

# Improving Wind Forecasts in the Lower Stratosphere by Distilling an Analog Ensemble into a Deep Neural Network

Salvatore Candido<sup>1</sup>, Aakanksha Singh<sup>1</sup>, and Luca Delle Monache<sup>2</sup>

<sup>1</sup>Loon

<sup>2</sup>University of California San Diego

November 22, 2022

## Abstract

We discuss improving forecasts of winds in the lower stratosphere using machine learning to post-process the output of the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated Forecast System. We post-process global three-dimensional predictions, and demonstrate distilling the analog ensemble (AnEn) method into a deep neural network which reduces post-processing latency to near zero maintaining increased forecast skill. This approach reduces the error with respect to ECMWF high-resolution deterministic prediction between 2-15% for wind speed and 15-25% for direction, and is on par with ECMWF ensemble (ENS) forecast skill to hour 60. Verifying with Loon data from stratospheric balloons, AnEn has 20% lower error than ENS for wind speed and 15% for wind direction, despite significantly lower real-time computational cost to ENS. Similar performance patterns are reported for probabilistic predictions, with larger improvements of AnEn with respect to ENS. We also demonstrate that AnEn generates a calibrated probabilistic forecast.

1 **Improving Wind Forecasts in the Lower Stratosphere**  
2 **by Distilling an Analog Ensemble into a Deep Neural**  
3 **Network**

4 **Salvatore Candido<sup>1</sup>, Aakanksha Singh<sup>1</sup>, and Luca Delle Monache<sup>1,2</sup>**

5 <sup>1</sup>Loon, Mountain View, California, USA

6 <sup>2</sup>Center for Western Weather and Water Extremes, Scripps Institution of Oceanography, University of  
7 California San Diego, La Jolla, California, USA

8 **Key Points:**

- 9 • An analog ensemble generates accurate predictions of lower-stratosphere winds and  
10 reliably quantifies the prediction uncertainty.  
11 • A cloud-based distributed computing implementation builds global three-dimensional  
12 predictions in tens of minutes.  
13 • Distilling the analog ensemble into a deep neural network allows scaling histor-  
14 ical forecasts without slowing post-processing speed.

## Abstract

We discuss improving forecasts of winds in the lower stratosphere using machine learning to post-process the output of the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated Forecast System. We post-process global three-dimensional predictions, and demonstrate distilling the analog ensemble (AnEn) method into a deep neural network which reduces post-processing latency to near zero maintaining increased forecast skill. This approach reduces the error with respect to ECMWF high-resolution deterministic prediction between 2-15% for wind speed and 15-25% for direction, and is on par with ECMWF ensemble (ENS) forecast skill to hour 60. Verifying with Loon data from stratospheric balloons, AnEn has 20% lower error than ENS for wind speed and 15% for wind direction, despite significantly lower real-time computational cost to ENS. Similar performance patterns are reported for probabilistic predictions, with larger improvements of AnEn with respect to ENS. We also demonstrate that AnEn generates a calibrated probabilistic forecast.

## Plain Language Summary

We demonstrate improvements in predicting winds in the stratosphere using machine learning. Our approach uses predictions and analyses from the European Centre for Medium-Range Weather Forecasts (ECMWF). By comparing how previous forecasts differed from what the winds ultimately were over many data points, we are able to modify the current forecast in a way that improves prediction of the winds observed by Loon high altitude balloons in the stratosphere.

A common barrier to using approaches like this to generate global predictions is processing a large amount of information quickly enough to be useful. We demonstrate that by using machine learning we are able to perform many of the slow calculations ahead of time, and that these forecast improvements can be deployed in real applications.

## 1 Introduction

This paper discusses forecasting stratospheric winds by post-processing numerical weather prediction models using machine learning techniques. Specifically, a new variant of the analog ensemble (AnEn; Delle Monache et al., 2013) algorithm that heavily leverages deep neural networks is proposed. The methodology is tested against analysis data and a dataset of observations of stratospheric winds from Loon (<http://www.loon.com>) high altitude balloons (Candido, 2020).

Our focus on winds in the lower stratosphere is driven by Loon’s need to predict the trajectory of high altitude balloons drifting through the stratosphere. Loon is a company that provides connectivity to people in underserved (often remote and rural) locations by placing telecommunications on these balloons. These high altitude platforms can change altitude, and are navigated using a machine learning approach to synthesize in situ observations of winds (from the balloon movements) and wind forecasts. Improved forecast accuracy and reliable uncertainty quantification of the forecasts, which are both key results of the approach we present, determine the navigation efficiency of balloons. Because this navigation system is a real-time operational system (that has navigated balloons for over 1 million hours of flight through the stratosphere), the amount of data to be downloaded from operational forecast centers, the compute needed to utilize that data in real-time operations, and the length of time required to do the post-processing, are also important factors that drive the quality of the system’s operation. These concerns led to the development of a less expensive model in post-processing and distilling the computational burden of the post-processing process into a neural network. It is expected that the approach proposed here to allow the real-time execution of postprocessing methods as the analog ensemble across millions of grid points and several lead times for global

64 predictions, can be applied to several other atmospheric variables and parameters. (See  
65 below for the range of applications for which the analog ensemble method has already  
66 been implemented.)

67 Recently, with the availability of increased computation resources suitable for the  
68 execution of neural networks (e.g., on graphics processing units) and access to large train-  
69 ing data sets, machine learning algorithms have been successfully explored to generate  
70 weather predictions and to postprocess numerical weather predictions (e.g., Tao et al.,  
71 2016; Gagne II et al., 2017; Rasp & Lerch, 2018; Scher, 2018; Chapman et al., 2019; Lagerquist  
72 et al., 2019; Burke et al., 2020). It has also been shown that machine learning can sup-  
73 port the decision-making process associated with high-impact weather phenomena (McGovern  
74 et al., 2017) and it can be leveraged to enhance our physical understanding of atmospheric  
75 processes (Gagne II et al., 2019; McGovern et al., 2019).

76 Analog-based methods, which are a type of machine learning, have been explored  
77 for decades (Lorenz, 1969) to develop predictions for a range of weather parameters. The  
78 basic idea is to find situations from the past similar to the current one and use what un-  
79 folded in these situations to estimate the future evolution of a parameter (Klausner et  
80 al., 2009; Panziera et al., 2011) or to infer the errors of today’s prediction from a dynam-  
81 ical model’s past performance (Delle Monache et al., 2013), an ensemble of model runs  
82 (Hamill & Whitaker, 2006), or other methods (Mahoney et al., 2012; Cervone et al., 2017).

83 One of the challenges of finding these similar situations is the size of the historical  
84 dataset available to the algorithm. Van den Dool (1994) estimated that when match-  
85 ing fields over large spatial domains (e.g., the northern hemisphere) a training dataset  
86  $10^{30}$  years long would be needed to find matches with a degree of analogy below obser-  
87 vational errors. However, Van den Dool (1994) also indicated that if the matching prob-  
88 lem can be reduced to a few degrees of freedom, a much shorter historical dataset can  
89 be sufficient.

90 We apply one such approach, the AnEn (Delle Monache et al., 2011, 2013), to the  
91 prediction of lower-stratosphere winds. In our case, matching to analogous situations is  
92 performed independently at each grid location and lead time over two parameters: wind  
93 speed and direction. Forecast improvements are demonstrated with only two years of pre-  
94 vious forecasts. Versions of the AnEn have been applied successfully for the prediction  
95 of weather parameters (Delle Monache et al., 2013; Nagarajan et al., 2015; Eckel & Delle Monache,  
96 2016; Frediani et al., 2017; Keller et al., 2017; Sperati et al., 2017; Plenkovi et al., 2018;  
97 Yang et al., 2018), tropical cyclone intensity (Alessandrini et al., 2018), air quality (Djalalova  
98 et al., 2015; Huang et al., 2017; Delle Monache et al., 2020), and renewable energy (Mahoney  
99 et al., 2012; Alessandrini, Delle Monache, Sperati, & Nissen, 2015; Alessandrini, Delle Monache,  
100 Sperati, & Cervone, 2015; Vanvyve et al., 2015; Junk et al., 2015; Cervone et al., 2017;  
101 Davò et al., 2016; Ferruzzi et al., 2016; Shahriari et al., 2020), but this is the first ap-  
102 plication of the approach to stratospheric winds.

103 A common issue with real world use of an AnEn-based system is achieving the post-  
104 processing speed that is needed in an operational environment. We outline how a dis-  
105 tributed computing system can apply the conventional AnEn globally using the past two  
106 years of forecasts in around 20 minutes. We demonstrate that this can be even more ef-  
107 ficient by distilling the entire AnEn into a deep neural network (DNN). Distilling, in the  
108 machine learning community, refers to training a DNN to memorize and thus mimic an-  
109 other model. It has been used in reinforcement learning (Rusu et al., 2015), to compress  
110 an ensemble of predictions into a single model (Hinton et al., 2015; Bucilu et al., 2006),  
111 and to approximate a more complex neural network with a simpler one (Ba & Caruana,  
112 2014). In all cases, the idea is to achieve a more computationally efficient version of a  
113 skillful, but perhaps inconvenient model.

114 Since the distilling process is performed offline (in advance), it does not impact real-  
 115 time operations regardless of the size of the historical dataset. This is a key factor given  
 116 that the skill of the AnEn tends to improve with a larger historical dataset.

## 117 2 Methods

118 In this section, we review the AnEn algorithm and discuss how it can be implemented  
 119 at a global scale using distributed computing. We then discuss distilling the AnEn into  
 120 a DNN. We use the former method to demonstrate that the much more efficient latter  
 121 method achieves equivalent performance despite being significantly more desirable for  
 122 use in a production system.

### 123 2.1 Conventional Analog Ensemble Algorithm

124 The AnEn estimates a probability distribution over a forecast parameter, such as  
 125 wind speed or direction, given a forecast, previous forecasts made by the same model,  
 126 and corresponding ground truth for those previous forecasts. A search for analogous sit-  
 127 uations, i.e., previous forecasts we consider to be similar to the current forecast, is per-  
 128 formed and ground truth corresponding to these analogous forecasts is used to construct  
 129 an ensemble (Delle Monache et al., 2013). We report (below) the skill of the analog en-  
 130 semble and its mean, which we use to generate probabilistic and deterministic predic-  
 131 tions, respectively.

132 Let  $f(y|x^f)$  be the probability distribution of the observed value  $y$  of some predicted  
 133 quantity given a model prediction  $x^f$ . The vector  $x^f = (x_1^f, x_2^f, \dots, x_k^f)$  contains  $k$  pre-  
 134 dictors from the model forecast, typically including a forecast value for  $y$  and other fields  
 135 considered to be related or providing context on similarity. In the results reported be-  
 136 low,  $x^f$  includes wind speed and direction.

137 AnEn is a nearest-neighbor algorithm using a learned distance function. The clos-  
 138 est analogs to  $x^f$  from previous forecasts are selected, typically restricting to  $x^i$  at  
 139 the same grid point, i.e., forecasts for the same latitude, longitude, and pressure and made  
 140 for the same lead time. Each forecast has a corresponding ground truth referred to as  
 141  $y^i$ . We denote the set of forecast and observation tuples at a grid point as  $\mathcal{P}$ . We rank  
 142 every  $\mathbf{x}^i \in \mathcal{P}$  by a distance function

$$d(x^f, x^i) = \sum_{j=0}^k \frac{w_j^{\mathcal{P}}}{\sigma_j^{\mathcal{P}}} |x_j^f - x_j^i| \quad (1)$$

143 where  $\sigma_j^{\mathcal{P}}$  is a normalization factor, e.g., the standard deviation, to bring all elements  
 144 of  $x$  into a uniform numeric range and  $w_j^{\mathcal{P}}$  is per-feature weight. The weight and nor-  
 145 malization factors are chosen independently for every grid point to optimize the root-  
 146 mean square-error (RMSE) of the ensemble mean on the training dataset using a leave  
 147 one out cross-validation, with the removed  $(\mathbf{x}^i, y^i)$  used as  $(\mathbf{x}^f, y)$ .

