"Using Machine Learning to Model Uncertainty for Water-Vapor Atmospheric

Motion Vectors"

Teixeira et al.

Responses to Referee 1

We would like to thank the referee for the careful read of the paper and comments. Please see our responses below:

GENERAL COMMENTS

The second version of this paper is greatly improved. I would like to thank the authors for the substantial effort they made to answer referees concerns.

I especially appreciated the addition of the Fig8-11 that illustrate the physical and geographical properties of the identified clusters. The general presentation of the method as a 'proof of concept' also improved a lot the overall coherency of the paper with the actual results.

I can now accept this paper for publication.

We thank the reviewer for their detailed and helpful comments on the first review of the paper. We are glad to see that we have met the reviewer's acceptance for publication.

SPECIFIC COMMENTS

1) Line 228 'to the Nature Run winds within each cluster (),' Probably missing something inside the parenthesis.

Thannk you for noticing this. We have fixed it by including the figure reference needed.

2) Line 274-275

I do not understand the sentence: '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 R.' What is 'R'?

'R' refers to the R programming language, a language used for statistical analysis. This has been noted in the paper.

3) Line 380-381

'In this paper we demonstrate the application of a machine learning uncertainty modeling tool to AMVs derived from hyperspectral sounder water vapor profiles.' This is not exactly correct. No data from hyperspectral sounders have been used in this study. Please precise or correct.

Thank you for this comment. This has been corrected to 'water vapor profiles intended to mimic hyper-spectral sounder retrievals.'

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

Teixeira et al.

Responses to Referee 2

We would like to thank the referee for the careful read of the paper and for the detailed comments. Please see our responses below:

1. While the manuscript has improved in several areas, I feel that some of my earlier more substantive comments have not been sufficiently addressed. These include comments on the usefulness of the algorithm in practical situations (earlier general points 1 and partially 3), the robustness of the performance (particularly of the clustering; general point 1, specific points 6 and 10), and the evaluation of the performance (specific points 16 and 17). In some of their replies to the comments the authors acknowledge limitations and issues of their algorithm, yet this is not reflected in the revised manuscript (especially general point 3, specific points 6 and 10). I therefore continue to be unconvinced that this algorithm is useful for the intended application.

For the paper to be acceptable for publication, I think these points need to be more thoroughly addressed. In my view this means that the results are either more critically evaluated and the limitations more clearly outlined, or the algorithm has to be refined and made more robust so that applicability to real-life applications is indeed ensured. In my view either of this requires another major revision. I do not think that the present results serve as a "proof-of-concept" of a useful algorithm that could be applied to practical situations with real data, as the authors claim. I could accept if the paper was written from the perspective of introducing a novel conceptual framework which has been explored in an initial implementation, which shows some skill in assigning regime-specific observation errors, but has revealed a range of issues that would need to be addressed for this to be a viable and useful algorithm for real applications (and which may not be possible to address).

We thank the reviewer for their insightful comments. We agree with the reviewer that what we intend to present is more of a 'conceptual framework' than a 'proof-of-concept'. We never intended to present this work as a 'ready-to-go' algorithm in this particular implementation; instead, we laid out the foundation for an uncertainty modelling approach which we plan to implement at a larger scale in subsequent work. To summarize, the conceptual foundations to this framework are the following:

1. We intend to model uncertainty in the AMV algorithm relative to the underlying value it is trying to capture. As described in previous responses and detailed in the paper, this is a departure from most machine learning based approaches to uncertainty modelling.

- 2. These uncertainties, and further their relationship with state vector elements, have shown to discretize themselves into different error regimes. We focus our approach on characterizing these regimes.
- 3. Following the work by Posselt et al (2019), we believe that these uncertainties are state-dependent. As such, our framework explicitly intends to examine these uncertainties as function of the relationship between state values. This is not to say that implementations of this framework should exclude additional information; on the contrary, we believe that the addition of context-dependent information could greatly enhance an implementation of this framework. However, at its core, we attempt to model state-dependent uncertainties with a state-dependent model.
- 4. We believe that the most versatile framework, in terms of potential application, is one which at its base is context-agnostic. A purely data-driven approach, as we lay out in the paper, provides the platform for context-dependent tuning when scaling the methodology.
- 2. A clear path of how these issues could be addressed should be provided, including an outline how the algorithm could be derived without the truth being available . The latter requires recognition and discussion of the issues that will be encountered when substituting the truth with proxy data beyond the superficial suggestion that has been added in the conclusion section. These proxy data (the authors suggest reanalysis data or collocated reference observations) will introduce their own regime-specific errors, and there is no mechanism apparent in the present implementation of the algorithm that would be able to separate the errors in the proxy data (or the collocation errors), rather than in the AMVs. In addition, real-life applications would encounter errors in the humidity retrievals used for the tracking, and these will affect the characteristics of the errors in the tracked winds as well as the performance of the algorithm.

We introduce our framework in an environment that is limited and well-behaved, but which nonetheless we believe provides insight into how such an approach would perform at a larger scale. The uncertainty model shows skill in discerning regime specific uncertainties in this scenario. Of course, there are issues when moving from the controlled environment of the simulation study to large scale applications. We understand these to be: (1) the existence of uncertainty on the tracked humidity values, and (2) the ability of the training dataset to adequately capture both the range of conditions of water vapor and wind speed, and their inherent relationship. We have rewritten and expanded the conclusion in the paper to include a lengthier discussion of these topics, as well as recommendations on how to address them. As a reference, this includes:

- 1. A discussion of how humidity retrieval uncertainty would impact our model and approaches to address it. L417-431.
- A discussion of considerations future users must make when developing training data for the model. L432-466; L481-486.
- A list of specific prescriptions and recommendations for improving the regime classification approach. L467-480.

4. A discussion of the random forest emulator and suggestions for future users. L487-498.

Some additional specific points:

1) Abstract: "... that is purely data driven and requires few tuning parameters": I think this statement reflects more a limitation of the present study than a feature of the conceptual algorithm. There are significant tuning parameters (e.g., number of clusters, choice of predictors), but they are not explored in this study. I would argue that the method still requires substantial tuning of these parameters to be robust and useful compared to other methods, and some of the responses of the authors appear to agree with this. So I suggest to rephrase this statement to avoid giving the impression that this is a positive feature of the algorithm and that it can be run "out of the box".

We thank the reviewer for this comment. In order to resolve any confusion, and in light of the fact that we have highlighted several parts of the algorithm for future users to adjust and improve, we have removed any reference to 'few' tuning parameters.

2) Abstract, last sentence "... and it is shown to adequately capture the error features of the tracked wind.": I don't think it is clear what "adequate" means in this context, and I think it would be useful to be frank that the truth was available for training. I suggest to rephrase "..., and it is shown to capture some of the regime-dependent error features of the tracked wind in a setting where the truth is available for training the algorithm."

