, these will be specific to the AMV extraction algorithm
502 being used. Certainly, many of the additional variables suggested above could be useful towards improving the
503 random forest. Users could also investigate replacing the random forest altogether with a different emulator, such as
504 a neural net or a gaussian process. Indeed, at its most general, our methodology consists of two parts: an emulator
505 and a clustering algorithm. In this implementation, random forest and Gaussian mixture modelling are the
506 approaches; in theory, these two steps could be accomplished using other algorithms belonging to the appropriate
507 class.

508 Thorough domain knowledge, both of the AMV extraction algorithm and the context in which it will be applied, is
509 critical in developing methods to improve it. As discussed previously, the bare bones implementation of our
510 methodology in this paper is intended as a structural presentation of the conceptual framework, not necessarily a
511 finalized model. However, it is also the case that the investigation by Posselt et. al (2019) showed that the variables
512 used in this implementation of the model are those most strongly related with AMV uncertainty in this particular
513 application. The state-dependent errors identified by Posselt et al. (2019) are also expected to apply to other water
514 vapor AMVs. This is because, in general, AMV algorithms have difficulty tracking fields with very small gradients,
515 and will produce systematic errors in situations for which isolines in the tracked field (e.g., contours of constant water
516 vapor mixing ratio) lie parallel to the flow. To the extent that our algorithm represents a general class of errors, the
517 results may be applicable to other geophysical scenarios and other AMV tracking methodologies. As mentioned in the
518 introduction, robust estimates of uncertainty are important for data assimilation, and we expect that our methodology
519 could be used to provide more accurate uncertainties for AMVs used in data assimilation for weather forecasting and
520 reanalysis.

521 **Author Contribution**

522 Teixeira conceived of the idea with inputs from Nguyen. Teixeira performed the computation. Wu provided the
523 experimental datasets along with data curation expertise. Posselt and Su provided subject matter expertise. All authors
524 discussed the results. Teixeira wrote the initial manuscript and updated the draft with inputs from co-authors.

525 **Competing Interest:** The Authors declare no conflict of interest.

526 **Funding Acknowledgment:** The research was carried out at the Jet Propulsion Laboratory, California Institute of
527 Technology, under a contract with the National Aeronautics and Space Administration (80NM0018D0004). © 2020.
528 California Institute of Technology. Government sponsorship acknowledged

529 **References**

- 530 Bormann, N., Hernandez-Carrascal, A., Borde, R., Lutz, H.J., Otkin, J.A. and Wanzong, S.: Atmospheric motion
531 vectors from model simulations. Part I: Methods and characterization as single-level estimates of wind, *Journal of*
532 *Applied Meteorology and Climatology*, 53(1), 47-64. <https://doi.org/10.1175/JAMC-D-12-0336.1>, 2014.
- 533 Breiman, L.: Random forests. *Machine learning*, 45(1), 5-32, 2001.
- 534 Cassola, F. and Burlando, M.: Wind speed and wind energy forecast through Kalman filtering of Numerical Weather
535 Prediction model output, *Applied Energy*, 99, 154-166, 2012.
- 536 Coulston, J.W., Blinn, C.E., Thomas, V.A. and Wynne, R.H., 2016. Approximating prediction uncertainty for
537 random forest regression models. *Photogrammetric Engineering & Remote Sensing*, 82(3), pp.189-197.
- 538 Tibshirani, R.J. and Efron, B., 1993. An Introduction to the Bootstrap. *Monographs on statistics and applied*
539 *probability*, 57, pp.1-436.
- 541 Fraley, C. and Raftery, A.E.: MCLUST: Software for model-based clustering, density estimation and discriminant
542 analysis (No. TR-415). Washington University, Seattle Department of Statistics, 2002.
- 543 Fraley, C., Raftery, A.E., Murphy, T.B. and Scrucca, L: mclust version 4 for R: normal mixture modeling for model-
544 based clustering, classification, and density estimation, Washington University, Seattle Department of Statistics, 2012
- 545 Gneiting, T. and Katzfuss, M., 2014. Probabilistic forecasting. *Annual Review of Statistics and Its Application*, 1,
546 pp.125-151
- 547 Hernandez-Carrascal, A. and Bormann, N.: Atmospheric motion vectors from model simulations. Part II:
548 Interpretation as spatial and vertical averages of wind and role of clouds, *Journal of Applied Meteorology and*
549 *Climatology*, 53(1), 65-82, 2014.
- 550 Holmlund, K., Velden, C. S., & Rohn, M.: Enhanced automated quality control applied to high-density satellite-
551 derived winds, *Monthly Weather Review*, 129(3), 517-529, 2001.
- 552 Kawa, S.R., Erickson, D.J., Pawson, S. and Zhu, Z.: Global CO2 transport simulations using meteorological data
553 from the NASA data assimilation system, *Journal of Geophysical Research: Atmospheres*, 109,
554 <https://doi.org/10.1029/2004JD004554>, 2004.
555
- 556 Kwon, Y., Won, J.H., Kim, B.J. and Paik, M.C., 2020. Uncertainty quantification using Bayesian neural networks in
557 classification: Application to biomedical image segmentation. *Computational Statistics & Data Analysis*, 142,
558 p.106816.
- 559 Le Marshall, J., Rea, A., Leslie, L., Seecamp, R., & Dunn, M.: Error characterisation of atmospheric motion vectors,
560 *Australian Meteorological Magazine*, 53(2), 2004.
- 561 Mueller, K.J., Wu, D.L., Horváth, Á., Jovanovic, V.M., Muller, J.P., Di Girolamo, L., Garay, M.J., Diner, D.J.,
562 Moroney, C.M. and Wanzong, S.: Assessment of MISR cloud motion vectors (CMVs) relative to GOES and
563 MODIS atmospheric motion vectors (AMVs), *Journal of Applied Meteorology and Climatology*, 56(3), 555-572,
564 <https://doi.org/10.1175/JAMC-D-16-0112.1>, 2017.
565
- 566 Nguyen, Hai, Noel Cressie, and Jonathan Hobbs. "Sensitivity of Optimal Estimation Satellite Retrievals to
567 Misspecification of the Prior Mean and Covariance, with Application to OCO-2 Retrievals." *Remote Sensing* 11.23
568 (2019): 2770.