148 The  $N$  analogs with the smallest distance to  $x^f$  form an ensemble forecast. We use  
 149 25 analogs in the results below. The weighted ensemble mean can be used as a deter-  
 150 ministic prediction (Delle Monache et al., 2011). We sort the candidate analogs by  $d(x^f, x^i)$   
 151 and compute the weighted mean on the first  $N$  analogs

$$\hat{y}_{wm} = \alpha \sum_{j=0}^N \frac{y^i}{\max(d(x^f, x^i), \epsilon)} \quad (2)$$

152 where  $\alpha$  is one over the sum of the weights and  $\epsilon$  is a very small constant which guards  
 153 against almost exact matches producing larger weights than can be represented numer-  
 154 ically.

155 This procedure is designed for cases where there is a plurality of analogous situ-  
 156 ations, but in the case of a rare forecast that is, e.g., larger than most samples in the train-  
 157 ing, then the AnEn will predict a reversion to the mean and likely not produce a skill-  
 158 ful forecast. Similar to Alessandrini et al. (2019) we apply a bias correction term to our  
 159 forecast of wind speed.

$$\hat{y}_{bc} = \alpha \sum_{j=0}^N \frac{y^j}{\max(d(x^f, x^j), \epsilon)} + (y^f - \hat{y}_{wm}) m \quad (3)$$

160 where  $m$  is a learned parameter to correct for systematic forecast bias.

## 161 2.2 Global Scale with Distributed Computing

162 While the calculation described above at a particular grid point is tractable, a bar-  
 163 rier to operationalizing a global AnEn system is processing the corpus of analogs, which  
 164 can easily grow to 100's of terabytes of data for three-dimensional global predictions over  
 165 several years. The AnEn algorithm provides a natural partitioning as execution is in-  
 166 dependent for each grid point and lead time. However, the data is not natively parti-  
 167 tioned as both every historical forecast and the current prediction contain a piece of data  
 168 needed to post-process every grid point. The challenge is to organize the data so that  
 169 the calculations can be efficiently executed across many datacenter computers. We use  
 170 the MapReduce paradigm (Dean & Ghemawat, 2004), which allows the computation to  
 171 run on a distributed computing (cloud) infrastructure like Google's Flume (Chambers  
 172 et al., 2010). We describe the mechanics of this technique and provide pseudo-code in  
 173 the supporting information.

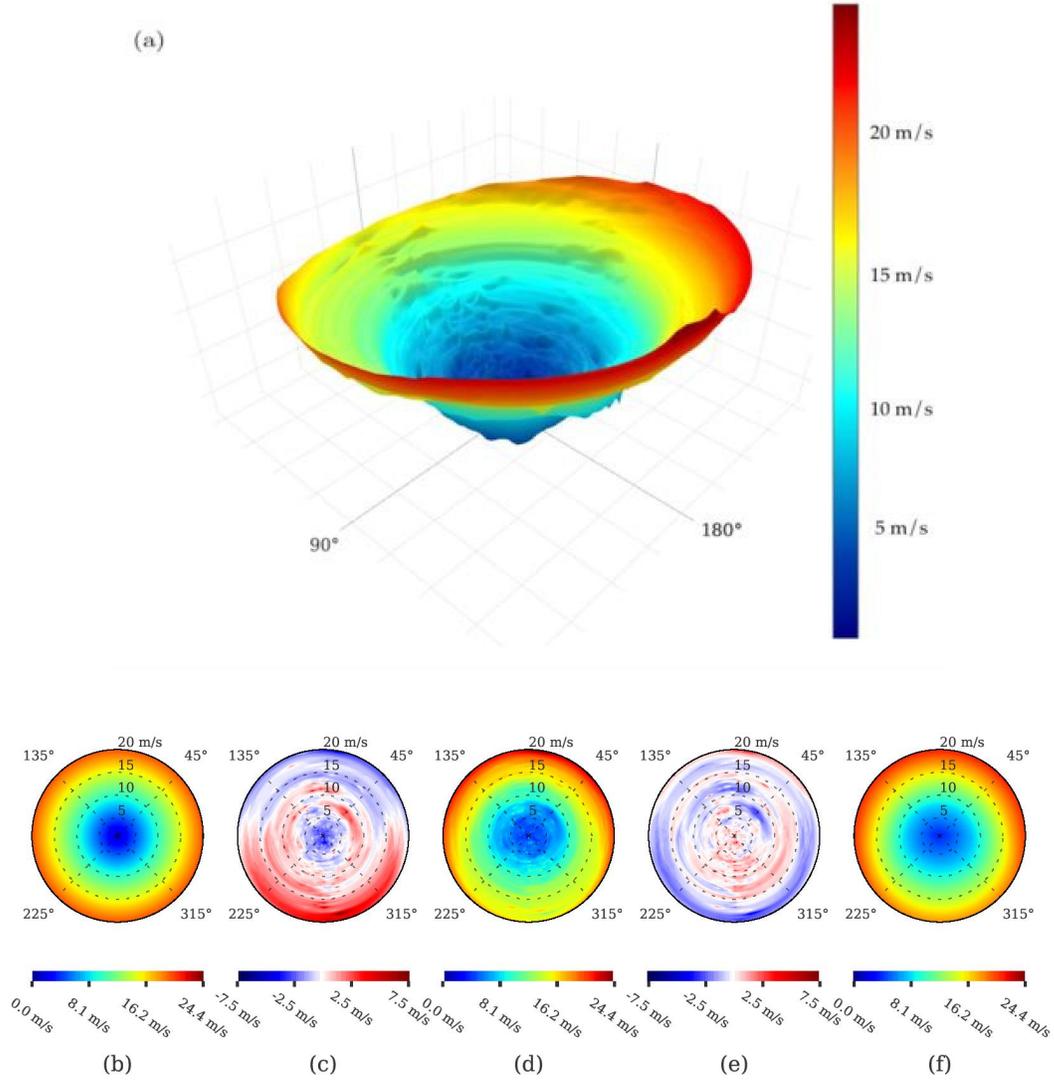
174 Using this technique at appropriate scale, one can post-process a stratospheric wind  
 175 forecast in 10-20 minutes. In our case, we use 100's to 1000's of datacenter machines.  
 176 We create a 3D forecast with 20 pressure levels and 0.5-degree resolution in latitude and  
 177 longitude over 20 lead times. This adds up to the AnEn being applied at around 100 mil-  
 178 lion grid points with analogs from around 3 years of prior forecasts, e.g., around 2196  
 179 candidate analogs per grid point. A rough estimate (ignoring inter-process overhead) of  
 180 trying to do this work on a single machine by multiplying the number of workers by the  
 181 10-20 minute compute time highlights why an implementation on a single machine is likely  
 182 easily too slow for an operational post-processing system.

183 Despite being able to achieve appropriate scale, this is an expensive computation  
 184 that grows proportional to corpus size. Post-processing would take significantly longer  
 185 in the case of a much larger historical corpus. In the next section we discuss distilling  
 186 this computation into a DNN to address this issue.

## 187 2.3 Distilling the Analog Ensemble Into a Deep Neural Network

188 Every value of the analog ensemble mean,  $\hat{y}_{bc}$ , corresponds to an HRES prediction  
 189 of wind speed and direction,  $x^f$ , which has been used to generate the analogs included  
 190 in the set  $\mathcal{P}$ . A DNN can be used to learn, a.k.a., to memorize or distill, the function  
 191 mapping the wind speed and direction of the HRES forecast to the resulting analog en-  
 192 semble mean. An example function for a particular grid point is shown in Figure 1. The  
 193 AnEn mean wind speed values ( $\hat{y}_{bc}$ ; color shading on the isosurface), are shown for each  
 194 HRES forecast  $x^f$  wind speed (distance from the origin) and direction forecast (rotation  
 195 around  $z$ -axis). The plot is roughly conical, and would be exactly conical if AnEn post-  
 196 processing had no effect. Some deformation from the perfect cone is introduced by the  
 197 AnEn algorithm, which we denote by  $h$ , i.e.,  $\hat{y}_{bc} = h_{\mathcal{P}}(x^f)$ .

198 Generating this figure does not require actual new, unseen HRES forecasts. We in-  
 199 stead plot the response of the AnEn in anticipation of potential HRES forecasts. Much  
 200 the same in the learning process, the response curve can be learned by the DNN in ad-



**Figure 1.** The output of the AnEn is a function mapping  $x^f$  to  $\hat{y}_{bc}$ . The plot in (a) shows the speed prediction for a particular  $\mathcal{P}$  swept over speed (distance from origin) and direction (rotation around  $z$ -axis). For speed, this function will typically look like a cone. Taking a top-down view of this plot, we see in (b) the identity operator, i.e., no AnEn post-processing. In (d) we see the same cone as (a) from the top down. Finally (f) shows the output of the distilled AnEn. Note that it resembles the transformation applied by the conventional AnEn but is not expected to be identical as the function is generalized across multiple  $\mathcal{P}$ . The plots in (c) and (e) show the difference ( $\text{m s}^{-1}$ ) between the adjacent plots.

vance of receiving a forecast and when the times comes to post-process the operational forecast we do not require access to the corpus of potential analogs. This makes the distilled AnEn significantly faster than AnEn and more efficient when handling a new forecast in real-time.

In this study we distill  $h_{\mathcal{P}}$  by training over all grid points. We train the DNN to learn the function  $\hat{y}_{distilled} = \hat{h}(k, x^f)$ , where  $k$  is a specific grid point. Because  $h_{\mathcal{P}}$  varies from grid point to grid point, we add the grid point parameters (latitude, longitude, pressure altitude, and lead time) as arguments to  $\hat{h}$  so that the DNN can learn different post-processing transformations at different grid points.

While we discuss results on ensemble mean below, this procedure is not specific to the mean. For example, we have distilled both the ensemble mean of speed and direction (analyzed below), and ensemble forecasts for both quantities into a single DNN with multiple outputs (not shown). The results presented below use a DNN with 10 trainable fully-connected ReLu layers 50 units wide trained with stochastic gradient descent in TensorFlow (Abadi et al., 2015). Full details on the DNN architecture, training parameters, and non-standard data flow, which is conceptually similar to a replay buffer in deep reinforcement learning (Lin, 1992; Mnih et al., 2015), can be found in the supporting information.

We are able to demonstrate a good approximation (next section) for forecast speed and direction within a few billion training examples. Because the training procedure can be performed once (or perhaps periodically, but infrequently) prior to using the network in the operation pipeline, the training time is not particularly important to optimize. Our unoptimized implementation was able to train the network used in our results within a few days on a single CPU (being fed from a distributed data flow). The specific DNN architecture and the data flow to supply training with examples are outlined in greater detail in the supporting information.

Once a network is trained it can be applied point-based, i.e., at a particular place and time with the HRES forecast as input. It adds only a few milliseconds to the real-time computational cost needed to look up forecast data, because the computation performed is a forward (inference) pass through a deep network, i.e., a simple mathematical expression is executed. More study is required to optimize the balance between generalization across different grid points and fitting the particular nuance of a given dataset.

### 3 Results

We present the forecast skill of the AnEn and distilled AnEn aggregated over a half year of forecasts from July to December, 2019, compared against the ECMWF Integrated Forecast System’s high-resolution forecast (HRES) and the ensemble forecast (ENS). We also provide a year long comparison at a different time period (October, 2017 to September, 2018) against the HRES and a persistence ensemble that provides an equivalent result in the supporting information (see Figure S8).

Comprehensive ground truth measurements of winds throughout the stratosphere are not currently available, so to evaluate the quality of the various forecasts we use two proxies for ground truth. The first proxy is the HRES analysis which provides an ‘observation’ comprehensively across all grid points. The second proxy is true observations from Loon high altitude balloons. This dataset of 10.5 million observations, largely concentrated in the lower latitudes, is significantly more sparse as it only allows us to compare forecasts at places and times where a Loon balloon was present. Taken together, these two comparisons characterize the quality of our method.

To summarize the detailed results that follow, the AnEn and distilled AnEn improve the ECMWF Integrated Forecast System’s high-resolution forecast (HRES) of winds

250 in the lower stratosphere. The AnEn methods also produce a skillful probabilistic forecast  
 251 that is able to quantify the forecast uncertainty, which is an advantage over using  
 252 the raw deterministic HRES forecast. The ENS ensemble mean outperforms the AnEn  
 253 methods when evaluating using the HRES analysis as ground truth, but underperforms  
 254 the AnEn methods on the sparser observations from real Loon flights. The AnEn method  
 255 has a significantly reduced computational cost of creating or using a 51-member ensemble  
 256 forecast. Overall the results that follow indicate the AnEn methods are very com-  
 257 petitive when both considering practical implications, and on the merits of forecast skill  
 258 alone.