We thank the reviewer for this comment. We slightly disagree with the idea that the 'truth' was available for training. As discussed in the revised conclusions (L441-451), we trained the model on the first 1.5 months of the Nature Run data, and presented results applied on the last .5 months of the Nature Run; while this is a situation where we expect the training data to fairly accurately reflect the domain of the testing data, it is not in fact the same data and thus not what we would deem the 'truth'. Nevertheless, the term 'adequate' is inherently subjective and we have replaced this with 'provide accurate and useful error features of the tracked wind', which is shown in the analysis in section 4.

3) Reply to my earlier point 11: Thanks for providing the information on the significance of the standard deviation and bias statistics for the clusters included in the reply. This should be included in the revised manuscript, as it would help to substantiate the claim that the clusters identify statistically significantly different error behaviours.

Thank you for this comment. We have included our reply to point 11 in the paper at the end of section 4, L373-392 (as well as Figure 18)

List of relevant changes in manuscript (in order as they appear):

- 1. Abstract and Title:
 - a. Rewording of final sentence of abstract.
- 2. Section 4:
 - a. L: 376-396. Included a discussion on the statistical significance of the bias and standard errors values derived from the uncertainty model.
- 3. Conclusion and Discussion:
 - a. Re-written and greatly expanded conclusion section to include discussion of potential problem areas for those wishing to implement this methodology. Including, but not limited to, a discussion of:
 - i. A discussion of how humidity retrieval uncertainty would impact our model and approaches to address it. L417-431.
 - ii. A discussion of considerations future users must make when developing training data for the model. L432-466; L481-486.
 - iii. A list of specific prescriptions and recommendations for improving the regime classification approach. L467-480.
 - iv. A discussion of the random forest emulator and suggestions for future users. L487-498.
- 4. References:
 - a. Included reference to (Efron and Tibshirani, 1993).
- 5. Figures:
 - a. New Figure 18, which accompanies the new additions to Section 4.

Using Machine Learning to Model Uncertainty for Water-Vapor 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 Deleted: and requires few tuning parameters

Deleted: adequately capture the error features of the tracked win...

37 methods (specifically, optimal estimation retrievals), then the posterior estimates would also be biased. Moreover, 38 they proved that when the input standard error is 'correct' (that is, it is as close to the unknown truth as possible), then 39 the resulting Bayesian estimate is 'efficient'; that is, it has the smallest error among all possible choices of prior 40 standard error. Additionally, multiple active and passive technologies are being developed to measure 3D winds, such 41 as Doppler wind lidar (DWL), radar, and infrared/microwave sensors that derive AMVs using feature-tracking of 42 consecutive images. Therefore, an accurate and robust methodology for modeling uncertainty will allow for more 43 accurate assessments of mission impacts, and the eventual propagation of data uncertainties for these instruments.

44 Velden and Bedka (2009) and Salonen et al. (2015) have shown that height assignment contributes a large component 45 of uncertainty in AMVs tracked from cloud movement and from sequences of infrared satellite radiance images. 46 However, with AMVs obtained from water vapor profiling instruments (e.g., infrared and microwave sounders), 47 height assignment error cannot be directly assessed purely through analysis of the AMV extraction algorithm. Height 48 assignment is instead an uncertainty in the water vapor profile itself. Unfortunately, without the quantified 49 uncertainties on the water vapor profile necessary to pursue such a study, that is well beyond the scope of this paper. 50 As such, this study will focus on errors in the AMV estimates at a given height. Previous work has demonstrated 51 several different approaches for characterizing AMV vector error. One common approach is to employ quality 52 indicator thresholds, as described by Holmund et al (2001), which compare changes in AMV estimates between 53 sequential timesteps and neighboring pixels, as well as differences from model predictions, to produce a quality 54 indicator to which a discrete uncertainty is assigned. The Expected Error approach, developed by Le Marshal et al. 55 (2004), builds a statistical model using linear regression against AMV-radiosonde values to estimate the statistical 56 characteristics of AMV observation error.

57 In this study, we outline a data-driven approach for building an AMV uncertainty model using observing system 58 simulation experiment (OSSE) data. We build on the work by Posselt et al. (2019) in which a water vapor feature-59 tracking AMV algorithm was applied to a high-resolution numerical simulation, thus providing a global set of AMV 60 estimates which can be compared to the reference winds produced by the simulation. In this case, a synthetic "true" 61 state is available with which AMVs can be compared and errors are quantified, and it is shown that errors in AMV 62 estimates are state dependent. Our approach will use a conjunction of machine learning (random forest) and 63 unsupervised parametric clustering (Gaussian mixture models) to build a model for the uncertainty structures found 64 by Posselt et al. (2019). The realism and robustness of the resulting uncertainty estimates depend on the realism and 65 representativeness of the reference dataset. This work builds upon the work of Bormann et al. (2014) and Hernandez-66 Carrascal and Bormann (2014), who showed that wind tracking could be divided into distinct geophysical regimes by 67 clustering based on cloud conditions. This study supplements that approach with the addition of machine learning, 68 which, compared with traditional linear modeling approaches, should allow the model to capture more complex non-69 linear processes in the error function.

70 Traditional techniques for modeling uncertainty through machine learning have focused on characterizing bias but 71 often struggle when estimating standard error. By pairing a random forest algorithm with unsupervised parametric 72 clustering, we propose a data-driven, cluster-based approach for quantifying both bias and standard error from

73 experimental data. According to the theory developed by Nguyen et al. (2019), these improved error characterizations

74 should then lead to improved error characteristics (e.g., lower bias, more accurate uncertainties) in subsequent analyses

75 such as flux inversion or data assimilation.

76 This paper does not purport that the specific algorithm detailed here should supplant error characterization approaches 77 for all AMVs; indeed, most commonly assimilated AMVs are based on tracking cloud features, not water vapor 78 profiles. In addition, this algorithm is trained and developed for a specific set of AMVs extracted from a water vapor 79 field associated with a particular range of flow features. As such, application of our algorithm to modeled or observed 80 AMVs will be most appropriate in situations with similar dynamics to our training set. However, we intend in this 81 paper to demonstrate that the methodology is successful in characterizing errors for this set of water vapor AMVs and 82 suggest that this approach- that is, capturing state-dependent uncertainties in feature-tracking algorithms through a 83 combination of clustering and random forest- could be implemented in other feature-tracking AMV extraction 84 methods and situations.

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

92 2. Experimental Set-up

93 2.1 Simulation and Feature-Tracking Algorithm

94 We trained our model on the simulated data used by Posselt et al. (2019), which applied an AMV algorithm to outputs 95 from the NASA Goddard Space Flight Center (GSFC) Global Modeling and Assimilation Office (GMAO) GEOS-5 96 Nature Run (G5NR; Putman et al. 2014). The Nature Run is a global dataset with ~7 km horizontal grid spacing that 97 includes, among other quantities, three-dimensional fields of wind, water vapor concentration, clouds, and 98 temperature. Note that throughout the text we will use the term 'Nature Run wind' to refer to reference winds in the 99 simulation dataset used to train the uncertainty model. The AMV algorithm is applied on four pressure levels (300hPa, 100 500hPa, 700hPa, and 850hPa) at 6-hourly intervals, using three consecutive global water vapor fields spaced one hour 101 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 102 local-area pattern matching algorithm that approximates wind speed and direction from movement of the matched 103 patterns. The algorithm searches a pre-set number of nearby pixels to minimize the sum-of-absolute-differences 104 between aggregated water vapor values across the pixels. Posselt et al. (2019) describes the sensitivity of the tracking algorithm and the dependency of the tracked winds on atmospheric states in detail. The coordinates of the data are on a 5758 x 2879 x 240 spatio-temporal grid for the longitude, latitude, and time dimension, respectively.