569 Posselt, D. J., L. Wu, K. Mueller, L. Huang, F. W. Irion, S. Brown, H. Su, D. , and C. S. Velden: Quantitative
570 Assessment of State-Dependent Atmospheric Motion Vector Uncertainties. *J. Appl. Meteor. Clim.*, In Press.
571 <https://doi.org/10.1175/JAMC-D-19-0166.1>, 2019.

572 Putman, W., A.M. da Silva, L.E. Ott and A. Darmanov: Model Configuration for the 7-km GEOS-5 Nature Run,
573 Ganymed Release (Non-hydrostatic 7 km Global Mesoscale Simulation). GMAO Office Note No.5 (Version 1.0),
574 18, 2014.

575 Salonen, K., J. Cotton, N. Bormann, and M. Forsythe: Characterizing AMV Height-Assignment Error by Comparing
576 Best-Fit Pressure Statistics from the Met Office and ECMWF Data Assimilation Systems, *J. Appl. Meteor.*
577 *Climatol.*, 54, 225–242, <https://doi.org/10.1175/JAMC-D-14-0025.1>, 2015.

578 Staffell, I. and Pfenninger, S.: Using bias-corrected reanalysis to simulate current and future wind power output,
579 *Energy*, 114,1224-1239, 2016.

580 Swail, V.R. and Cox, A.T.: On the use of NCEP–NCAR reanalysis surface marine wind fields for a long-term North
581 Atlantic wave hindcast, *Journal of Atmospheric and oceanic technology*, 17(4), 532-545, 2000.

582 Tran, D., Dusenberry, M., van der Wilk, M. and Hafner, D., 2019. Bayesian layers: A module for neural network
583 uncertainty. In *Advances in Neural Information Processing Systems* (pp. 14660-14672).
584

585 Tripathy, R.K. and Bilonis, I., 2018. Deep UQ: Learning deep neural network surrogate models for high
586 dimensional uncertainty quantification. *Journal of computational physics*, 375, pp.565-588.

587 Velden, C.S. and K.M. Bedka,: Identifying the Uncertainty in Determining Satellite-Derived Atmospheric Motion
588 Vector Height Attribution. *J. Appl. Meteor. Climatol.*, 48, 450–463, <https://doi.org/10.1175/2008JAMC1957.1>, 2009.

589 Zeng, X., S. Ackerman, R.D. Ferraro, T.J. Lee, J.J. Murray, S. Pawson, C. Reynolds, and J. Teixeira:
590 Challenges and opportunities in NASA weather research. *Bull. Amer. Meteor. Soc.*, 97, 137–140, 2016.

591

592

593

594

595

596

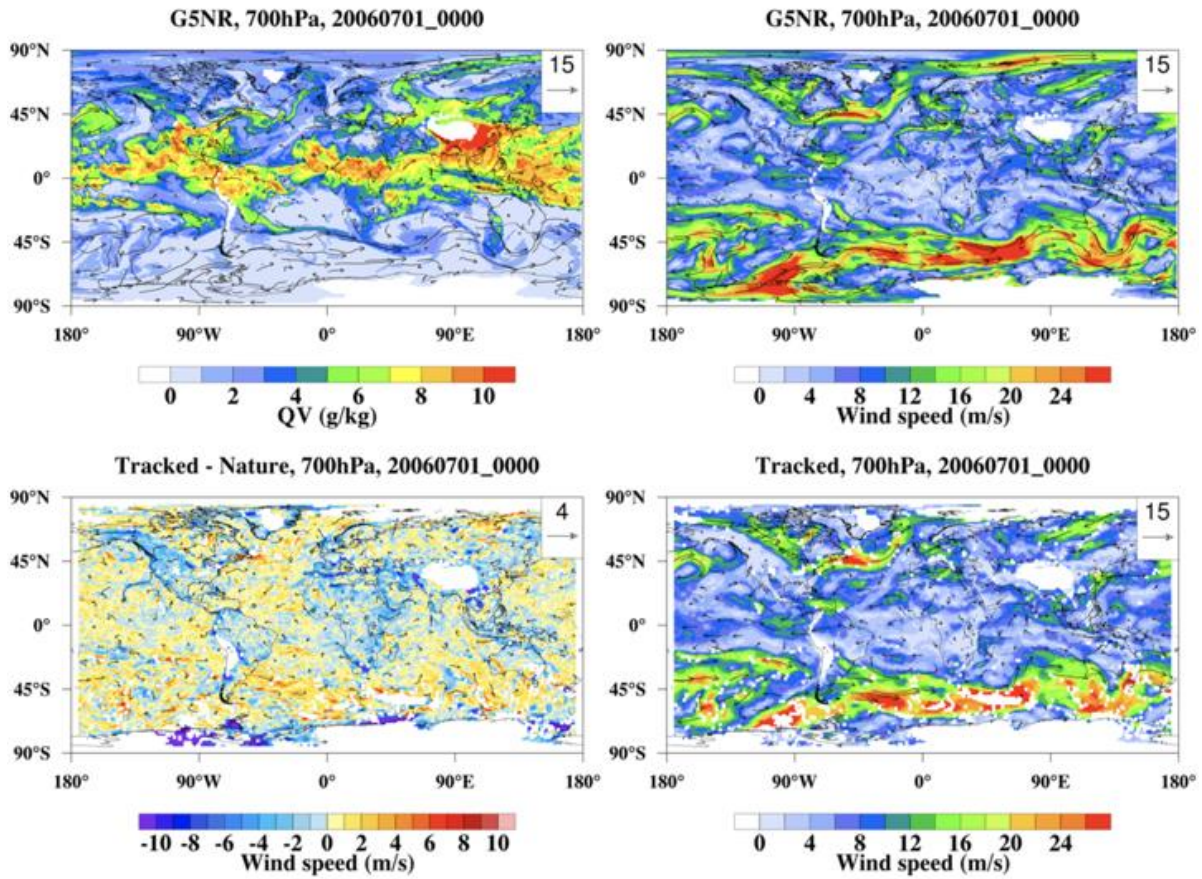
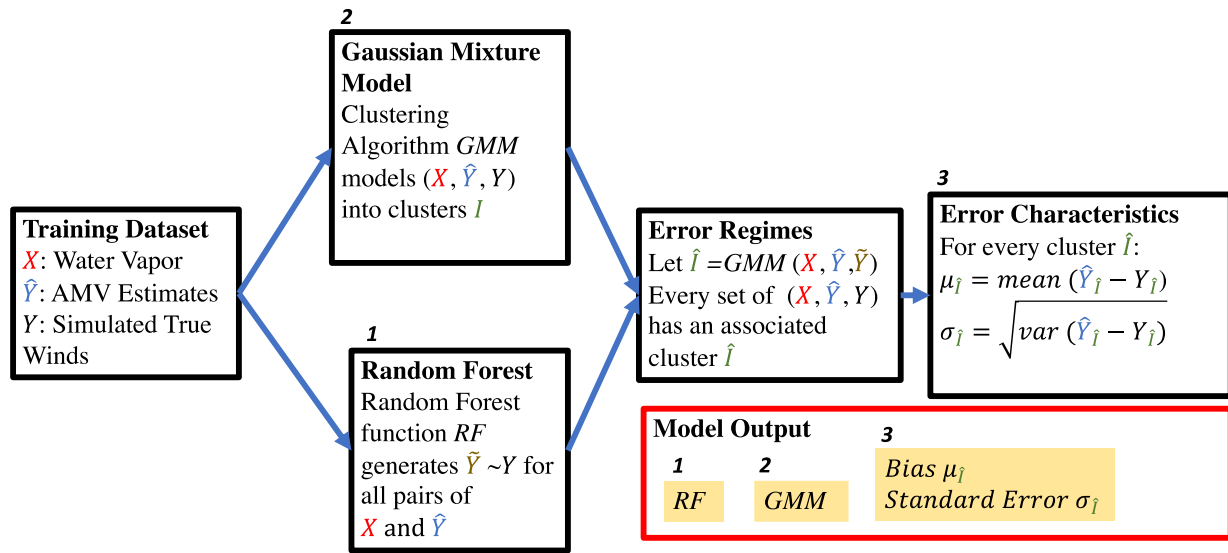
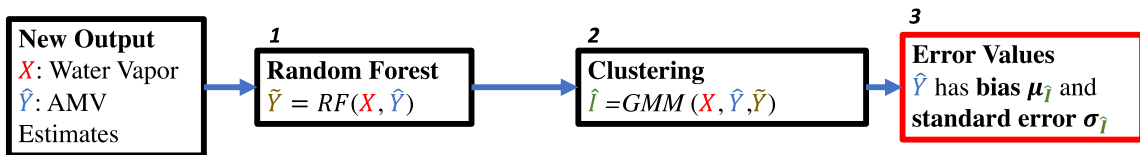