259 Our region of interest is the lower stratosphere, from around 48 to 145 hPa. We  
 260 apply the technique globally and consider the lead times forecast in the HRES which range  
 261 from 12 hours to 10 days in the future. The results reported in this section are in lat-  
 262 itudes below 70 degrees. Results at higher latitudes are similar, but not shown. Our train-  
 263 ing dataset is the HRES forecasts produced from July, 2016, to June, 2019. We use this  
 264 to choose weights used in the analog matching process. The validation period is over the  
 265 HRES forecasts produced from July, 2019, to December, 2019. The data available in the  
 266 AnEn matching includes all the forecasts in the training dataset plus any additional fore-  
 267 casts between the beginning of the validation time period but prior to the current fore-  
 268 cast. This simulates operational use of an AnEn system. To evaluate the distilled AnEn  
 269 we only use a DNN distilled from the training dataset. In practice, one would distill the  
 270 AnEn into a new DNN from time to time to incorporate additional forecasts into the train-  
 271 ing corpus, but that has not been attempted in this study.

272 Figure 2 shows a comparison of the aggregated add[ldm]deterministic forecast er-  
 273 ror of the HRES, ENS, AnEn, and distilled AnEn grouped by lead time. Note that 90%  
 274 bootstrap confidence intervals are omitted because they are very small because for each  
 275 metric computed and for each lead time we have almost 2 billions and more than 10.5  
 276 millions ground truth / prediction pairs when using HRES and Loon data, respectively.  
 277 The reader can find a view of the these confidence intervals in Figures S4 and S5 of the  
 278 supporting information. Figure 2(a) shows the evaluation performed using the HRES  
 279 operational analysis as the ground truth field. The centered root-mean-square (CRMSE)  
 280 is the portion of the RMSE measuring the random (or anomaly) differences between two  
 281 fields (Taylor, 2001). The AnEn methods have a lower CRMSE than HRES across all  
 282 lead times for wind direction, and after hour 84 for wind speed. The AnEn methods have  
 283 the same skill as ENS up to hour 60 and are competitive for longer lead times, which  
 284 is remarkable considering that AnEn realtime computation cost, given that it is based  
 285 on HRES, is significantly lower than ENS. The correlation between the fields and the  
 286 ground truth is either preserved or improved with the analog-based methods when com-  
 287 pared to HRES. The remaining portion of RMSE is the bias, which in this study is sig-  
 288 nificantly lower than CRMSE for all the prediction systems analyzed (not shown). The  
 289 large reductions of CRMSE for both wind speed and direction obtained with AnEn con-  
 290 firm the ability to tackle conditional biases, which is a result of the algorithm being de-  
 291 signed to learn the error of the current prediction from the errors of analogous past fore-  
 292 casts. The ability of the distilled approach to reproduce AnEn deterministic skill is re-  
 293 markable, as shown by the minimal differences between the two AnEn versions across  
 294 the different metrics and cases considered.

295 Figure 2(b) shows the results when the measurements from Loon stratospheric bal-  
 296 loons are used as ground-truth. This is a much smaller dataset and lacks global cover-  
 297 age, but is real in situ observations from the stratosphere. (see Figure S2 of the support-  
 298 ing information for the geographical distribution of Loon’s measurements). For the con-  
 299 venience of the reader, we provide basic statistical breakdowns and ranges of the obser-  
 300 vations in the dataset overlapping with our validation period in Figure S1 of the sup-  
 301 porting information. The AnEn methods exhibit lower CRMSE than HRES, and signif-  
 302 icantly lower than ENS for both wind speed and direction. AnEn correlation is signif-

303 icantly higher than ENS for wind speed and better than HRES for wind direction. The  
 304 better performance of AnEn compared to ENS when using Loon data can be explained  
 305 by the fact that AnEn, by design, is an excellent downscaling method. This is more ev-  
 306 ident when making a comparison with data that has a high spatial and temporal reso-  
 307 lution, like Loon in situ observations. On the other hand, that is a disadvantage for the  
 308 coarser ENS.

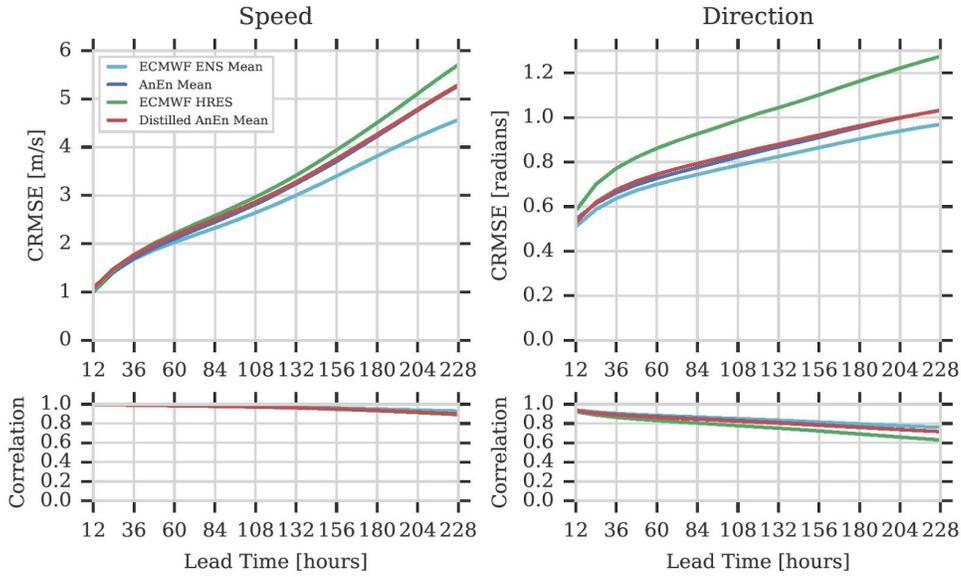
309 We turn our attention to probabilistic forecasts. We compare the ensemble fore-  
 310 cast generated by the AnEn on stratospheric winds to the ENS. Figure 3 shows the con-  
 311 tinuous ranked probability score (CRPS), rank histogram, and binned-spread/skill plot  
 312 across different lead times for the AnEn and ENS. We show these metrics for wind di-  
 313 rection forecasts using the HRES analysis (left) and Loon data (right) as ground truth.  
 314 Results for wind speed are qualitatively similar, and are shown in Figure S3 of the sup-  
 315 porting information.

316 The CRPS provides an assessment of the quality of a probabilistic forecast that is  
 317 not necessarily of a binary event (Hersbach, 2000). It is the probabilistic equivalent of  
 318 the mean absolute error for deterministic predictions, and a zero indicates a perfect fore-  
 319 cast. Similarly to the deterministic results with HRES analysis as the ground truth, AnEn  
 320 is competitive with ENS up to hour 60 and better than HRES at all lead times. How-  
 321 ever, when this performance metric is calculated against the Loon data, AnEn is signif-  
 322 icantly better even than ENS, reducing the latter CRPS between 7 and 70%.

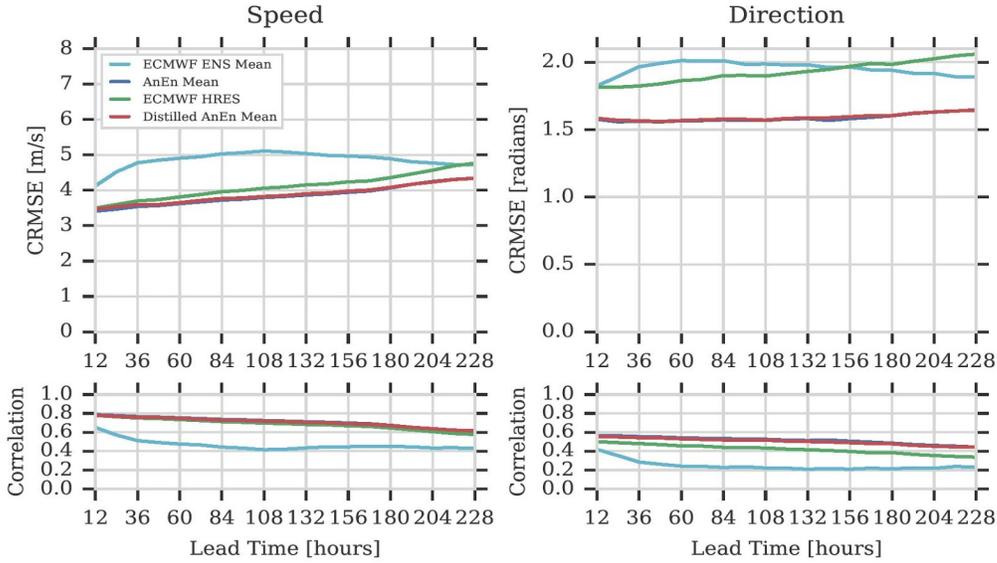
323 The rank histogram estimates the statistical consistency of an ensemble (Anderson,  
 324 1996). For a perfect ensemble, the observation will appear to be drawn from the same  
 325 distribution as the ensemble members. The rank histogram is flat in that case. The ENS  
 326 has a U-shaped rank histogram with both ground-truth data sets, which indicates a lack  
 327 of spread. With HRES as the ground-truth, the AnEn rank histogram instead is closer  
 328 to the ideal flat shape, though it exhibits for the first few lead times a dome shape in-  
 329 dicating an excess spread. This may reflect that the AnEn is including a few analogs that  
 330 have a larger match distance at early lead times. Against Loon data, AnEn has a rank  
 331 histogram significantly closer to the ideal shape, being U-shaped but less so than ENS.

332 The binned-spread/skill plot (van den Dool, 1989; Wang & Bishop, 2003) (which  
 333 is only applicable to probabilistic predictions) characterizes, perhaps, the most impor-  
 334 tant attribute of an ensemble system: the ability to quantify uncertainty while account-  
 335 ing for the flow-dependent error characteristics. This is approximated by analyzing the  
 336 spread-skill relationship across different spread bins. A perfect ensemble results in a di-  
 337 agonal line. Against the HRES analysis, AnEn is closer to the diagonal than ENS, al-  
 338 though both systems exhibit a good spread-skill relationship. However, when Loon mea-  
 339 surements are used as ground-truth, AnEn exhibits a significantly better ability to char-  
 340 acterize the prediction uncertainty. The ENS diagram is horizontal for most bins and  
 341 lead times, which reflects a lack of a spread-skill relationship for the ECMWF ensem-  
 342 ble system when predicting wind direction.

343 Figure 4(a) shows an example of the difference in forecast wind speed between the  
 344 (distilled) AnEn-based forecast and the HRES across a constant-pressure slice of the strato-  
 345 sphere. Figure 4(b) shows the percent change. In this particular example, which was ar-  
 346 bitrarily chosen at random, the largest percent changes are made in the tropics. This  
 347 tends to be a common pattern. Most regions we have analyzed see forecast improvements  
 348 with the AnEn when compared to HRES and the largest improvements are at latitudes  
 349 below 23 degrees. The arrows in Figure 4(b) indicate the flow of the wind direction vec-  
 350 tor field at this pressure level.