107 It is important to note that the AMV algorithm tracks water vapor on fixed pressure levels. In practice, these would be

108 provided by satellite measurements, whereas in this paper we use simulated water vapor from the GEOS-5 Nature

Run. In this simulation height assignment of the AMVs is assumed to be perfectly known. This assumption is far from

110 guaranteed in real world applications but, as previously discussed, its implications are not pursued in this paper. As

such, we focus solely on observational AMV error and not on height assignment error. We note that in practice, one

112 approach to understanding the behavior and accuracy of the wind-tracking algorithm is to apply it to modeled data

113 (e.g., Posselt et al., 2019). Our approach seeks to complement this approach by providing a machine-

114 learning/clustering hybrid approach that can further divide comparison domains into 'regimes' which may provide

115 further insights into the behavior of the errors and/or feedback into the wind-tracking algorithm.

A snapshot of the dataset at 700hPa is given in Figure 1, where we display the water vapor from Nature Run (top left

117 panel), the wind speed from Nature Run (top right panel), the tracked wind from the AMV-tracking algorithm (bottom

118 right panel), and the difference between the Nature Run and tracked wind (bottom left panel). Note that the wind-

119 tracking algorithm tends to have trouble in region where the Nature Run water vapor content is close to zero. It is clear

120 that while the wind-tracking algorithm tends to perform well in most regions (we can classify these regions as areas

121 where the algorithm is skilled), in some regions the algorithm is unable to reliably make a reasonable estimate of the

- 122 wind speed (unskilled). We will examine these skilled and unskilled regimes (and their corresponding contributing
- 123 factors) in section 3.

124 2.2 Importance of Uncertainty Representation in Data Assimilation

125 Proper error characterization for any measurement, including AMVs, is important in data assimilation. Data

126 assimilation often uses a regularized matrix inverse method based on Bayes' theorem, which, when all probability

127 distributions in Bayes' relationship are assumed to be Gaussian, reduces to minimizing a least-squares (quadratic) cost

128 function Eq (1):

129

$$\mathbf{J} = (\mathbf{x} - \mathbf{x}_{\mathbf{b}})\mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_{\mathbf{b}}) + ((\mathbf{y} - \mathbf{a}) - \mathbf{H}[\mathbf{x}])^{\mathrm{T}}\mathbf{R}^{-1}((\mathbf{y} - \mathbf{a}) - \mathbf{H}[\mathbf{x}])$$
(1)

where **x** represents the analysis value, $\mathbf{x}_{\mathbf{b}}$ represents the background field (first guess), **B** represents the background error covariance, **y** represents the observation, and **H** represents the forward operator that translates model space into observation space. This translation may consist of spatial and/or temporal interpolation if **x** and **y** are the same variable (e.g., if the observation of temperature comes from a radiosonde), or may be far more complicated (e.g., a radiative transfer model in the case of satellite observations). **R** represents the observation error covariance, and **a** represents the accuracy, or bias, in the observations. The right hand side of Eq. (1) can be interpreted as a sum of the contribution of information from the data ($\mathbf{y} - \mathbf{H}[\mathbf{x}] - \mathbf{a}$) and the contribution from the prior ($\mathbf{x} - \mathbf{x}_{\mathbf{b}}$), which are weighted by their Formatted: Font: Not Bold, Not Italic Deleted: Figure 1 respective covariance matrices. In our analysis, the AMVs obtained from the wind-tracking algorithm is used as 'data'
in subsequent analysis. That is, the tracked wind data y is a biased and noisy estimator of the true wind y, and might

140 be assumed to follow the model Eq. (2):

$$y = y + \epsilon$$
 (2)

142 where ϵ is an error term, commonly assumed to be Gaussian with mean **a** and covariance matrix **R** (i.e., $\epsilon \sim N(\mathbf{a}, \mathbf{R})$), 143 which are the same two terms that appear in Equation (1). As such, for data assimilation to function, it is essential to 144 correctly specify the bias vector **a** and the standard error matrix **R**. Incorrect characterizations of either of these

145 components could have adverse consequences on the resulting data assimilation analyses with respect to bias and/or

146 the standard error (Nguyen et al., 2019).

147 3 Methodology

148 3.1 Generalized Error Characterization Model

149 An overview of our approach is outlined in Figure 2, Given a set of training predictors X, training responses \hat{Y} , and 150 simulated true response Y, our approach begins with two independent steps. In one step, a Gaussian mixture model is 151 trained on the set of X, \hat{Y} , and Y. This clustering algorithm identifies geophysical regimes where the nonlinear 152 relationships between the three variables differ. In the other step, a random forest is used to model Y based on X and

153 \hat{Y} . This step produces an estimate of the true response (we call this Y) using only the training predictors and response.

154 We then employ the Gaussian mixture model to estimate the clusters which the set of X, \hat{Y} , and Y pertain to.

155 Subsequently, we compute the error characteristics of each cluster of X, \hat{Y} , and Y in the training dataset. Thereafter,

156 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

157 characteristics.

158 In this paper, we are primarily interested in the distribution of a retrieved quantity versus the truth. That is, given a 159 retrieved value \hat{Y}_i , we are interested in the first and second moments (i.e., E($\hat{Y}_i - Y$) and var($\hat{Y}_i - Y$)), respectively. 160 We note that there is a large body of existing work on uncertainty modeling in the machine learning literature (e.g., 161 Coulston et al., 2016; Tripathy et al., 2018; Tran et al., 2019; Kwon et al., 2020), although these approaches primarily 162 define the uncertainty of a prediction as var(\hat{Y}_i), or quantify how sensitive that prediction is to tiny changes in the 163 models/inputs. Our approach, on the other hand, characterizes the error as var($\hat{Y}_i - Y$), which describes how accurate 164 a prediction is relative to the true value. For this reason, our methodology is more stringent in that it requires 165 knowledge of the true field (which comes naturally within OSSE framework) or some proxies such as independent 166 validation data or reanalysis data. In return, the error estimates from our methodology fit naturally within the data 167 assimilation framework (that is, it constitutes the parameter R in Eq. (1)).

168 What follows in this paper is an implementation of the error characterization model obtained for a subsample of the 169 GEOS-5 Nature Run at a fixed height of 700hPa. In particular, we trained the error characterization on a random Formatted: Font: Not Bold

Deleted: Figure 2

171 subsample from the first 1.5 months of the Nature Run, and show the results obtained when applying it to a test

172 subsample drawn from the subsequent 0.5 months of the Nature Run.