Figure 1: Map of Nature Run at one timestep at 700hPa (A): Water Vapor (B): Nature Run Wind Speed (C): Difference between Nature Run Wind Speed and AMV Estimate (D): AMV Estimate.

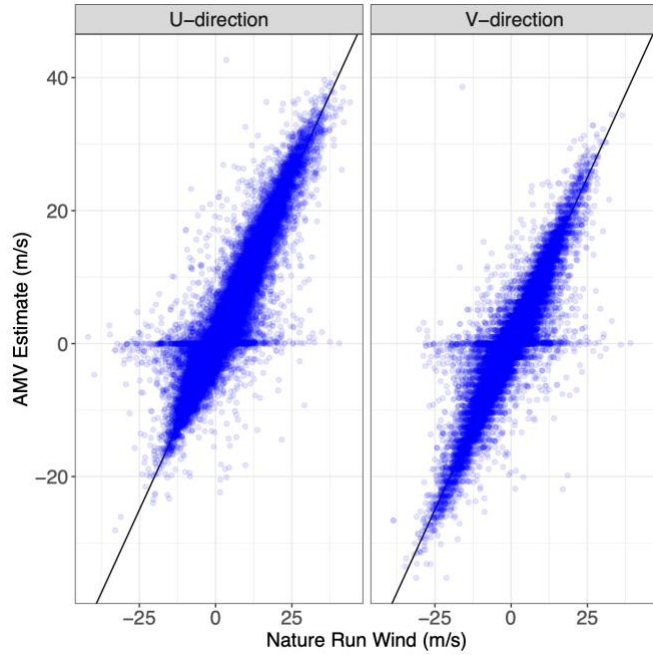
1. Training



2. Implementation

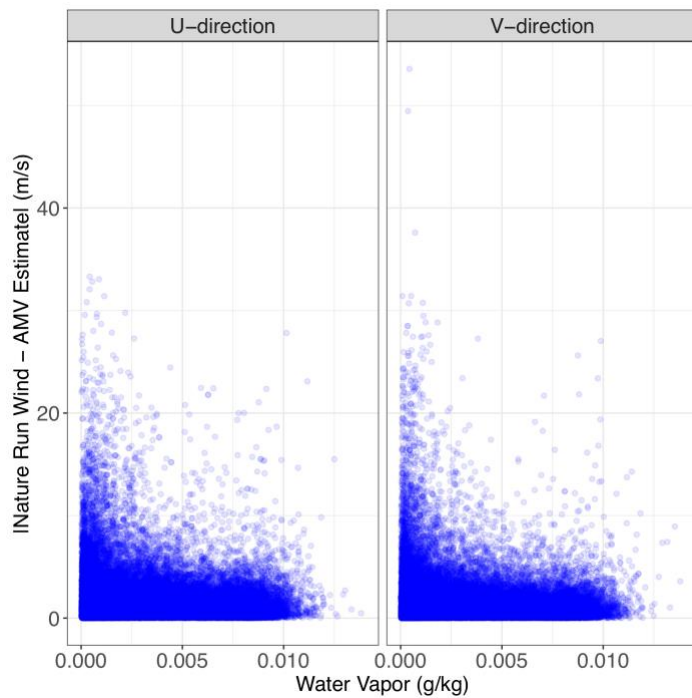


600 Figure 2: Diagram of Training Approach and Diagram of Implementation steps.



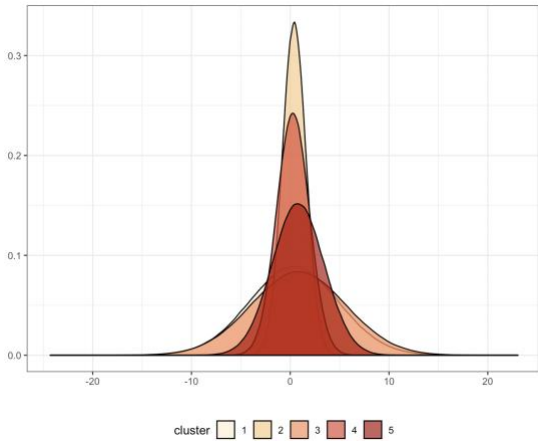
601

602 **Figure 3: Scatter plot of the simulated Nature Run wind vs AMV estimates for u and v wind in the training**
 603 **dataset.**