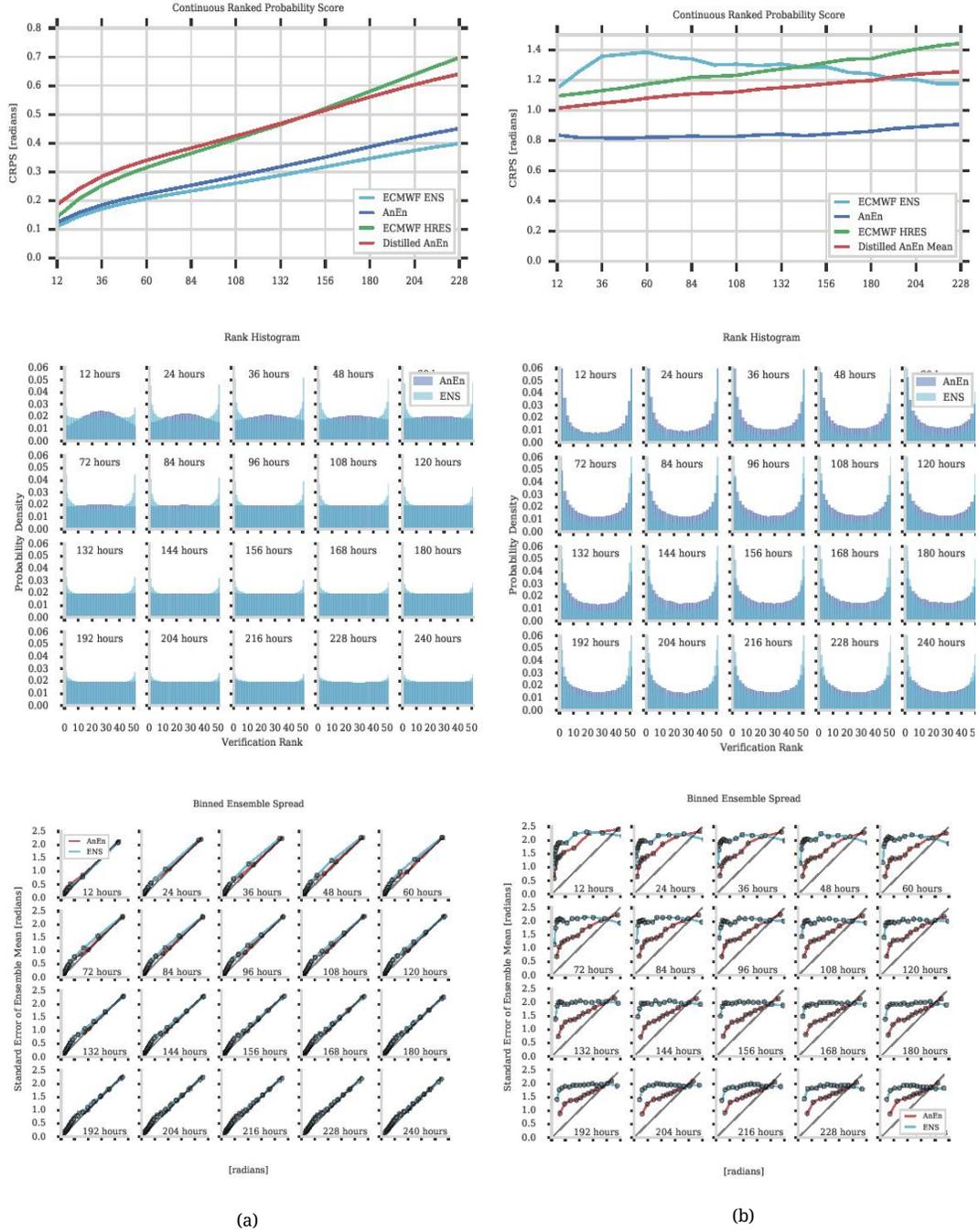


(a)

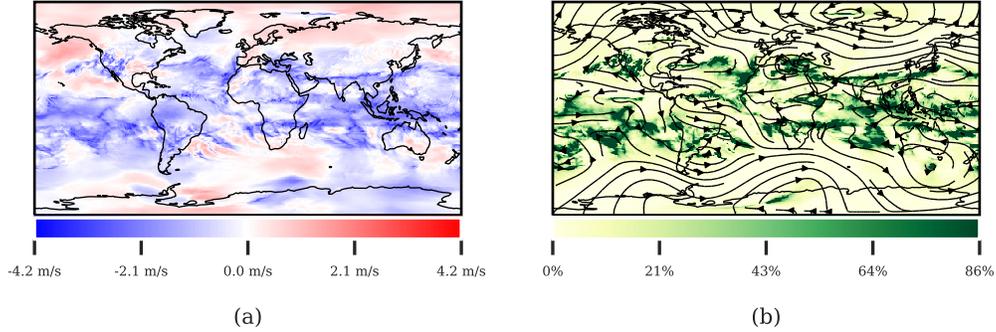


(b)

**Figure 2.** A deterministic wind speed and direction forecast skill comparison between HRES and the means of ECMWF ENS, AnEn, and Distilled AnEn over all lead times is shown using as ground truth (a) HRES analysis and (b) Loon observations of stratospheric winds. The metrics are computed for each lead time across the available observation-prediction pairs from all the grid points in latitudes below 70 degrees.



**Figure 3.** Probabilistic forecast evaluation metrics comparing the AnEn forecast of wind direction to forecasts produced by a ENS. Results with HRES analysis as ground truth are shown on the left (a), while results against Loon's measurements are on the right (b). From top to bottom, the metrics shown are CRPS, rank histogram, and binned-spread skill.



**Figure 4.** Differences between the HRES and distilled AnEn forecast of wind speed worldwide at 50 hPa for 2019-10-20 18:00 UTC with a 5 day lead time. (a) shows the difference between the two forecasts and (b) shows the absolute relative change (absolute value of change between the two forecasts as a percentage of the HRES forecast) between the two forecasts with the direction field overlaid.

#### 4 Discussion

The analog ensemble (AnEn) and distilled AnEn improve the European Centre for Medium-Range Weather Forecasts (ECMWF) high-resolution (HRES) deterministic forecast of winds in the lower stratosphere in our evaluation over half a year of forecasts using both global ECMWF analyses and a smaller set of observations from Loon high altitude balloons as ground truth. The AnEn is also competitive with ECMWF ensemble (ENS) system up to hour 60 for deterministic and probabilistic forecasts when HRES analysis is used as ground truth and significantly better when the performance metrics are computed against Loon’s dataset of true ground truth observations. In particular, AnEn is able to quantify the prediction uncertainty, as evident from the analysis of the probabilistic systems spread-skill relationship, while ENS lacks such attribute, particularly for wind direction predictions. This is true, despite AnEn being computationally cheaper in real-time.

Physics-based numerical weather models, such as the ECMWF’s HRES, are marvels of engineering and science and produce high quality forecasts of many meteorological fields in a coupled and principled manner. However, improvements can sometimes come at great cost, both in research time and in computation and power. Pure machine learning techniques, i.e., end to end learned model-free forecasting, hold promise but are limited due to training on a small number of observations and a limited ability to extrapolate beyond that training data.

For example, a weakness of an analogs-based approach is new situations. If not handled properly, post-processing can reduce forecast skill. We found a specific example of this in our experiments which covered a period of vortex breakdown over North America during February, 2018. Because there was only a single Northern hemisphere winter in our training corpus and it did not exhibit a large vortex breakdown over North America, the algorithm was not able to find analogs with sufficiently high wind speed. When testing the method without the bias correction term of Equation (3), the method decreased forecast skill. While bias correction acts as a stop-gap in this scenario, the desired approach would be to extend the historical corpus to be long enough to find analogous vortex breakdown scenarios.

381 Recently there have been several contributions exploring the potential of machine  
 382 learning for weather and climate predictions (e.g., Tao et al., 2016; Gagne II et al., 2017;  
 383 McGovern et al., 2017; Rasp & Lerch, 2018; Scher, 2018; Chapman et al., 2019; Gagne II  
 384 et al., 2019; Lagerquist et al., 2019; McGovern et al., 2019; Burke et al., 2020). However,  
 385 although there have been encouraging attempts to develop pure machine learning weather  
 386 forecasting methods (e.g., Weyn et al., 2019), those may still be out of reach given the  
 387 relatively low number of available learning examples compared to the number of degrees  
 388 of freedom in the atmosphere. Currently, successful attempts have been reported only  
 389 in replacing individual physical processes (e.g., O’Gorman & Dwyer, 2018).

390 The AnEn distilling procedure can be seen through two lenses. One can consider  
 391 the distilled AnEn as an approximation of the conventional AnEn, i.e., a highly efficient  
 392 implementation of the conventional technique. A second lens is that the DNN is the learn-  
 393 ing technique and the process of distilling the AnEn is a data augmentation method to  
 394 increase the number of examples used to train the network. One may prefer to distill an  
 395 AnEn over directly training a DNN to improve forecasts because DNNs have a high ca-  
 396 pacity (the complexity of the function the model can encode) and, unfortunately, there  
 397 are limited numbers of forecast-ground truth pairs that are available for training. The  
 398 lack of training data is exacerbated by growing the number of outputs we want the DNN  
 399 to produce, e.g., a probability distribution over our forecast field. The AnEn has been  
 400 shown to generalize well as a machine learning algorithm, i.e., to provide an improved  
 401 forecast when deployed on long validation periods on unseen meteorological forecasts.  
 402 The distilled AnEn bootstraps training a DNN off the AnEn, effectively combining the  
 403 AnEn’s strength of being able to generate forecasts with a relatively small corpus of train-  
 404 ing examples with the DNNs ability to memorize this complex correction function with  
 405 a significantly smaller amount of data.

406 This may be a pragmatic compromise. It seems there is a large opportunity for ma-  
 407 chine learning by relying on the extremely high quality numerical weather models and  
 408 making improvements in post-processing. The authors believe there is potential in this  
 409 fused approach. This paper provides an example of how machine learning can contribute  
 410 to increasing forecast skill and uncertainty quantification. As forecasts are asked to be  
 411 simultaneously faster, more granular, and more accurate, the physics-based models can  
 412 continue to do the heavy lifting and machine learning post-processing can improve fore-  
 413 cast quality to alleviate some issues of scale.

## 414 Acknowledgments

415 For more information on Loon stratospheric wind observations contact Loon at stratospheric-  
 416 data@loon.com. The data used in this this paper is portion of the publicly available Loon  
 417 observations data set (Candido, 2020).

418 We are grateful to James Antifaev, Marc Bellemare, Rob Carver, Justin Garofoli,  
 419 Max Kamenetsky, Sameera Ponda, and Aneesh Subramanian for providing useful com-  
 420 ments and edits to an earlier version of this paper.

## 421 References

- 422 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., . . . Zheng, X.  
 423 (2015). *TensorFlow: Large-scale machine learning on heterogeneous systems*.  
 424 Retrieved from <https://www.tensorflow.org/> (Software available from  
 425 tensorflow.org)
- 426 Alessandrini, S., Delle Monache, L., Rozoff, C. M., & Lewis, W. E. (2018). Prob-  
 427 abilistic prediction of tropical cyclone intensity with an analog ensemble.  
 428 *Monthly Weather Review*, *146*(6), 1723-1744. Retrieved from [https://](https://doi.org/10.1175/MWR-D-17-0314.1)  
 429 [doi.org/10.1175/MWR-D-17-0314.1](https://doi.org/10.1175/MWR-D-17-0314.1) doi: 10.1175/MWR-D-17-0314.1