173 3.2 Error Regime

174 When examining the relationship between AMVs and Nature Run winds in Figure 3, it is clear that there are two 175 distinct 'error-regimes' present in the dataset. The majority of AMV estimates can be categorized as 'skilled', wherein 176 their estimate lies clearly along a one-to-one line with the Nature Run wind. However, there is also clearly an 177 'unskilled' regime, for which the AMV estimate is very close to zero when there are actually moderate or large Nature 178 Run wind values present. Our goal is to provide unique error characterizations for each error regime, because the error 179 dynamics are different within each regime. Furthermore, when we analyze this error and its relationship to water 180 vapor, we see that 'unskilled' regime correlates highly with areas of low water vapor in Figure 4. This matches the 181 error patterns discussed in Posselt et al. (2019). We note that the division between skilled and unskilled regimes does 182 not need to be binary. For instance, in some regions the wind-tracking algorithm might be unbiased with high-183 correlation with the true winds, and in other regions the algorithm might still be unbiased relative to the true winds, 184 but with higher errors. The second situation is clearly less skilled than the first, although it might still be considered 185 'skilled', and this separation of the wind-tracking estimates into various 'grades' of skill forms the basis of our error 186 model.

187 3.3 Gaussian Mixture Model

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

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

$$P(\theta) = \sum_{i=1}^{K} \pi_{i} N(\mu_{i}, \Sigma_{j})$$

200

199

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

Formatted: Font: Not Bold, Not Italic

Deleted: Figure 3

Deleted: Figure 4

Formatted: Font: Not Bold, Not Italic

(3)

- 203 K is the number of clusters, and π_j is the mixture proportion.
- 204 4. Determine π_j, μ_j, Σ_j for K clusters using the Expectation–Maximization Algorithm 205 5. From 3. and 4., estimate the probability of a given x_i belonging to the j-th cluster as $P(x_i \in k_j) = p_{ij}$
- $206 \qquad \qquad 6. \quad \text{Assign point } x_i \text{ to the cluster with the maximum probability } p_{ij}$
- 207 The mixture model clustering is based on the R package 'Mclust' developed by Fraley et al. (2012), which builds upon
- 208 the theoretical work of Fraley and Raftery (2002) for model-based clustering and density estimation. The process uses 209 an Expectation-Maximization algorithm to cluster the dataset, estimating a variable number of distinct multivariate
- 210 Gaussian distributions from a sample dataset. Training the Gaussian mixture model on this dataset provides a
- 211 clustering function which outputs a unique cluster for any data point with the same number of variables.

212 In one dimension, a Gaussian mixture model looks like the distributions depicted in Figure 5, instead of modeling a

213 population as a single distribution (Gaussian or otherwise), the GMM algorithm fits multiple Gaussian distributions 214 to a population. One key aspect of this algorithm is the capability of assigning a new point to the most likely 215 distribution. For example, in the 1-D figure, a normalized AMV estimate with a value of 10 would be more likely to 216 originate from the broad cluster '2' than the narrow cluster '4'. In this case, we model the population as a Gaussian 217 mixture model in five-dimensional space, which consists of two Nature Run wind vector components (u and v), two 218 AMV estimates of these wind components (u and v), and the simulated water vapor values, all of which have been 219 standardized to have mean 0 and standard deviation of 1. Each cluster has a 5-dimensional mean vector for the center 220 and a 5x5 covariance matrix defining their multivariate Gaussian shape. The estimation of a covariance matrix allows 221 for the characterization of the relationships between the different dimensions within each cluster, and as such the 222 gaussian mixture model approach provides greater potential for understanding the geophysical basis of error regimes 223 than other unsupervised clustering approaches.

224 We note that the choice of inputs to the clustering methodology is limited, and that a more successful clustering may 225 be achieved by including additional meteorological or geographic information. However, the intention of this paper 226 is to study the ability of a purely data-driven approach, where no additional information or assumptions are passed to 227 the machine learning model outside of the inputs and outputs to the AMV algorithm itself. Posselt et al. (2019) showed 228 that state dependent uncertainties are a major source of error in water vapor AMVs; introducing further information 229 may cloud our ability to discern these specific uncertainties. While scaling this methodology to other applications may 230 incentivize tailoring to specific conditions, this paper aims to demonstrate that modifications are encouraged for 231 improvement, but not necessary for success.

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

Formatted: Font: Not Bold, Not Italic Deleted: Figure 5

Deleted: Figure 6

238 but divide the AMV estimates by cluster. We see that, for example, clusters 4, 5, and 7 clearly represent cases in which 239 the feature-tracking algorithm provides an accurate estimate of the Nature Run winds, with very low variance around 240 the one-to-one line (i.e., low estimation errors). Clusters 1, 2, 3, and 9 are somewhat noisier than the low-variance 241 clusters, with error characteristics similar to those of the entirety of the dataset. That is, they are considered less skilled, 242 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 243 clearly unskilled in different ways. Cluster 6 is a noisy regime, which captures much of the more extreme differences 244 between the AMV estimates and the Nature Run winds. Cluster 8, on the other hand, represents the low AMV estimate, 245 high Nature Run wind regime. This cluster is returning AMVs with values of zero where the Nature Run wind is 246 clearly non-zero because of the very low water vapor present. We further see the stratification of the regimes when 247 analyzing the absolute AMV error in relation to the water vapor content (Figure 7). We see that clusters that have

248 similar behaviors in the error pattern (such as 1, 2, and 3) represent different regimes of water vapor content.

249 We specified 9 individual clusters due to a combination of quantitative and qualitative reasons. Quantitatively, the 250 'Mclust' package uses the Bayesian Information Criterion (BIC), a model selection criterion based on the likelihood 251 function which attempts to penalize overfitting, to select the optimal number of clusters given an input range. Using 252 an input range of one through nine, the BIC indicated the highest number of clusters would be optimal. More 253 importantly, however, the 9 clusters can be physically distinguished and interpreted. Plots of the geophysical variables 254 in the testing set associated with each of the clusters are shown in Figures 8-11. Specifically, Figure 8 plots the 255 distribution of water vapor for each cluster, while Figure 9 plots the mean wind magnitude in each direction by cluster. 256 Figure 10 plots the correlation matrix for each cluster and Figure 11 show the geographic distribution of each cluster. 257 In looking at these in combination, we see discernable and discrete clusters with unique characteristics. For example, 258 cluster 1 captures the very dry, high-wind regime in the southern hemisphere visible in Figure 2. Cluster 7 259 encompasses the tropics, while cluster 3 captures mid-latitude storm systems. Clusters 6, 8, and 9 are all characterized 260 by a much worse performance of the AMV tracking algorithm, exhibited both in Figure 7 and in Figure 8 but all 261 encompass different geographic and geophysical regimes. We see that the clustering algorithm succeeds in capturing 262 physically interpretable clusters without having any knowledge of the underlying physical dynamics. We note that in 263 other applications, the optimal number of clusters will change and the researcher will need to explore various choices 264 of this parameter in their modeling, although this tuning process should be greatly simplified by the inclusion of an 265 information criterion (e.g., BIC) in the GMM algorithm.