604

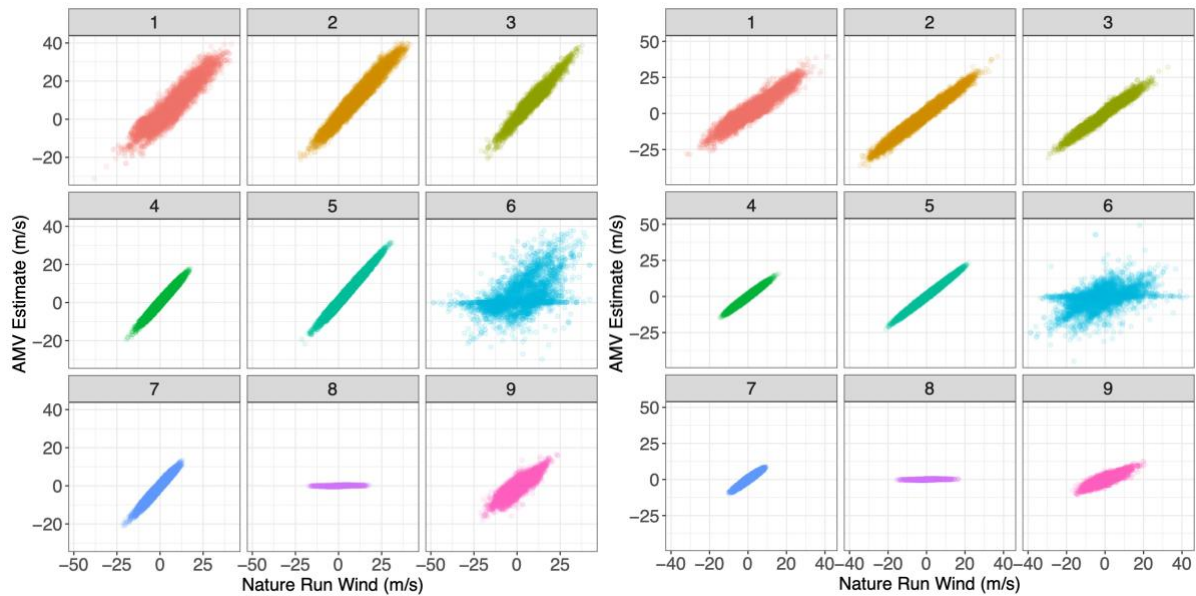
605 **Figure 4: Simulated water vapor vs the absolute value of the difference between Nature Run and tracked**
 606 **winds in the training dataset.**



607

608 **Figure 5: Example of Gaussian Mixture Model in one dimension. Density Figures for the U-Direction AMV**
 609 **Estimate dimension of fitted Gaussian mixture.**