- 430 Alessandrini, S., Delle Monache, L., Sperati, S., & Cervone, G. (2015). An analog  
 431 ensemble for short-term probabilistic solar power forecast. *Applied Energy*,  
 432 *157*, 95–110. Retrieved 2016-11-21, from [http://www.sciencedirect.com/  
 433 science/article/pii/S0306261915009368](http://www.sciencedirect.com/science/article/pii/S0306261915009368)
- 434 Alessandrini, S., Delle Monache, L., Sperati, S., & Nissen, J. N. (2015). A novel ap-  
 435 plication of an analog ensemble for short-term wind power forecasting. *Renew-  
 436 able Energy*, *76*, 768-781. Retrieved from [http://www.sciencedirect.com/  
 437 science/article/pii/S0960148114007915](http://www.sciencedirect.com/science/article/pii/S0960148114007915) doi: [https://doi.org/10.1016/j  
 438 .renene.2014.11.061](https://doi.org/10.1016/j.renene.2014.11.061)
- 439 Alessandrini, S., Sperati, S., & Delle Monache, L. (2019). Improving the analog en-  
 440 semble wind speed forecasts for rare events. *Monthly Weather Review*, *147*(7),  
 441 2677-2692. Retrieved from <https://doi.org/10.1175/MWR-D-19-0006.1> doi:  
 442 [10.1175/MWR-D-19-0006.1](https://doi.org/10.1175/MWR-D-19-0006.1)
- 443 Anderson, J. L. (1996). A method for producing and evaluating probabilistic fore-  
 444 casts from ensemble model integrations. *Journal of Climate*, *9*(7), 1518-1530.  
 445 Retrieved from [https://doi.org/10.1175/1520-0442\(1996\)009<1518:  
 446 AMFPAE>2.0.CO;2](https://doi.org/10.1175/1520-0442(1996)009<1518:AMFPAE>2.0.CO;2) doi: [10.1175/1520-0442\(1996\)009<1518:AMFPAE>2.0.CO;  
 447 2](https://doi.org/10.1175/1520-0442(1996)009<1518:AMFPAE>2.0.CO;2)
- 448 Ba, J., & Caruana, R. (2014). Do deep nets really need to be deep? In Z. Ghahra-  
 449 mani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.),  
 450 *Advances in neural information processing systems 27* (pp. 2654–2662).  
 451 Curran Associates, Inc. Retrieved from [http://papers.nips.cc/paper/  
 452 5484-do-deep-nets-really-need-to-be-deep.pdf](http://papers.nips.cc/paper/5484-do-deep-nets-really-need-to-be-deep.pdf)
- 453 Bucilu, C., Caruana, R., & Niculescu-Mizil, A. (2006). Model compression. In  
 454 *Proceedings of the 12th acm sigkdd international conference on knowledge  
 455 discovery and data mining* (pp. 535–541).
- 456 Burke, A., Snook, N., Gagne II, D. J., McCorkle, S., & McGovern, A. (2020).  
 457 Calibration of machine learningbased probabilistic hail predictions for op-  
 458 erational forecasting. *Weather and Forecasting*, *35*(1), 149-168. Re-  
 459 trieved from <https://doi.org/10.1175/WAF-D-19-0105.1> doi: [10.1175/  
 460 WAF-D-19-0105.1](https://doi.org/10.1175/WAF-D-19-0105.1)
- 461 Candido, S. (2020, April). *Loon stratospheric sensor data*. Zenodo. Retrieved from  
 462 <https://doi.org/10.5281/zenodo.3763022> doi: [10.5281/zenodo.3763022](https://doi.org/10.5281/zenodo.3763022)
- 463 Cervone, G., Clemente-Harding, L., Alessandrini, S., & Monache, L. D. (2017).  
 464 Short-term photovoltaic power forecasting using artificial neural networks  
 465 and an analog ensemble. *Renewable Energy*, *108*, 274 - 286. Retrieved from  
 466 <http://www.sciencedirect.com/science/article/pii/S0960148117301386>  
 467 doi: <https://doi.org/10.1016/j.renene.2017.02.052>
- 468 Chambers, C., Raniwala, A., Perry, F., Adams, S., Henry, R., Bradshaw, R., &  
 469 Nathan. (2010). Flumejava: Easy, efficient data-parallel pipelines. In *Acm  
 470 sigplan conference on programming language design and implementation (pldi)*  
 471 (p. 363-375). 2 Penn Plaza, Suite 701 New York, NY 10121-0701. Retrieved  
 472 from <http://dl.acm.org/citation.cfm?id=1806638>
- 473 Chapman, W. E., Subramanian, A. C., Delle Monache, L., Xie, S. P., & Ralph,  
 474 F. M. (2019). Improving atmospheric river forecasts with machine learn-  
 475 ing. *Geophysical Research Letters*. Retrieved from [https://agupubs  
 476 .onlinelibrary.wiley.com/doi/abs/10.1029/2019GL083662](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019GL083662) doi:  
 477 [10.1029/2019GL083662](https://doi.org/10.1029/2019GL083662)
- 478 Davò, F., Alessandrini, S., Sperati, S., Delle Monache, L., Airoldi, D., & Vespucci,  
 479 M. T. (2016). Post-processing techniques and principal component analysis for  
 480 regional wind power and solar irradiance forecasting. *Solar Energy*, *134*, 327 -  
 481 338. Retrieved from [http://www.sciencedirect.com/science/article/pii/  
 482 S0038092X16300962](http://www.sciencedirect.com/science/article/pii/S0038092X16300962) doi: <https://doi.org/10.1016/j.solener.2016.04.049>
- 483 Dean, J., & Ghemawat, S. (2004). Mapreduce: Simplified data processing on large  
 484 clusters. In *Osd'04: Sixth symposium on operating system design and imple-*

- 485        *mentation* (pp. 137–150). San Francisco, CA.
- 486 Delle Monache, L., Alessandrini, D. I., S., Wilczak, K. J. C., James, & Kumar, R.
- 487        (2020). Improving air quality predictions over the united states with an analog
- 488        ensemble. *Weather and Forecasting*.
- 489 Delle Monache, L., Eckel, F. A., Rife, D. L., Nagarajan, B., & Searight, K.        (2013).
- 490        Probabilistic weather prediction with an analog ensemble. *Monthly Weather*
- 491        *Review*, *141*(10), 3498-3516. Retrieved from [https://doi.org/10.1175/](https://doi.org/10.1175/MWR-D-12-00281.1)
- 492        [MWR-D-12-00281.1](https://doi.org/10.1175/MWR-D-12-00281.1) doi: 10.1175/MWR-D-12-00281.1
- 493 Delle Monache, L., Nipen, T., Liu, Y., Roux, G., & Stull, R.        (2011). Kalman fil-
- 494        ter and analog schemes to postprocess numerical weather predictions. *Monthly*
- 495        *Weather Review*, *139*(11), 3554–3570.
- 496 Djalalova, I., Delle Monache, L., & Wilczak, J.        (2015). Pm 2.5 analog forecast and
- 497        kalman filter post-processing for the community multiscale air quality (cmaq)
- 498        model. *Atmospheric Environment*, *108*, 76–87.
- 499 Eckel, F. A., & Delle Monache, L. (2016, February). A hybrid nwp-analog ensemble.
- 500        *Monthly Weather Review*, *144*, 897–911. Retrieved 2016-02-16, from [https://](https://journals.ametsoc.org/doi/full/10.1175/MWR-D-15-0096.1)
- 501        [journals.ametsoc.org/doi/full/10.1175/MWR-D-15-0096.1](https://journals.ametsoc.org/doi/full/10.1175/MWR-D-15-0096.1) doi: 10.1175/
- 502        [MWR-D-15-0096.1](https://doi.org/10.1175/MWR-D-15-0096.1)
- 503 Ferruzzi, G., Cervone, G., Delle Monache, L., Graditi, G., & Jacobone, F.        (2016).
- 504        Optimal bidding in a day-ahead energy market for micro grid under uncer-
- 505        tainty in renewable energy production. *Energy*, *106*, 194 - 202. Retrieved from
- 506        [://www.sciencedirect.com/science/article/pii/S0360544216302432](http://www.sciencedirect.com/science/article/pii/S0360544216302432)
- 507        doi: <https://doi.org/10.1016/j.energy.2016.02.166>
- 508 Frediani, M. E. B., Hopson, T. M., Hacker, J. P., Anagnostou, E. N., Delle Monache,
- 509        L., & Vandenberghe, F.        (2017, December). Object-based analog forecasts
- 510        for surface wind speed. *Monthly Weather Review*, *145*(12), 5083–5102. Re-
- 511        trieved 2017-12-20, from [http://journals.ametsoc.org/doi/10.1175/](http://journals.ametsoc.org/doi/10.1175/MWR-D-17-0012.1)
- 512        [MWR-D-17-0012.1](http://journals.ametsoc.org/doi/10.1175/MWR-D-17-0012.1) doi: 10.1175/MWR-D-17-0012.1
- 513 Gagne II, D. J., Haupt, S. E., Nychka, D. W., & Thompson, G.        (2019). Inter-
- 514        pretable deep learning for spatial analysis of severe hailstorms. *Monthly*
- 515        *Weather Review*, *147*(8), 2827-2845. Retrieved from [https://doi.org/](https://doi.org/10.1175/MWR-D-18-0316.1)
- 516        [10.1175/MWR-D-18-0316.1](https://doi.org/10.1175/MWR-D-18-0316.1) doi: 10.1175/MWR-D-18-0316.1
- 517 Gagne II, D. J., McGovern, H. S. H., Amy, Sobash, R. A., Williams, J. K., & Xue,
- 518        M.        (2017). Stormbased probabilistic hail forecasting with machine learning
- 519        applied to convectionallowing ensembles. *Weather and Forecasting*, *32*(5),
- 520        1819-1840. Retrieved from <https://doi.org/10.1175/wafd170010.1>
- 521 Hamill, T. M., & Whitaker, J. S.        (2006). Probabilistic quantitative precipita-
- 522        tion forecasts based on reforecast analogs: Theory and application. *Monthly*
- 523        *Weather Review*, *134*(11), 3209-3229. Retrieved from [https://doi.org/](https://doi.org/10.1175/MWR3237.1)
- 524        [10.1175/MWR3237.1](https://doi.org/10.1175/MWR3237.1) doi: 10.1175/MWR3237.1
- 525 Hersbach, H.        (2000). Decomposition of the continuous ranked probability score for
- 526        ensemble prediction systems. *Weather and Forecasting*, *15*(5), 559-570. Re-
- 527        trieved from [https://doi.org/10.1175/1520-0434\(2000\)015<0559:DOTCRP>](https://doi.org/10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2)
- 528        [2.0.CO;2](https://doi.org/10.1175/1520-0434(2000)015(0559:DOTCRP)2.0.CO;2) doi: 10.1175/1520-0434(2000)015(0559:DOTCRP)2.0.CO;2
- 529 Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural net-
- 530        work. In *Nips deep learning and representation learning workshop*. Retrieved
- 531        from <http://arxiv.org/abs/1503.02531>
- 532 Huang, J., McQueen, J., Wilczak, J., Djalalova, I., Stajner, I., Shafran, P., ...
- 533        Delle Monache, L.        (2017). Improving noaa naqfc pm2.5 predictions
- 534        with a bias correction approach. *Weather and Forecasting*, *32*(2), 407-
- 535        421. Retrieved from <https://doi.org/10.1175/WAF-D-16-0118.1> doi:
- 536        [10.1175/WAF-D-16-0118.1](https://doi.org/10.1175/WAF-D-16-0118.1)
- 537 Junk, C., Delle Monache, L., Alessandrini, S., Cervone, G., & von Bremen, L.
- 538        (2015). Predictor-weighting strategies for probabilistic wind power forecasting
- 539        with an analog ensemble. *Meteorologische Zeitschrift*, *24*(4), 361–79.