266 3.5 Random Forest

The clustering algorithm requires the Nature Run wind vector component values (u and v) in order to classify the AMV error. When applying the algorithm in practice to tracked AMV wind from real observations, the true winds are unknown. To represent the fact that we will not know the true winds in practice, we develop a proxy for the Nature Run winds using only the AMV estimates and the simulated water vapor itself. This is an instance in which the application of machine learning is desirable, since machine learning excels at learning high-dimensional non-linear 272 relationships from large training datasets. In this case, we specifically use random forest to create an algorithm which 273 predicts the Nature Run wind values as a function of the tracked wind values and water vapor.

274 Random forest is a machine learning regression algorithm which, as detailed by Breiman (2001), employs an ensemble

275 of decision trees to model a nonlinear relationship between a response and a set of predictors from a training dataset.

276 Here, we chose random forest specifically because it possesses certain robustness properties that are more appropriate

277 for our applications than other machine learning methods. For instance, random forest will not predict values that are

278 outside the minimum and maximum range of the input dataset, whereas other methods such as neural networks can

- 279 exceed the training range, sometimes considerably so. Random forest, due to the sampling procedure employed during
- 280 training, also tends to be robust to overtraining in addition to requiring fewer tuning parameters compared with
- 281 methods such as neural networks.

282 We trained a random forest with 50 trees on a separate set of tracked winds and water vapor values to predict Nature

283 Run winds using the 'randomForest' package in the R programming language. While the random forest estimate as a

284 whole does not perform much better than the AMV values in estimating the Nature Run wind (2.89 RMSE for random

285 forest vs 2.91 RMSE for AMVs), as shown in Figure 12, it does not display the same discrete regimentation as the

286 AMV estimates in Figure 3. As such, the random forest estimates can act as a proxy for Nature Run wind values in

287 our clustering algorithm — they remove the regimentation which is a critical distinction between the AMV estimates

288 and the Nature Run wind values.

289 3.6 Finalized Error Characterization Model

290 The foundation of the error characterization approach is to combine the random forest and clustering algorithm. We 291 apply the Gaussian mixture model, as trained on the Nature Run winds (in addition to the AMVs and water vapor), to

292 each point of water vapor, AMV estimate, and associated random forest estimate. This produces a set of clusters

293 which, when implemented, require no direct knowledge of the actual Nature Run state (Figure 13).

294 Naturally, the clustering algorithm performs better when applied to the dataset with the Nature Run winds, as

295 opposed to winds generated from the random forest algorithm. The former is created with direct knowledge of the

296 Nature Run winds, and any approximation will lead to increased uncertainties. In practice, the performance of the 297

- cluster analysis can be improved by enhancing the performance of the random forest itself. As with any machine 298 learning algorithm, the random forest contains hyperparameters that can be optimized for specific applications. In
- 299
- addition, performance could be improved by including additional predictor variables. Our intent is not to use the

300 random forest as a wind tracking algorithm; rather, the random forest is presented in this paper as a proof of concept.

- 301 Nonetheless, we see in Figure 13 and Figure 14 that the error characterization still discretizes the testing data set into 302 meaningful error regimes. The algorithm manages to separate the AMV estimates into appropriate error clusters. Once 303 again, clusters 6 and 8 manage to capture unskilled regimes, and cluster 7, and to a lesser extent clusters 4 and 5,
- 304 remain skillful. By taking the mean and standard deviation of the difference between AMV estimates and Nature Run

Formatted: Font: Not Bold, Not Italic

Deleted: Figure 12

Deleted: Figure 13 Formatted: Font: Not Bold, Not Italic 307 winds in each cluster, we develop error characteristics for each cluster (Figure 15); these quantities are precisely the

 $\frac{1}{308}$ bias and uncertainty that we require for the cost function J in Eq (1). We see that the unskilled clusters have very high

309 standard errors and they correspond roughly to the areas of unskilled regimes in Figure 3. Similarly, skilled clusters

310 5, 4 and 7 have standard errors below that of the entire dataset. Since each cluster now has associated error

311 characteristics (e.g., bias and standard deviation), it is then straightforward to assign the bias and uncertainty for any

312 new tracked wind observation by computing which regime it is likely to belong to.

313 3.7 Experimental Set up

314 In this section we will describe our experimental setup for training our model on the GEOS-5 Nature Run data and

- 315 testing its performance on a withheld dataset. We divide the dataset into two parts: a training set consisting of the first
- 316 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
- 317 training/testing procedure for the simulation data and tracked wind is as follows:
- Divide the simulation data and tracked wind into two sets: training set of 1,000,000 points from the first 1.5
 months of the Nature Run and a testing set of 1,000,000 points from the final 0.5 months of the Nature Run.
 We train a Gaussian Mixture Model on a normalized random sample of observations from the training dataset
- 321 of Nature Run winds (u and v direction), tracked winds (u and v direction), and water vapor with n=9 clusters.
- We train two separate random forests on a different random sample of 750,000 observations from the training
 dataset. We use tracked wind (u and v direction) and water vapor to model, separately, Nature Run winds in
 both the u and v directions.
- We apply the random forests to the dataset used for the Gaussian Mixture Model. This provides a randomforest estimate for each point, which is used as a substitute for Nature Run wind values in the next step.
- 327 5. We predict the Gaussian mixture component assignment for each point of water vapor, tracked winds, and328 random forest estimate using the GMM parameters estimated in Step 2.
- We compute the mean and standard deviation of the difference between the tracked winds and the Nature
 Run winds, per direction, for each Gaussian mixture model cluster assignment. This provides a set of error
 characteristics that are specific to each cluster.
- We can apply the random forest, and then the cluster estimation, to any set of water vapor and tracked AMV
 estimates. Thusly, any set of tracked AMV estimates and water vapor can be mapped to a specific cluster,
 and therefore its associated error characteristics.

335 4 Results and Validation

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

Formatted: Font: Not Bold, Not Italic Deleted: Figure 15

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

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

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

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

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

Thus far we have shown that our method produces more accurate error-characterization than an alternative method based on marginal means and variance. Now, we assess whether our methodology provides valid probabilistic prediction; that is, we test whether the uncertainty estimates provided are consistent with the empirical distribution of the validation data. To assess this, we construct a metric in which we normalize the difference between the Nature Run wind and the tracked wind by the predicted variance. That is, for the *i*-th observation, we compute the normalized values for u_i and v_i using the following equations: Formatted: Font: Not Bold, Not Italic