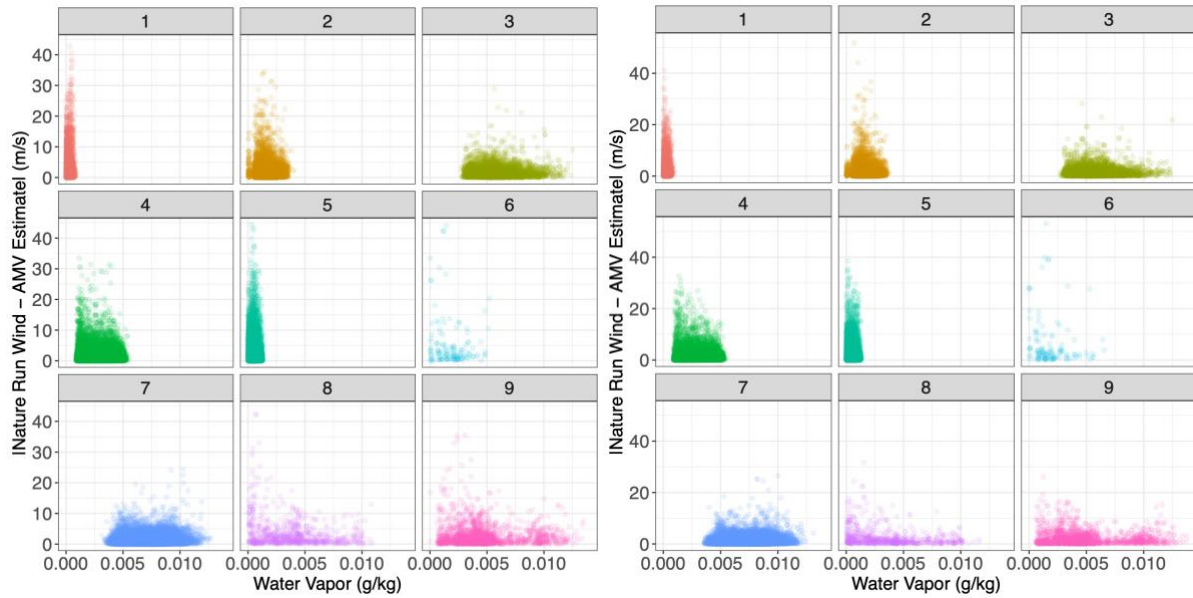
610



611

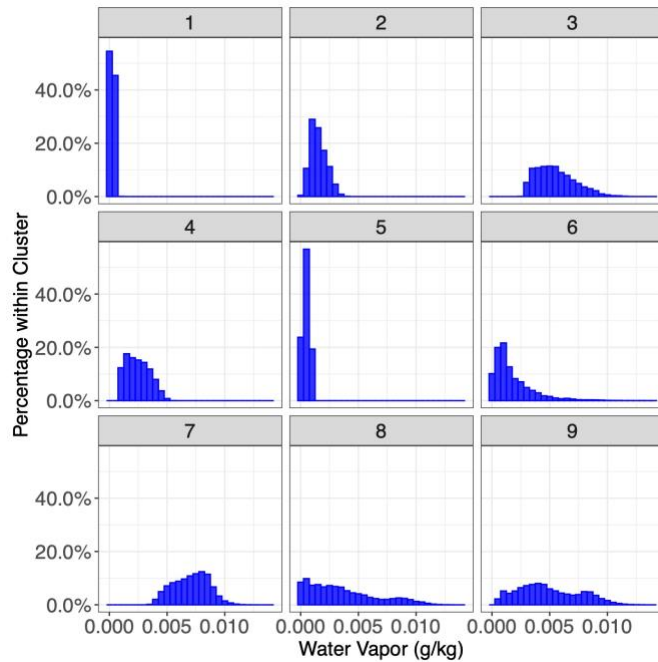
612 **Figure 6: Scatterplot of simulated Nature Run wind vs AMV Estimates, each sub-panel corresponding to the**
 613 **specific Gaussian mixture component to which each point in the testing set has been assigned. (A): U-**
 614 **Direction Wind (B): V-Direction Wind.**

615



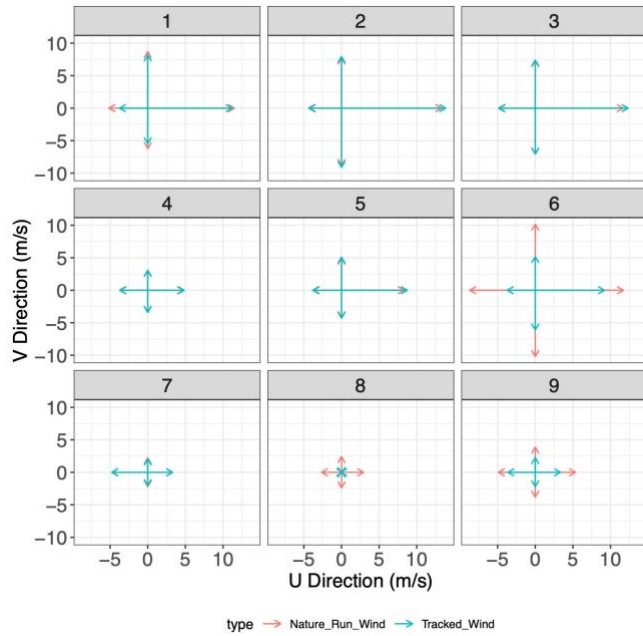
616

617 **Figure 7: Scatterplot of Water Vapor vs Absolute Tracked Wind Error, each sub-panel corresponding to the**
618 **specific Gaussian mixture component to which each point in the testing set has been assigned. (A): U-**
619 **Direction Wind (B): V-Direction Wind.**



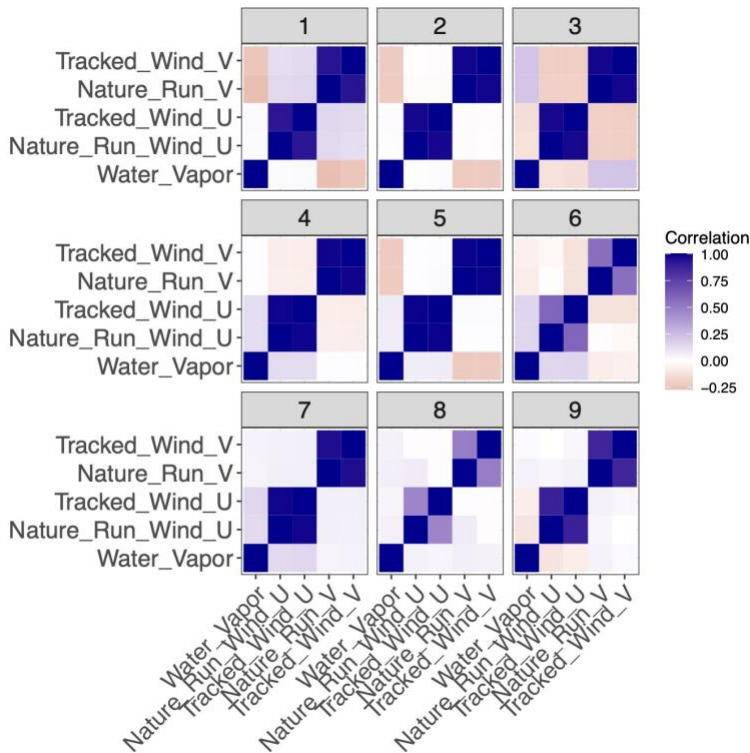
620

621 **Figure 8: Histogram of Nature Run water vapor for each cluster identified by the Gaussian mixture model,**
622 **applied to the testing set. Each sub-panel represents the cluster each point was assigned to.**