- 540 Keller, J. D., Delle Monache, L., & Alessandrini, S. (2017, July). Statistical down-  
 541 scaling of a high-resolution precipitation reanalysis using the analog ensemble  
 542 method. *Monthly Weather Review*, *56*, 2081–2095. Retrieved 2017-07-12, from  
 543 <https://journals.ametsoc.org/doi/full/10.1175/JAMC-D-16-0380.1>  
 544 doi: 10.1175/JAMC-D-16-0380.1
- 545 Klausner, Z., Kaplan, H., & Fattal, E. (2009). The similar days method for predict-  
 546 ing near surface wind vectors. *Meteorological Applications*, *16*(4), 569–579. Re-  
 547 trieved from [https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/met.158)  
 548 [met.158](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/met.158) doi: 10.1002/met.158
- 549 Lagerquist, R., McGovern, A., & Gagne II, D. J. (2019). Deep learning for spatially  
 550 explicit prediction of synoptic-scale fronts. *Weather and Forecasting*, *34*(4),  
 551 1137–1160. Retrieved from <https://doi.org/10.1175/WAF-D-18-0183.1> doi:  
 552 10.1175/WAF-D-18-0183.1
- 553 Lin, L.-J. (1992). Self-improving reactive agents based on reinforcement learning,  
 554 planning and teaching. *Machine learning*, *8*(3-4), 293–321.
- 555 Lorenz, E. N. (1969). The predictability of a flow which possesses many scales of  
 556 motion. *Tellus*, *21*(3), 289–307. Retrieved from [https://onlinelibrary](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.2153-3490.1969.tb00444.x)  
 557 [.wiley.com/doi/abs/10.1111/j.2153-3490.1969.tb00444.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.2153-3490.1969.tb00444.x) doi: 10.1111/  
 558 [j.2153-3490.1969.tb00444.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.2153-3490.1969.tb00444.x)
- 559 Mahoney, W. P., Parks, K., Wiener, G., Liu, Y., Myers, W. L., Sun, J., ... Haupt,  
 560 S. E. (2012, Oct). A wind power forecasting system to optimize grid in-  
 561 tegration. *IEEE Transactions on Sustainable Energy*, *3*(4), 670–682. doi:  
 562 10.1109/TSTE.2012.2201758
- 563 McGovern, A., Elmore, K. L., Gagne, D. J., Haupt, S. E., Karstens, C. D.,  
 564 Lagerquist, R., ... Williams, J. K. (2017). Using artificial intelligence to  
 565 improve real-time decision-making for high-impact weather. *Bulletin of the*  
 566 *American Meteorological Society*, *98*(10), 2073–2090. Retrieved from [https://](https://doi.org/10.1175/BAMS-D-16-0123.1)  
 567 [doi.org/10.1175/BAMS-D-16-0123.1](https://doi.org/10.1175/BAMS-D-16-0123.1) doi: 10.1175/BAMS-D-16-0123.1
- 568 McGovern, A., Lagerquist, R., John Gagne, D., Jergensen, G. E., Elmore, K. L.,  
 569 Homeyer, C. R., & Smith, T. (2019). Making the black box more trans-  
 570 parent: Understanding the physical implications of machine learning. *Bul-*  
 571 *letin of the American Meteorological Society*, *100*(11), 2175–2199. Re-  
 572 trieved from <https://doi.org/10.1175/BAMS-D-18-0195.1> doi: 10.1175/  
 573 [BAMS-D-18-0195.1](https://doi.org/10.1175/BAMS-D-18-0195.1)
- 574 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G.,  
 575 ... others (2015). Human-level control through deep reinforcement learning.  
 576 *Nature*, *518*(7540), 529.
- 577 Nagarajan, B., Delle Monache, L., Hacker, J. P., Rife, D. L., Searight, K., Knievel,  
 578 J. C., & Nipen, T. N. (2015, September). An evaluation of analog-  
 579 based post-processing methods across several variables and forecast mod-  
 580 els. *Weather and Forecasting*, *30*, 1623–1643. Retrieved 2015-12-01, from  
 581 <http://journals.ametsoc.org/doi/abs/10.1175/WAF-D-14-00081.1> doi:  
 582 10.1175/WAF-D-14-00081.1
- 583 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameter-  
 584 ize moist convection: Potential for modeling of climate, climate change, and  
 585 extreme events. *Journal of Advances in Modeling Earth Systems*, *10*(10), 2548-  
 586 2563. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351)  
 587 [10.1029/2018MS001351](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001351) doi: 10.1029/2018MS001351
- 588 Panziera, L., Germann, U., Gabella, M., & Mandapaka, P. V. (2011). Noranow-  
 589 casting of orographic rainfall by means of analogues. *Quarterly Journal of the*  
 590 *Royal Meteorological Society*, *137*(661), 2106–2123. Retrieved from [https://](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.878)  
 591 [rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.878](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.878) doi: 10.1002/qj  
 592 .878
- 593 Plenkovi, I. O., Delle Monache, L., Horvath, K., & Hrastinski, M. (2018). De-  
 594 terministic wind speed predictions with analog-based methods over complex

- 595 topography. *Journal of Applied Meteorology and Climatology*, 57(9), 2047-  
 596 2070. Retrieved from <https://doi.org/10.1175/JAMC-D-17-0151.1> doi:  
 597 10.1175/JAMC-D-17-0151.1
- 598 Rasp, S., & Lerch, S. (2018). Neural networks for postprocessing ensemble  
 599 weather forecasts. *Monthly Weather Review*, 146(11), 3885-3900. Re-  
 600 trieved from <https://doi.org/10.1175/MWR-D-18-0187.1> doi: 10.1175/  
 601 MWR-D-18-0187.1
- 602 Rusu, A. A., Colmenarejo, S. G., Gulcehre, C., Desjardins, G., Kirkpatrick, J.,  
 603 Pascanu, R., ... Hadsell, R. (2015). Policy distillation. *arXiv preprint*  
 604 *arXiv:1511.06295*.
- 605 Scher, S. (2018). Toward data-driven weather and climate forecasting: Approx-  
 606 imating a simple general circulation model with deep learning. *Geophys-  
 607 ical Research Letters*, 45(22), 12,616-12,622. Retrieved from [https://  
 608 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018GL080704](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018GL080704) doi:  
 609 10.1029/2018GL080704
- 610 Shahriari, M., Cervone, G., Clemente-Harding, L., & Delle Monache, L. (2020).  
 611 Using the analog ensemble method as a proxy measurement for wind power  
 612 predictability”, journal = ”renewable energy. , 146, 789 - 801. Retrieved from  
 613 <http://www.sciencedirect.com/science/article/pii/S0960148119309668>  
 614 doi: <https://doi.org/10.1016/j.renene.2019.06.132>
- 615 Sperati, S., Alessandrini, S., & Delle Monache, L. (2017). Gridded probabilistic  
 616 weather forecasts with an analog ensemble. *Quarterly Journal of the Royal Me-  
 617 teorological Society*, 143(708), 2874–2885.
- 618 Tao, Y., Gao, X., Hsu, K., Sorooshian, S., & Ihler, A. (2016). A deep neural network  
 619 modeling framework to reduce bias in satellite precipitation products. *Jour-  
 620 nal of Hydrometeorology*, 17(3), 931-945. Retrieved from [https://doi.org/  
 621 10.1175/JHM-D-15-0075.1](https://doi.org/10.1175/JHM-D-15-0075.1) doi: 10.1175/JHM-D-15-0075.1
- 622 Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single  
 623 diagram. *Journal of Geophysical Research: Atmospheres*, 106(D7), 7183–7192.
- 624 van den Dool, H. M. (1989). A new look at weather forecasting through analogues.  
 625 *Monthly Weather Review*, 117(10), 2230-2247. Retrieved from [https://doi  
 626 .org/10.1175/1520-0493\(1989\)117<2230:ANLAWF>2.0.CO;2](https://doi.org/10.1175/1520-0493(1989)117<2230:ANLAWF>2.0.CO;2) doi: 10.1175/  
 627 1520-0493(1989)117(2230:ANLAWF)2.0.CO;2
- 628 Van den Dool, H. M. (1994). Searching for analogues, how long must we wait? *Tel-  
 629 lus A*, 46(3), 314-324. Retrieved from [https://onlinelibrary.wiley.com/  
 630 doi/abs/10.1034/j.1600-0870.1994.t01-2-00006.x](https://onlinelibrary.wiley.com/doi/abs/10.1034/j.1600-0870.1994.t01-2-00006.x) doi: 10.1034/  
 631 j.1600-0870.1994.t01-2-00006.x
- 632 Vanvyve, E., Delle Monache, L., Monaghan, A. J., & Pinto, J. O. (2015). Wind  
 633 resource estimates with an analog ensemble approach. *Renewable En-  
 634 ergy*, 74, 761 - 773. Retrieved from [http://www.sciencedirect.com/  
 635 science/article/pii/S0960148114005308](http://www.sciencedirect.com/science/article/pii/S0960148114005308) doi: [https://doi.org/10.1016/  
 636 j.renene.2014.08.060](https://doi.org/10.1016/j.renene.2014.08.060)
- 637 Wang, X., & Bishop, C. H. (2003). A comparison of breeding and ensemble  
 638 transform kalman filter ensemble forecast schemes. *Journal of the Atmo-  
 639 spheric Sciences*, 60(9), 1140-1158. Retrieved from [https://doi.org/  
 640 10.1175/1520-0469\(2003\)060<1140:ACOBAE>2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060<1140:ACOBAE>2.0.CO;2) doi: 10.1175/  
 641 1520-0469(2003)060(1140:ACOBAE)2.0.CO;2
- 642 Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict  
 643 weather? using deep learning to predict gridded 500-hpa geopotential height  
 644 from historical weather data. *Journal of Advances in Modeling Earth Systems*,  
 645 11(8), 2680-2693. Retrieved from [https://agupubs.onlinelibrary.wiley  
 646 .com/doi/abs/10.1029/2019MS001705](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001705) doi: 10.1029/2019MS001705
- 647 Yang, J., Astitha, M., Delle Monache, L., & Alessandrini, S. (2018). An ana-  
 648 log technique to improve storm wind speed prediction using a dual nwp  
 649 model approach. *Monthly Weather Review*, 146(12), 4057-4077. Re-

650

trieved from <https://doi.org/10.1175/MWR-D-17-0198.1>

doi: 10.1175/

651

MWR-D-17-0198.1

Figure 1.

(a)