Deleted: Figure 16

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

375

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

376 Where u_i is the *i*-th Nature Run u wind from the Nature Run data, u_i is the tracked-wind, and $\sigma_{u,i}$ is the error as 377 assessed by our model (recall that it is a function of the regime index to which u_i has been assigned). The values for 378 the v-wind are defined similarly. The residuals in Eq (5) can be considered as a variant of the z-score, and it is 379 straightforward to see that if our error estimates are valid (i.e., accurate), then the normalized residuals in Eq. (5) 380 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 381 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 382 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 . 383 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 384 that our error characterization model is highly accurate when forecasting uncertainties. 385 A further validation of our methods encompasses an analysis of the statistical significance of the uncertainty in our 386 model. To this end, we constructed confidence intervals for the bias and standard deviation within each regime using 387 the bootstrap (Efron and Tibshirani, 1993). The procedure of our bootstrap is as follows, 388 Subset the data to retain only observations with regime index j. Let's assume that we have Nj observation 389 within this data subset 390 2. Sample with replacement N_j observations from this subset. This forms a bootstrap sample 391 3. From 2., compute an estimate of the bias and standard deviation. 392 4. Repeat step 2-3 for 1000 times, giving us 1000 estimates of the bias and 1000 estimates of the standard 393 deviation within regime j. 394 5. Compute 95% confidence intervals from the 1000 estimates of bias and standard deviation from 4, 395 The results for the confidence intervals (in graphical form in Figure 18. We note that the figure indicates that for 396 many of the biases, they can be considered unbiased since their confidence interval includes 0 (e.g., regimes 2-8 for 397 u-wind). However, the plot also clearly indicates that two regimes are statistically different from 0 (regime 1 and 9). 398 We also note that for the standard deviation maps, the CI's indicate that they are fairly stable (small narrow range) 399 and that most of the regimes have statistically different standard deviation (denoted here visually as CI's that do not 400 overlap one another). We also note that u and v wind direction tend to have very similar patterns, indicating that our 401 regime classification is persistent across u and v. To summarize, the CI plot above indicate that the differences in 402 standard deviation between different regimes are highly statistically significant (as evidenced by the small 403 confidence intervals and their spacing). For the biases, 3 of the regimes are statistically significantly different from 404 the rest (i.e., regimes 1, 6, and 9), while the rest are likely relatively unbiased (i.e., bias = 0).

Deleted: Figure 17Figure 17	
Formatted: Font: Not Bold	
Deleted: Figure 17	
Formatted: Font: Not Bold, Not Italic	
,	

Formatted: Font color: Blac	rmatted: For	t color:	Black
-----------------------------	--------------	----------	-------

Formatted: Space Before:	Auto, After:	Auto, Line	spacing:
1.5 lines			

(Formatted: Font: 10 pt, Font color: Text 1
Y	Formatted: List Paragraph, Space Before: Auto, After:
	Auto, Line spacing: 1.5 lines, Numbered + Level: 1 +
	Numbering Style: 1, 2, 3, + Start at: 1 + Alignment: Left
	+ Aligned at: 0.25" + Indent at: 0.5"
N/	

Formatted: Space Before:	Auto, After:	Auto, Line spacing:
1.5 lines		

407 **5** Conclusion and Discussion

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

414 In this paper we demonstrate the application of a machine learning uncertainty modeling framework to AMVs derived 415 from water vapor profiles intended to mimic hyper-spectral sounder retrievals. The methodology, based on a 416 combination of gaussian mixture model clustering and random forest, identified distinct geophysical regimes and 417 provided uncertainties specific to each regime. This was achieved in a purely data-driven framework; nothing was 418 known to the model except the specific inputs and outputs of the AMV algorithm, deducing the relationship between 419 regime and uncertainty from the underlying multivariate distribution of water vapor, Nature Run wind, and tracked 420 wind. Our algorithm does require one major tuning parameter in the number of clusters for the GMM algorithm, 421 although the search for the 'optimal' number of clusters can be aided by the inclusion of an information criterion (e.g., 422 the BIC) in the GMM model. This implementation is not intended as a 'ready-to-go' algorithm for general use. Instead, 423 we lay the foundation of an uncertainty modelling approach which we plan to implement at a larger scale in subsequent

424	work Nonetheless this bare bones implementation is sufficient to produce improved error estimates of state-dependent
425	uncertainties as detailed in Posselt et al. (2019).

426 We introduce this framework in an environment that is limited and well-behaved, but which nonetheless we believe

427 provides insight into how such an approach would perform at a larger scale. Of course, there are issues when moving

from the controlled environment of the simulation study to large scale applications. We understand these to be: (1) the

429 existence of uncertainty on the tracked humidity values, and (2) the ability of the training dataset to adequately capture

430 <u>both the range of conditions of water vapor and wind speed, and their inherent relationship.</u>

The simulation used for introducing this framework was a 'perfect-observation' environment; that is, the water
 vapor was assumed to be perfectly known to the wind tracking algorithm. In real world scenarios, this is obviously

- not the case. However, we believe that this is mitigated by two factors. Firstly, Posselt et al (2019) also conducted a
- 434 study where measurement noise was added to the water vapor measurement. This did not show to have an effect on
- the uncertainty in the AMV estimate, except where there was the presence of strong vertical wind shear, a situation
- 436 which can be identified a larger scale application. Secondly, given quantified uncertainties on the water vapor
- 437 retrievals themselves (the scope of which is decidedly outside the work of this paper), these could be assimilated
- 438 into the uncertainty modelling framework in a straightforward manner by adding them as a prediction variable in
- 439 both the regime classification and emulator. This would allow for the model to itself ascertain the relationship

Deleted: tool

Deleted:

Deleted: hyperspectral sounder

Formatted: Justified, Space After: Auto

Deleted: methodology Deleted: was

Formatted: Font color: Black

446 data-driven. 447 The reliability of the training dataset is the fundamental assumption of any machine learning approach. To reiterate, 448 we present a methodology which aims to characterize the uncertainty in the difference between a measurement X 449 and its true target X (that is, var (X - X)). As such, we require some proxy for the truth in the development of our 450 model (call this X^*). To expand further, we are modelling the relationship between X and X as a function of water 451 vapor Y, with f(Y) = X and g(Y) = X, where f represents the AMV algorithm and g the 'true' relationship 452 between wind speed and water vapor. Thus, we additionally require a proxy function g^* , which is the relationship 453 implied by the training data output of water vapor and reference winds. In the implementation presented in this 454 paper, q^* is represented by the underlying physical models that model the motion of water vapor and windspeed in 455 the GEOS-5 Nature Run. 456 The fidelity of our framework relies upon the assumption $X^* \sim X$ and $g^* \sim g$. In the simulation study, X^* is the first 457 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 458 nature run simulation. We have given the algorithm a training dataset with what we believe is a plausible range of 459 conditions which could occur in X. To the extent that errors may be seasonally and regionally dependent, it will be 460 more effective to train the error estimation algorithm on data that is expected to represent the specific flow regimes 461 and water vapor features valid for a particular forecast or assimilation period. A range of model data encompassing 462 enough seasonal variability should be a reasonable proxy for the possible range of true X. This would significantly 463 increase the computational demands of training the model (~1 day on a single processor, per pressure level to train 464 the current implementation of the algorithm and an average of 3 days per pressure level, on a non-optimized cluster 465 network to run the AMV extraction on the nature run), although such concerns could be mitigated by strategic 466 subsampling approaches. 467 On the other hand, in this implementation q^* is a perfectly known representation of q, which is the GEOS-5 model 468 that runs the simulation. This is where the simulation approach might create the largest source of uncertainty and 469 unreliability in the model. The true process g can only ever be approximated, and different attempts to do so will 470 involve different tradeoffs when implementing this framework. Users could, for example, use high quality validation 471 data such as matchups with radiosondes. In theory, this provides the best possible approximation of the true process 472 g, but could involve a sparsity of data such that the range of, X* supplied is too narrow for a useful model (indeed, 473 the data might be so sparse as to- from a pure machine learning aspect- reduce the overall fidelity of the model 474 itself). On the other hand, model or reanalysis data can provide dense and diverse training datasets, but rely on the 475 assumption that the underlying physical models in those simulations are an adequate representation of the true 476 process. At the core of atmospheric models such as GOES-5 are the laws of fluid dynamics and thermodynamics. In 477 this context, water vapor is advected by the mean wind and as such the wind and water vapor are intrinsically related

between water vapor uncertainty and AMV estimate uncertainty, without breaking the foundational aspect of being

445

478 in these models. This has been the case since the first atmospheric weather prediction models have been developed.