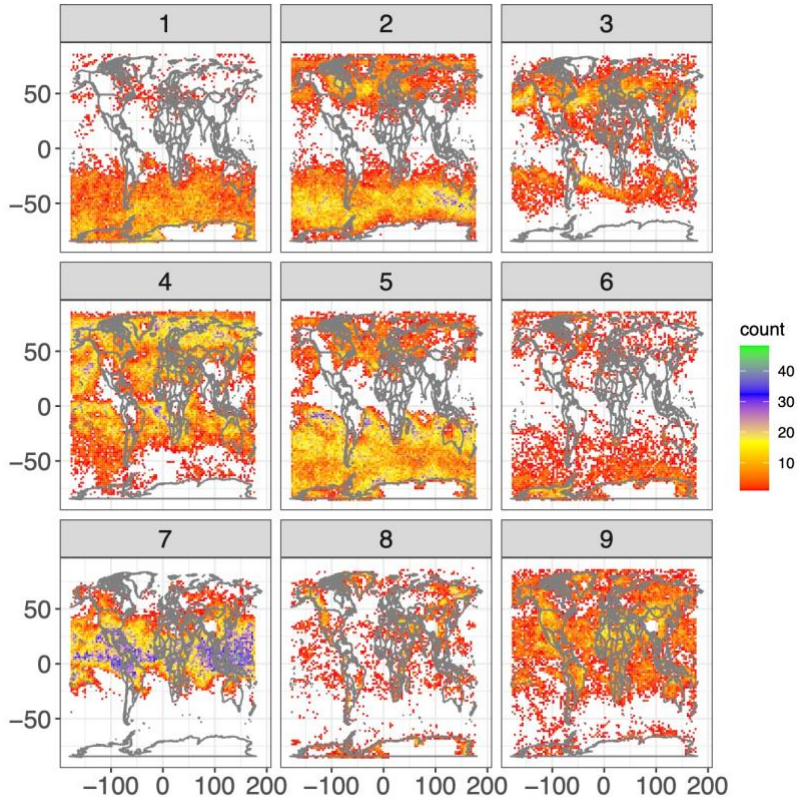
623
624
625
626

Figure 9: Mean tracked winds and Nature Run winds, in each direction, for each cluster applied to the test set. Each sub-panel represents the cluster each point was assigned to.



627
628
629

Figure 10: Correlation matrix between each clustered element for each identified cluster in the original training dataset. Each sub-panel refers to a specific cluster.

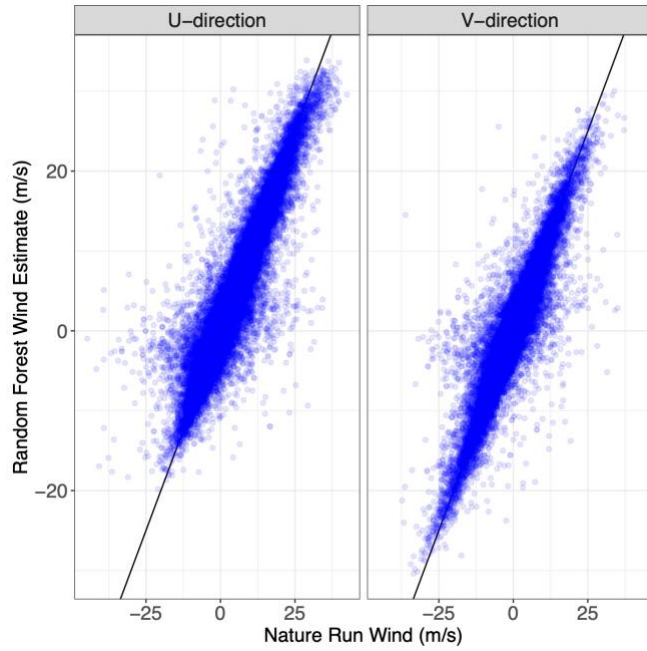


630
631
632

Figure 11: Geographic distribution by cluster of AMV retrieval locations in the testing dataset. Each sub-panel represents one cluster.

633

634

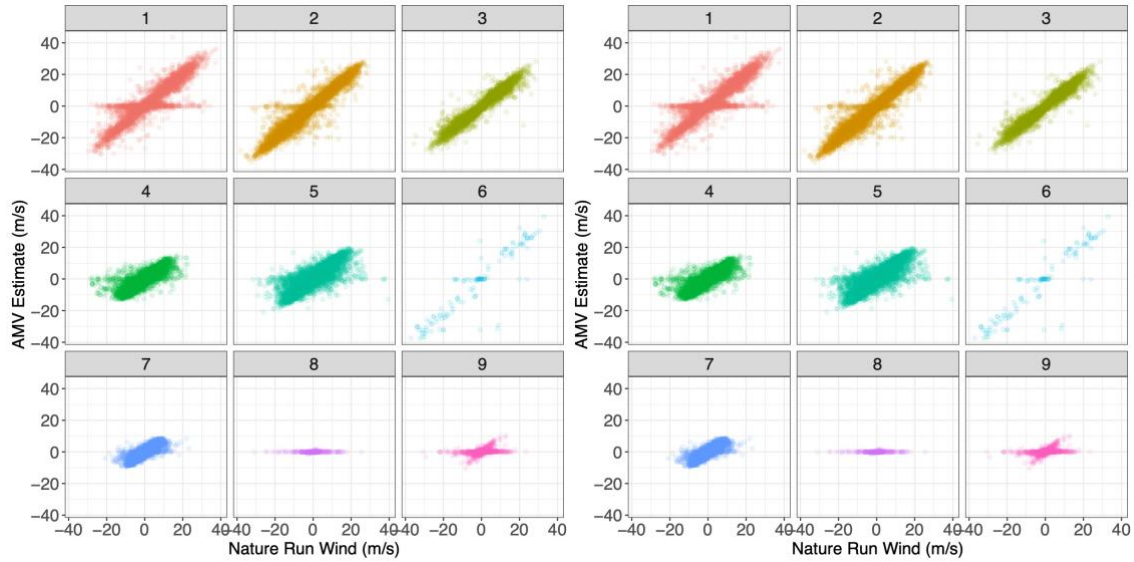


635

636 **Figure 12: Scatterplot of Nature Run wind estimate vs random forest produced estimate. (A): U Direction**