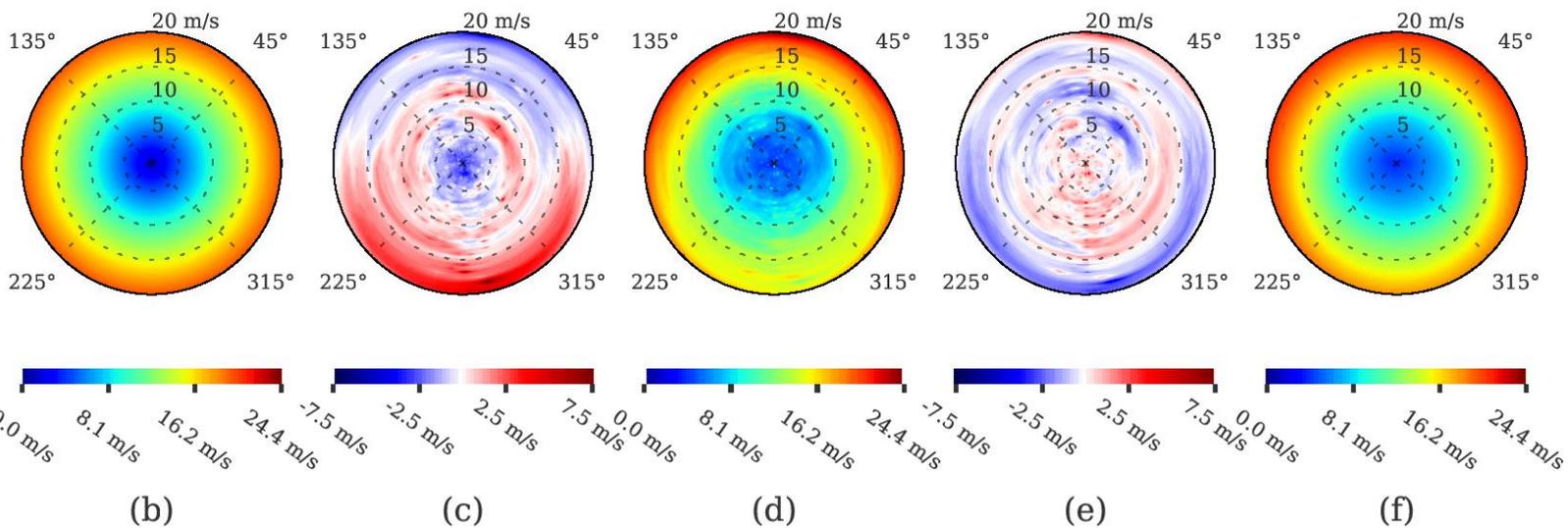
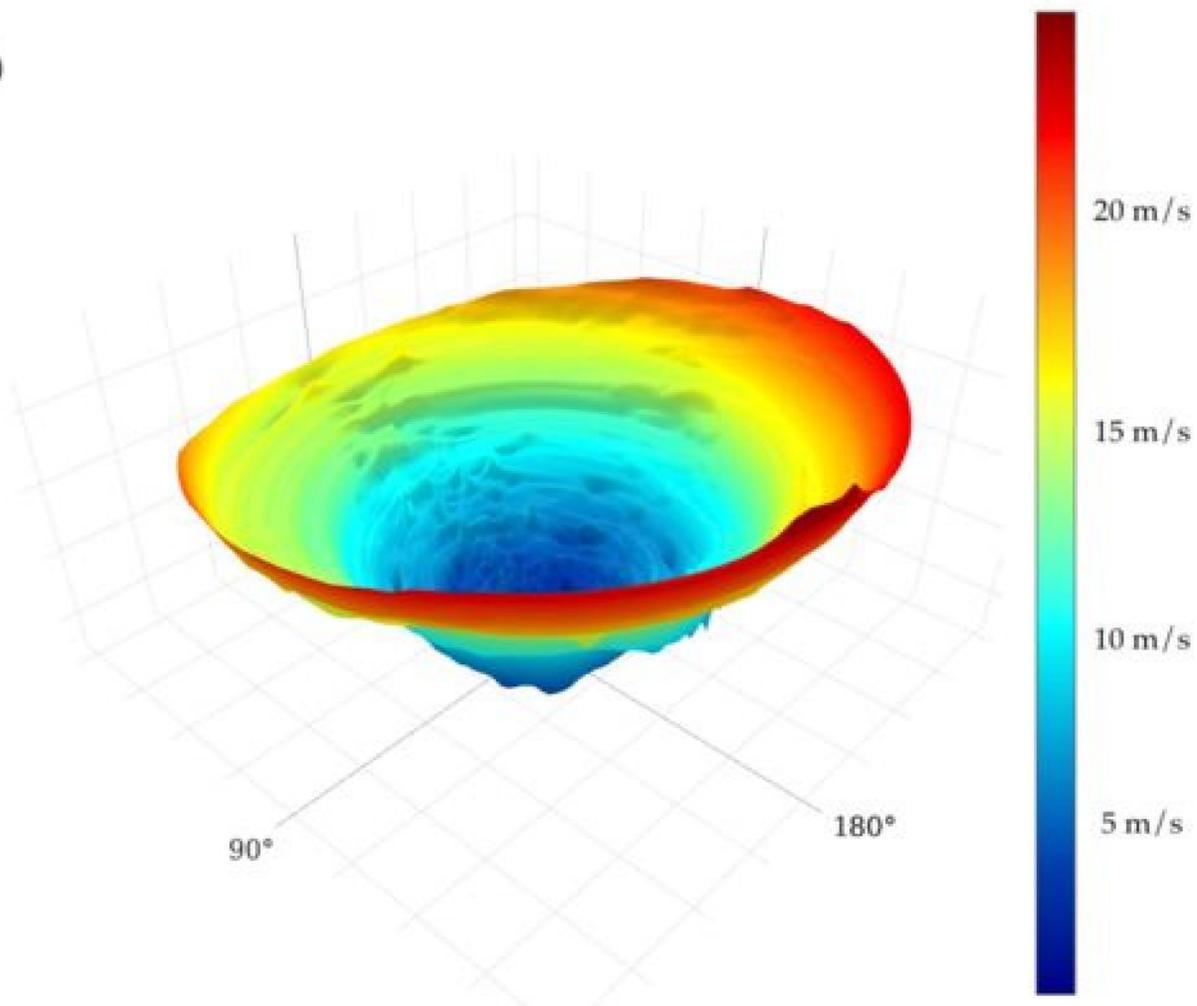
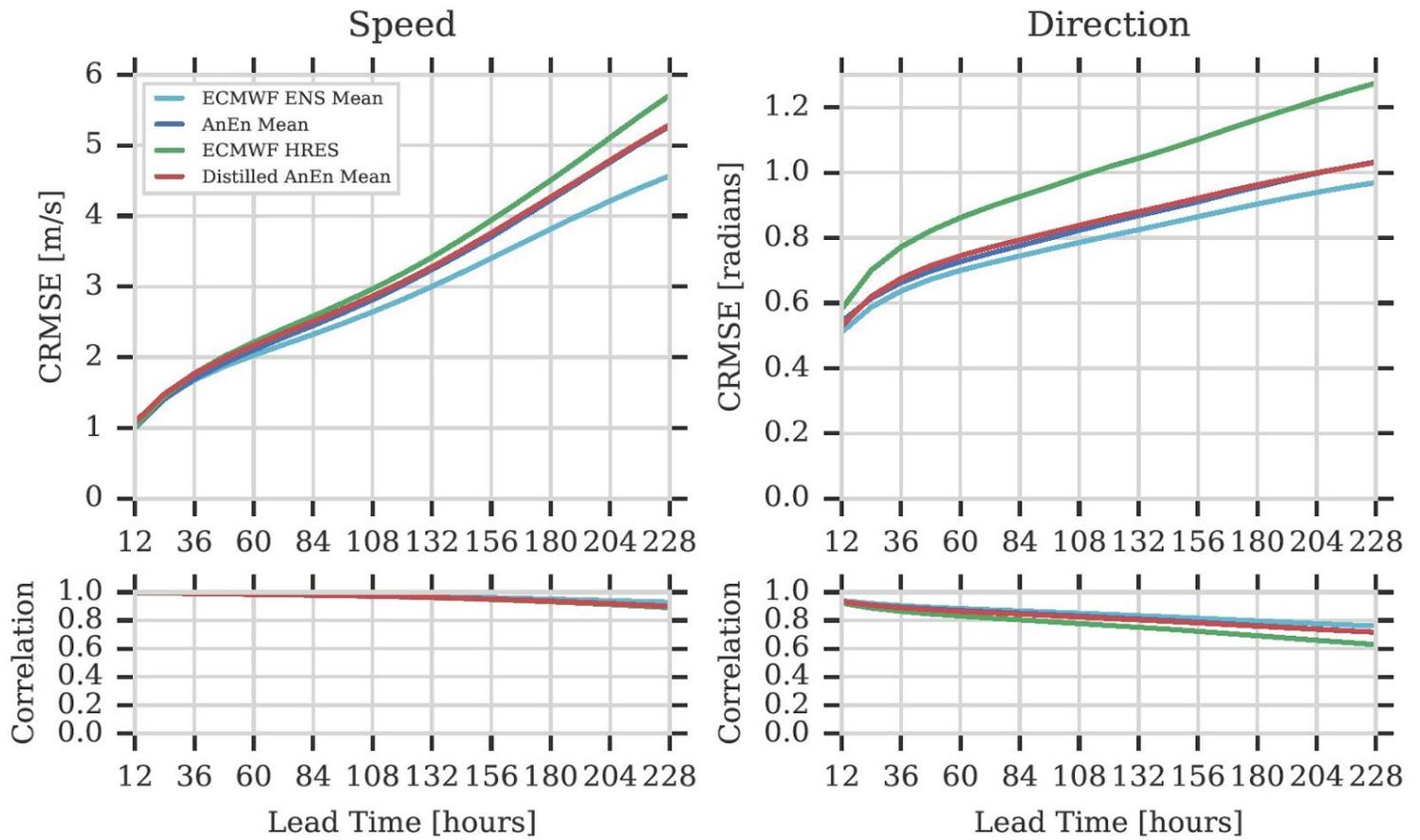
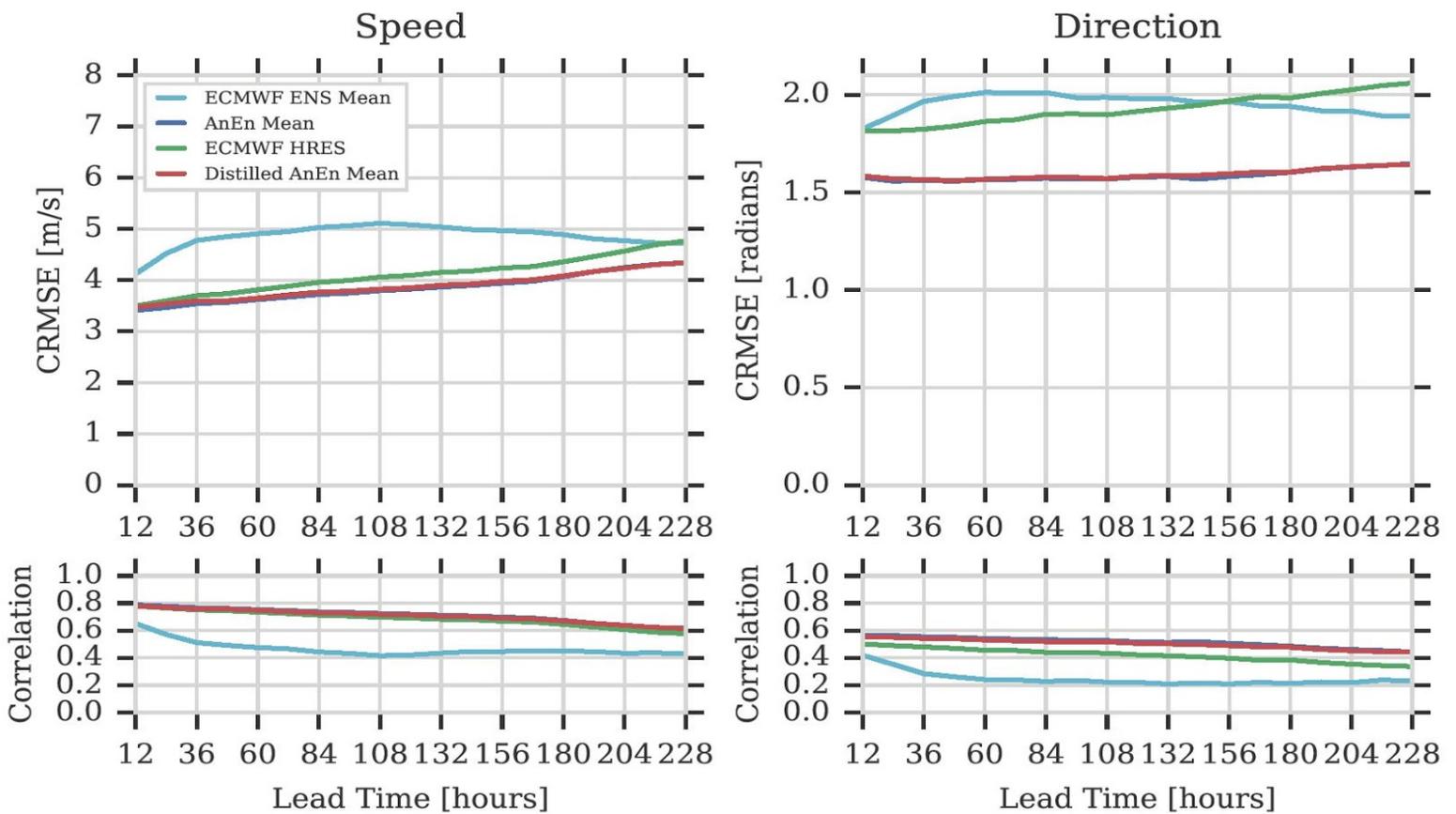


Figure 2.