479 There are of course uncertainties associated with the discretization of the fluid dynamics equations, and sometimes

vater vapor structures that are selected for the wind tracking algorithm.
n both these cases, the model could likely be improved by the inclusion of additional variables in the clustering
lgorithm. These could include a variety of parameters to address different potential problem areas in the model.
nentioned previously, including quantified values of uncertainty in water vapor estimates would algorithmically
ink the uncertainty in the humidity retrieval with the uncertainty in the AMV tracking. Similarly, including
parameters that correlate with geophysical phenomena where the AMV algorithm is known to perform poorly (su
is a marker for vertical wind shear or frontal features) would enable domain knowledge to inform the clustering
lgorithm and emulator. Finally, it is likely that the several parameters used in formulating both the Quality
ndicator (Holmlund et al. 1998) and Expected Error (Le Marshall et al. 2004) approaches would be informative
enhancing the algorithm. One critical aspect for users to consider is that these variables must be continuous
parameterizations, rather than discrete markers (which are often used in quality control); discrete variables canno
asily incorporated into a Gaussian mixture model, or indeed most clustering algorithms. Furthermore, we would
ecommend that users implement parameters that are readily available at the same measurement location and time
he AMV estimate itself. Part of the motivation for the purely state dependent approach in this framework is ease
mplementation; colocation and interpolation could add further uncertainty to the model.
iser. There will always be some degree of uncertainty imparted by the inability of reference dataset to perfectly effect reality (indeed this uncertainty itself could be regime dependent, further complicating a regime dependent
incertainty framework). To some extent, this is true of all validation and uncertainty modelling endeavors.
lowever, thoughtful and careful implementations by users, keeping in mind the prescriptions and concepts detail bove, should mitigate the training data dependent uncertainty.
³ uture users would also be wise to consider improvements in the random forest step of the framework. The
apability of this implementation in discerning accurate error regimes degrades substantially with the introduction
he random forest wind estimates. This work focused on the ability to capture regime dependent error, and as suc
· · · · · · · · · · · · · · · · · · ·
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator;
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith being used. Certainly, many of the additional variables suggested above could be useful towards improving the
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith being used. Certainly, many of the additional variables suggested above could be useful towards improving the andom forest. Users could also investigate replacing the random forest altogether with a different emulator, such
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith being used. Certainly, many of the additional variables suggested above could be useful towards improving the andom forest. Users could also investigate replacing the random forest altogether with a different emulator, such the uncertainty of a gaussian process. Indeed, at its most general, our methodology consists of two parts: an emulator
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith being used. Certainly, many of the additional variables suggested above could be useful towards improving the andom forest. Users could also investigate replacing the random forest altogether with a different emulator, such i neural net or a gaussian process. Indeed, at its most general, our methodology consists of two parts: an emulator ind a clustering algorithm. In this implementation, random forest and Gaussian mixture modelling are the
he random forest was not studied in depth. An improved emulator would certainly increase the accuracy of the incertainty estimates produced by this framework. There are a wide variety of ways to improve the emulator; iltimately, and even more so than the regime classification, these will be specific to the AMV extraction algorith being used. Certainly, many of the additional variables suggested above could be useful towards improving the andom forest. Users could also investigate replacing the random forest altogether with a different emulator, such a neural net or a gaussian process. Indeed, at its most general, our methodology consists of two parts: an emulator and a clustering algorithm. In this implementation, random forest and Gaussian mixture modelling are the upproaches; in theory, these two steps could be accomplished using other algorithms belonging to the appropriate

ed: ¶

515 Thorough domain knowledge, both of the AMV extraction algorithm and the context in which it will be applied, is 516 critical in developing methods to improve it. As discussed previously, the bare bones implementation of our 517 methodology in this paper is intended as a structural presentation of the conceptual framework, not necessarily a 518 finalized model. However, it is also the case that the investigation by Posselt et. al (2019) showed that the variables 519 used in this implementation of the model are those most strongly related with AMV uncertainty in this particular 520 application. The state-dependent errors identified by Posselt et al. (2019) are also expected to apply to other water 521 vapor AMVs. This is because, in general, AMV algorithms have difficulty tracking fields with very small gradients, 522 and will produce systematic errors in situations for which isolines in the tracked field (e.g., contours of constant water 523 vapor mixing ratio) lie parallel to the flow. To the extent that our algorithm represents a general class of errors, the 524 results may be applicable to other geophysical scenarios and other AMV tracking methodologies. As mentioned in the 525 introduction, robust estimates of uncertainty are important for data assimilation, and we expect that our methodology 526 could be used to provide more accurate uncertainties for AMVs used in data assimilation for weather forecasting and 527 reanalysis.

528 Author Contribution

529 Teixeira conceived of the idea with inputs from Nguyen. Teixeira performed the computation. Wu provided the

- 530 experimental datasets along with data curation expertise. Posselt and Su provided subject matter expertise. All authors
- 531 discussed the results. Teixeira wrote the initial manuscript and updated the draft with inputs from co-authors.
- 532 Competing Interest: The Authors declare no conflict of interest.
- 533 Funding Acknowledgment: The research was carried out at the Jet Propulsion Laboratory, California Institute of
- 534 Technology, under a contract with the National Aeronautics and Space Administration (80NM0018D0004). © 2020.
- 535 California Institute of Technology. Government sponsorship acknowledged

536 References

- 537 Bormann, N., Hernandez-Carrascal, A., Borde, R., Lutz, H.J., Otkin, J.A. and Wanzong, S.: Atmospheric motion vectors from model simulations. Part I: Methods and characterization as single-level estimates of wind, Journal of
- vectors from model simulations. Part I: Methods and characterization as single-level estimates of wind, Journal of Applied Meteorology and Climatology, 53(1), 47-64. https://doi.org/10.1175/JAMC-D-12-0336.1, 2014.
- 540 Breiman, L.: Random forests. Machine learning, 45(1), 5-32, 2001.
- 541 Cassola, F. and Burlando, M.: Wind speed and wind energy forecast through Kalman filtering of Numerical Weather
 542 Prediction model output, Applied Energy, 99, 154-166, 2012.
- 543 Coulston, J.W., Blinn, C.E., Thomas, V.A. and Wynne, R.H., 2016. Approximating prediction uncertainty for
 544 random forest regression models. *Photogrammetric Engineering & Remote Sensing*, 82(3), pp.189-197.
 545
- Tibshirani, R.J. and Efron, B., 1993. An Introduction to the Bootstrap. *Monographs on statistics and applied probability*, 57, pp.1-436.