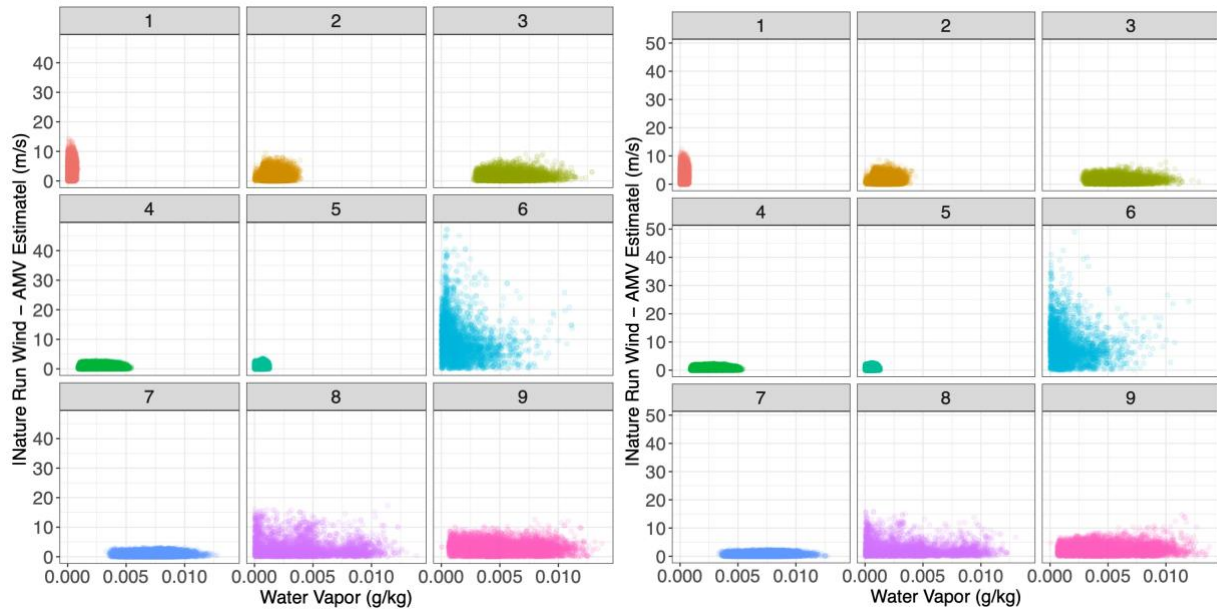
637 **(B): V Direction**

638



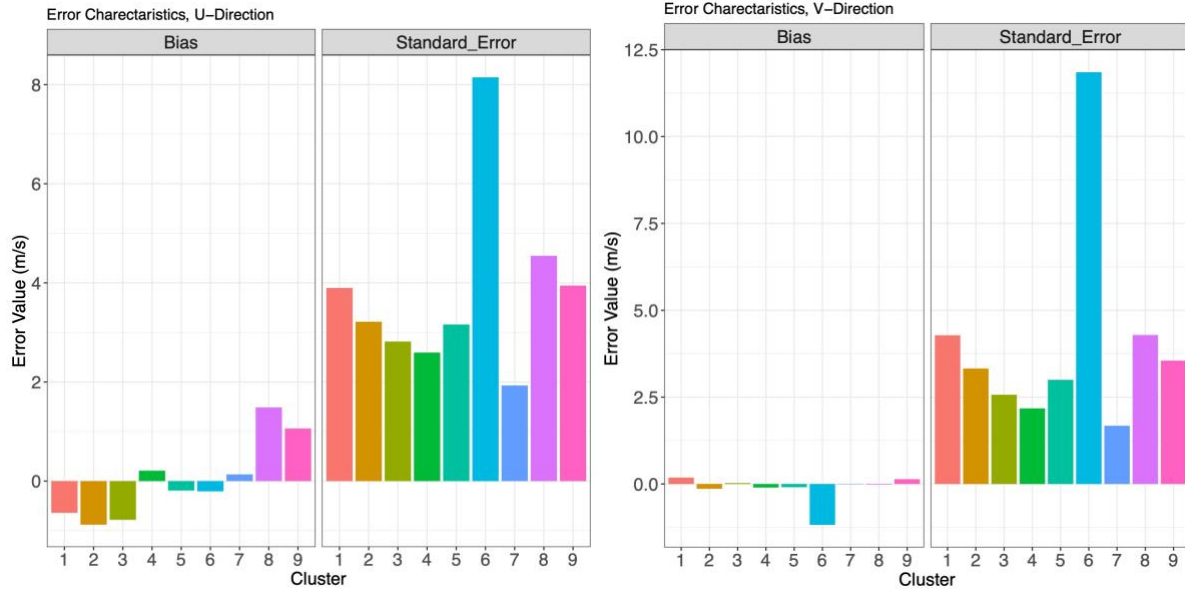
639

640 **Figure 13: Scatterplot of Nature Run wind vs AMV Estimates, each sub-panel corresponding to the specific**
 641 **Gaussian mixture component to which each point in the testing set has been assigned when the Nature Run**
 642 **wind value has been substituted by the random estimate. (A): U-Direction Wind (B): V-Direction Wind**



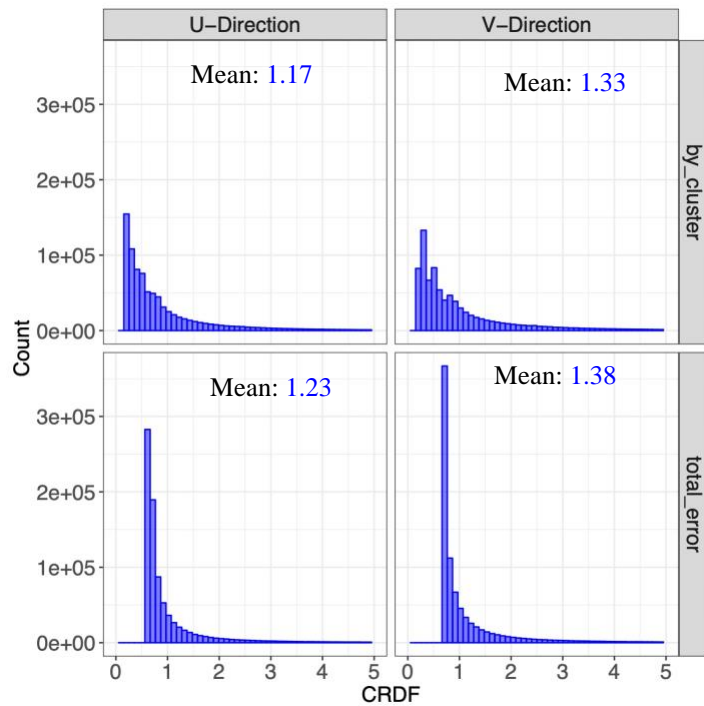
643

644 **Figure 14: Water Vapor vs Absolute Tracked Wind Error, each sub-panel corresponding to the specific**
 645 **Gaussian mixture component each point in the testing set has been assigned when the Nature Run wind**
 646 **value has been substituted by the random estimate. (A): U-Direction Wind (B): V-Direction Wind**



647

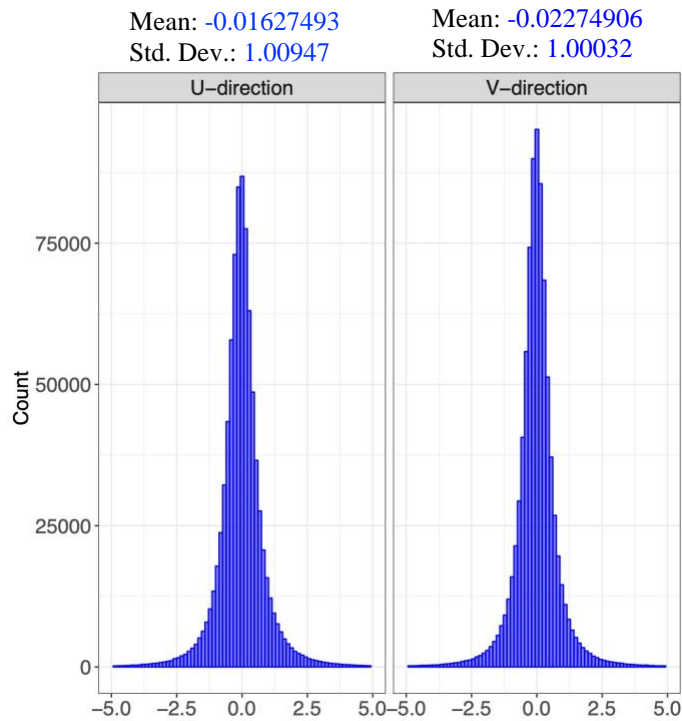
648 **Figure 15: (A): Bias (Left Panel) and Standard Error (Right Panel) for each Gaussian mixture cluster in**
 649 **figure 6, U direction. (B): Same as (A) for V-direction**



650

651 **Figure 16: CRSP applied to different error approaches. (A): Cluster Errors for U Winds (B): Total Errors**
 652 **for U Winds (C): Cluster Errors for V Winds (D): Total Errors for V Winds.**

653



654

$$z_u = \frac{u - \hat{u}}{\sigma_u}$$

$$z_v = \frac{v - \hat{v}}{\sigma_v}$$

655

656 **Figure 17: U and V winds normalized using the error characteristics developed by our methodology.**

657

658

659

660

661

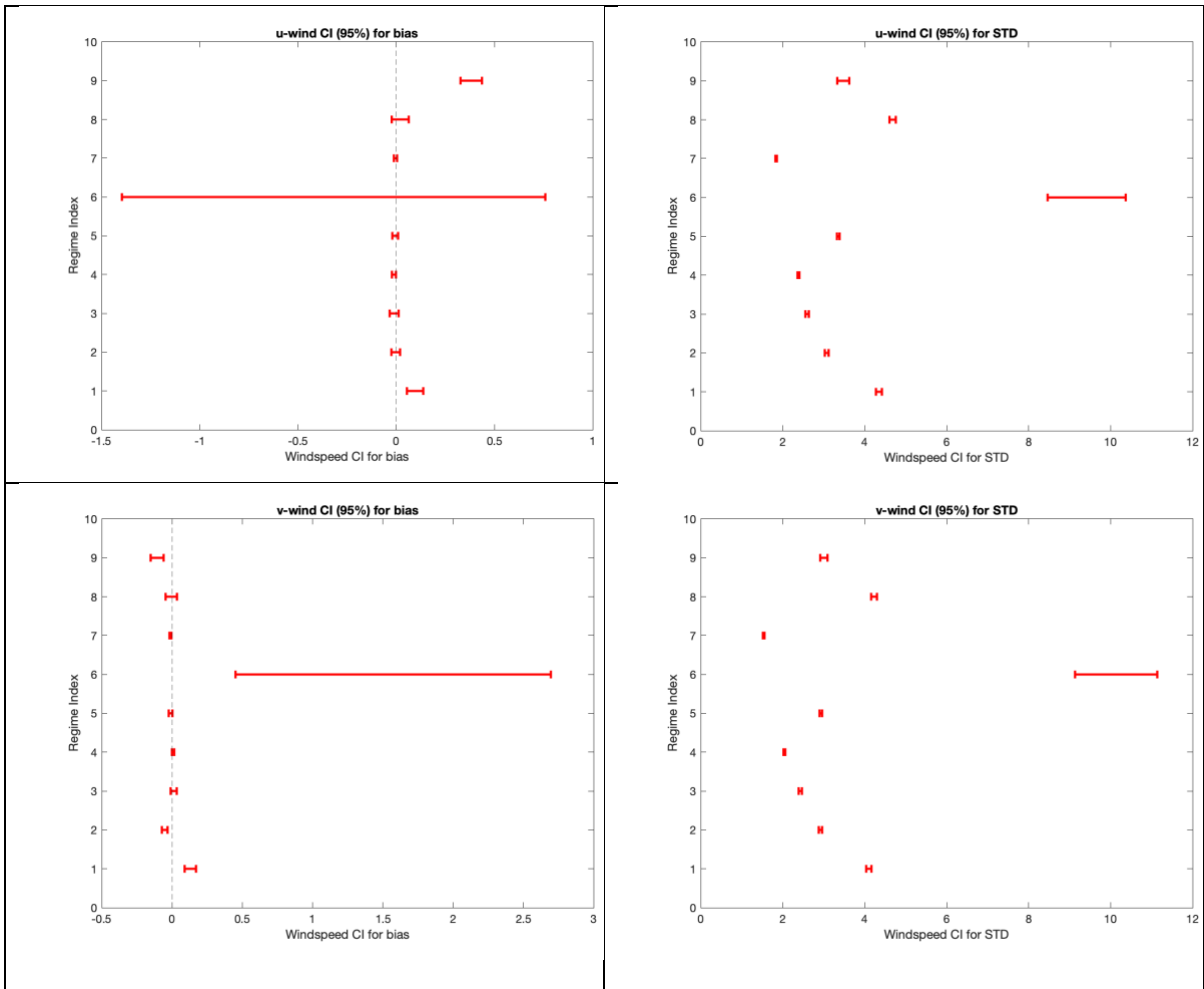
662

663

664

665

666



667 Figure 18: Top rows (bias and std confidence intervals for u-wind), bottom rows (bias and std confidence
 668 intervals for v-winds). The interval represent a 95% confidence interval.

669

670

671

672

673