Formatted: Font: (Default) Times New Roman
Formatted: Font: (Default) Times New Roman
Formatted: Font: (Default) Times New Roman

- Fraley, C. and Raftery, A.E.: MCLUST: Software for model-based clustering, density estimation and discriminant
 analysis (No. TR-415). Washington University, Seattle Department of Statistics, 2002.
- 550 Fraley, C., Raftery, A.E., Murphy, T.B. and Scrucca, L: mclust version 4 for R: normal mixture modeling for model-551 based clustering, classification, and density estimation, Washington University, Seattle Department of Statistics, 2012
- 552 Gneiting, T. and Katzfuss, M., 2014. Probabilistic forecasting. *Annual Review of Statistics and Its Application*, *1*, 553 pp.125-151
- Hernandez-Carrascal, A. and Bormann, N.: Atmospheric motion vectors from model simulations. Part II:
 Interpretation as spatial and vertical averages of wind and role of clouds, Journal of Applied Meteorology and
 Climatology, 53(1), 65-82, 2014.
- 557 Holmlund, K., Velden, C. S., & Rohn, M.: Enhanced automated quality control applied to high-density satellitederived winds, Monthly Weather Review, 129(3), 517-529, 2001.
- Kawa, S.R., Erickson, D.J., Pawson, S. and Zhu, Z.: Global CO2 transport simulations using meteorological data
 from the NASA data assimilation system, Journal of Geophysical Research: Atmospheres, 109,
 https://doi.org/10.1029/2004JD004554, 2004.

562

- Kwon, Y., Won, J.H., Kim, B.J. and Paik, M.C., 2020. Uncertainty quantification using Bayesian neural networks in
 classification: Application to biomedical image segmentation. *Computational Statistics & Data Analysis*, 142,
 p.106816.
- Le Marshall, J., Rea, A., Leslie, L., Seecamp, R., & Dunn, M.: Error characterisation of atmospheric motion vectors,
 Australian Meteorological Magazine, 53(2), 2004.
- Mueller, K.J., Wu, D.L., Horváth, Á., Jovanovic, V.M., Muller, J.P., Di Girolamo, L., Garay, M.J., Diner, D.J.,
 Moroney, C.M. and Wanzong, S.: Assessment of MISR cloud motion vectors (CMVs) relative to GOES and
 MODIS atmospheric motion vectors (AMVs), Journal of Applied Meteorology and Climatology, 56(3), 555-572,
 https://doi.org/10.1175/JAMC-D-16-0112.1, 2017.
- Nguyen, Hai, Noel Cressie, and Jonathan Hobbs. "Sensitivity of Optimal Estimation Satellite Retrievals to
 Misspecification of the Prior Mean and Covariance, with Application to OCO-2 Retrievals." *Remote Sensing* 11.23
 (2019): 2770.
- Posselt, D. J., L. Wu, K. Mueller, L. Huang, F. W. Irion, S. Brown, H. Su, D., and C. S. Velden: Quantitative
 Assessment of State-Dependent Atmospheric Motion Vector Uncertainties. J. Appl. Meteor. Clim., In Press.
 https://doi.org/10.1175/JAMC-D-19-0166.1., 2019.
- Putman, W., A.M. da Silva, L.E. Ott and A. Darmenov: Model Configuration for the 7-km GEOS-5 Nature Run,
 Ganymed Release (Non-hydrostatic 7 km Global Mesoscale Simulation). GMAO Office Note No.5 (Version 1.0),
 18, 2014.
- Salonen, K., J. Cotton, N. Bormann, and M. Forsythe: Characterizing AMV Height-Assignment Error by Comparing
 Best-Fit Pressure Statistics from the Met Office and ECMWF Data Assimilation Systems, J. Appl. Meteor.
 Climatol., 54, 225–242, https://doi.org/10.1175/JAMC-D-14-0025.1, 2015.
- -----<u>-</u>-----
- Staffell, I. and Pfenninger, S.: Using bias-corrected reanalysis to simulate current and future wind power output,
 Energy, 114,1224-1239, 2016.
- Swail, V.R. and Cox, A.T.: On the use of NCEP–NCAR reanalysis surface marine wind fields for a long-term North
 Atlantic wave hindcast, Journal of Atmospheric and oceanic technology, 17(4), 532-545, 2000.

- 590 Tran, D., Dusenberry, M., van der Wilk, M. and Hafner, D., 2019. Bayesian layers: A module for neural network uncertainty. In Advances in Neural Information Processing Systems (pp. 14660-14672).
- 593 Tripathy, R.K. and Bilionis, I., 2018. Deep UQ: Learning deep neural network surrogate models for high dimensional uncertainty quantification. *Journal of computational physics*, 375, pp.565-588.
- 595 Velden, C.S. and K.M. Bedka,: Identifying the Uncertainty in Determining Satellite-Derived Atmospheric Motion Vector Height Attribution. J. Appl. Meteor. Climatol., 48, 450–463, https://doi.org/10.1175/2008JAMC1957.1_2009.
- 597 Zeng, X., S. Ackerman, R.D. Ferraro, T.J. Lee, J.J. Murray, S. Pawson, C. Reynolds, and J. Teixeira: Challenges and opportunities in NASA weather research. Bull. Amer. Meteor. Soc., 97, 137-140, 2016.



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.



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





Figure 3: Scatter plot of the simulated Nature Run wind vs AMV estimates for u and v wind in the trainingdataset.



612 Figure 4: Simulated water vapor vs the absolute value of the difference between Nature Run and tracked

613 winds in the training dataset.



615 Figure 5: Example of Gaussian Mixture Model in one dimension. Density Figures for the U-Direction AMV

616 Estimate dimension of fitted Gaussian mixture.

617

614





- specific Gaussian mixture component to which each point in the testing set has been assigned. (A): U-
- Direction Wind (B): V-Direction Wind.



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



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





638 639 Figure 11: Geographic distribution by cluster of AMV retrieval locations in the testing dataset. Each sub-panel represents one cluster.



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

644 (B): V Direction



646

647 Figure 13: Scatterplot of Nature Run wind vs AMV Estimates, each sub-panel corresponding to the specific

648 Gaussian mixture component to which each point in the testing set has been assigned when the Nature Run 649 wind value has been substituted by the random estimate. (A): U-Direction Wind (B): V-Direction Wind





651 Figure 14: Water Vapor vs Absolute Tracked Wind Error, each sub-panel corresponding to the specific

652 Gaussian mixture component each point in the testing set has been assigned when the Nature Run wind value

653 has been substituted by the random estimate. (A): U-Direction Wind (B): V-Direction Wind



654

655 Figure 15: (A): Bias (Left Panel) and Standard Error (Right Panel) for each Gaussian mixture cluster in

656 figure 6, U direction. (B): Same as (A) for V-direction



657

658 Figure 16: CRSP applied to different error approaches. (A): Cluster Errors for U Winds (B): Total Errors

659 for U Winds (C): Cluster Errors for V Winds (D): Total Errors for V Winds.

660



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