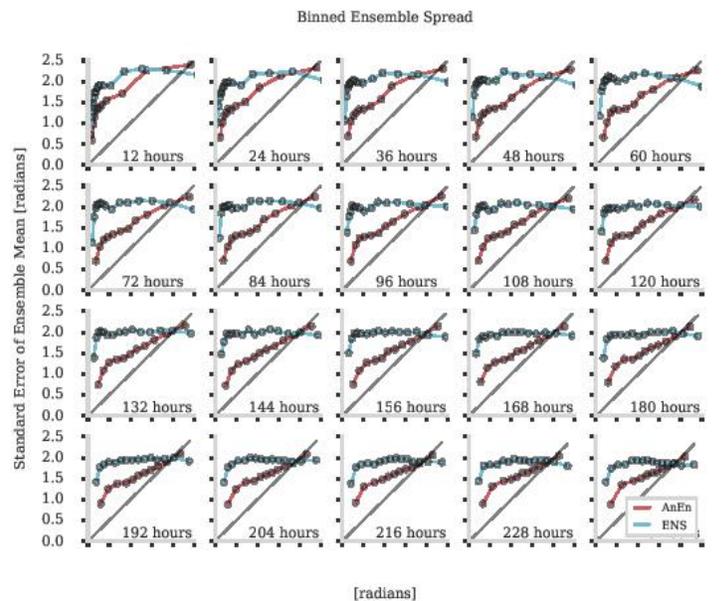
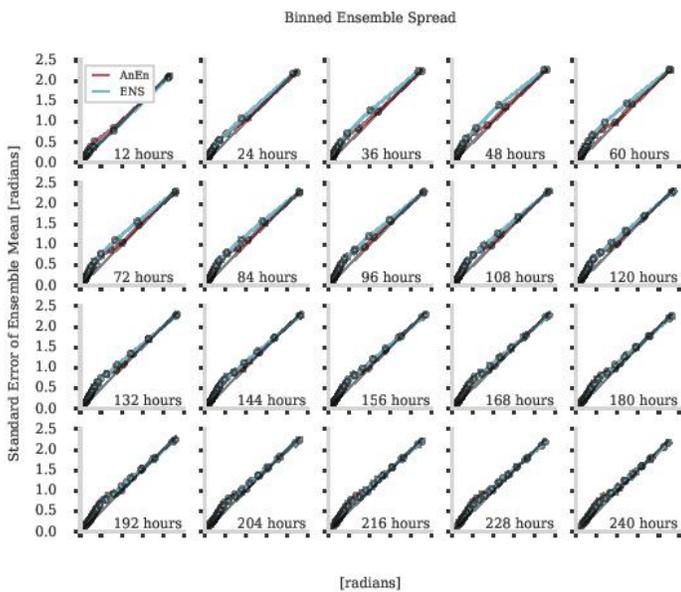
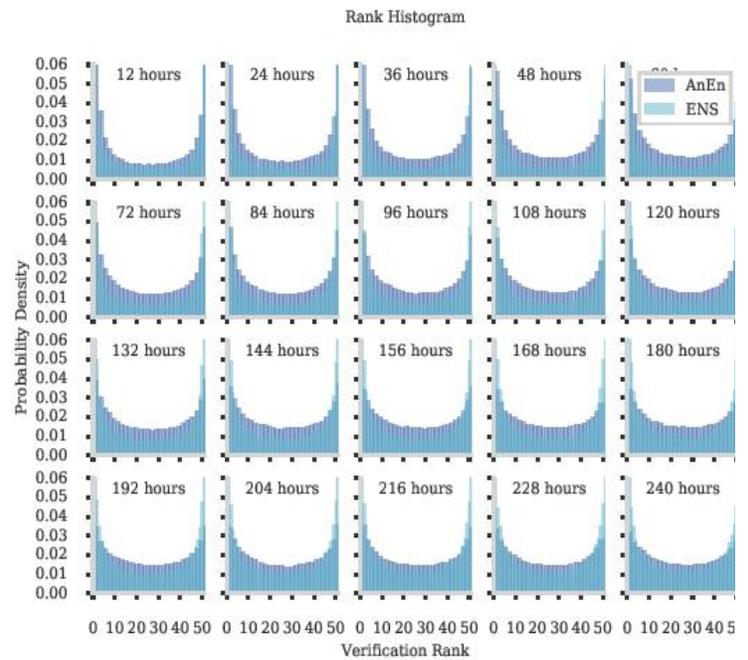
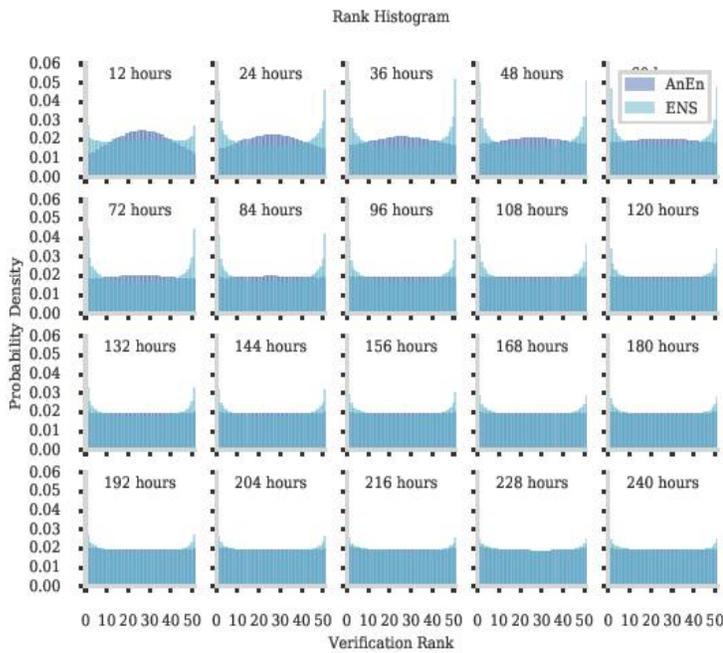
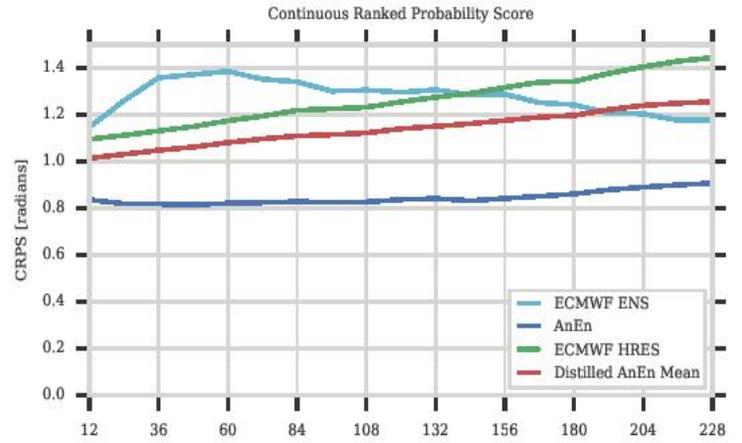
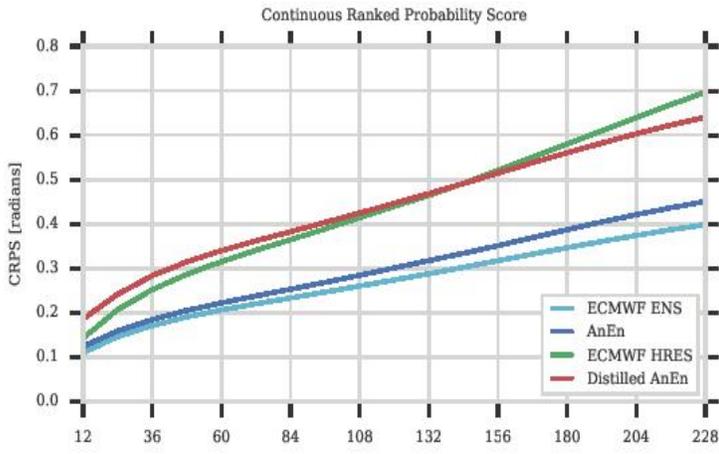


(a)



(b)

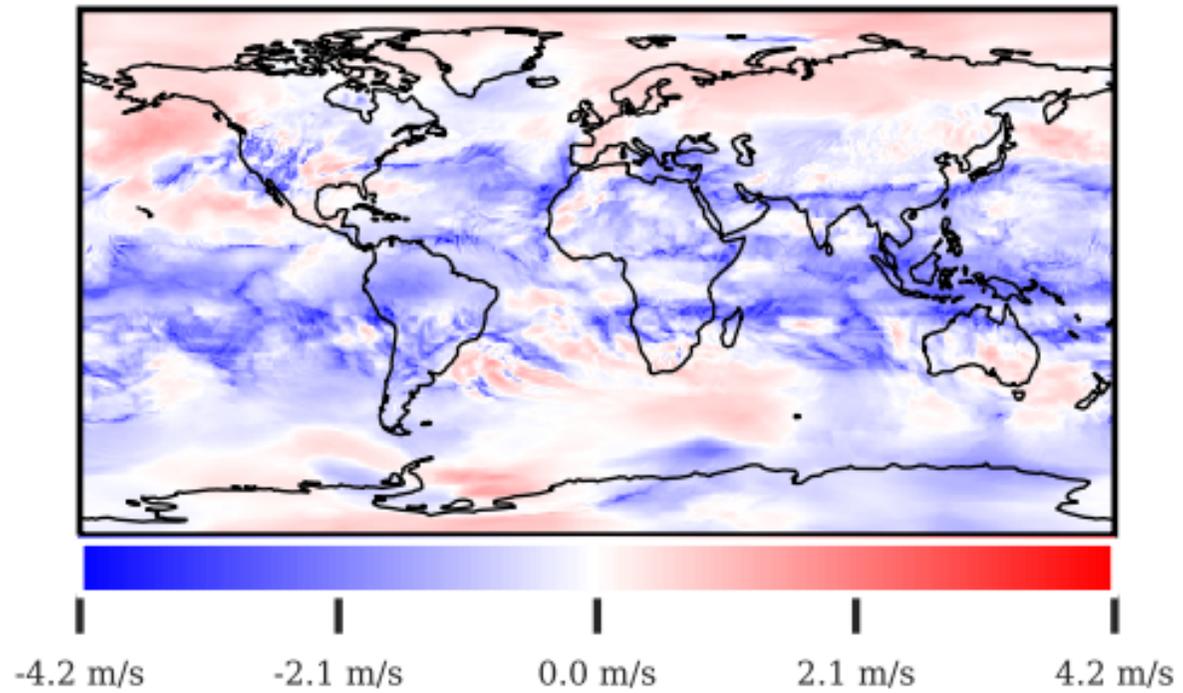
Figure 3.



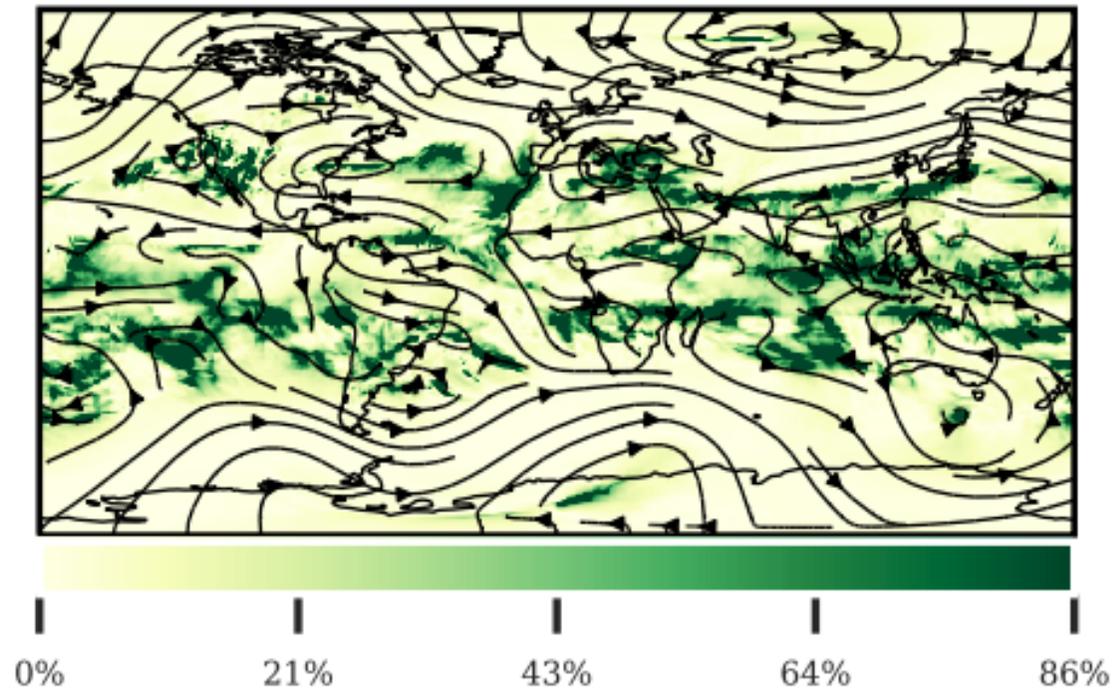
(a)

(b)

Figure 4.



(a)



(b)

# Supporting Information for Improving Wind Forecasts in the Lower Stratosphere by Distilling an Analog Ensemble into a Deep Neural Network

Salvatore Candido<sup>1</sup>, Aakanksha Singh<sup>1</sup>, and Luca Delle Monache<sup>1,2</sup>

<sup>1</sup>Loon, Mountain View, California, USA

<sup>2</sup>Center for Western Weather and Water Extremes, La Jolla, California, USA

## Contents of this file

1. Computing the Conventional Analog Ensemble with MapReduce
2. Details of Training a Distilled Analog Ensemble Model
3. Statistics of the Loon Data Set
4. Probabilistic Evaluation Metrics for Wind Speed
5. Confidence Intervals on Deterministic Evaluations
6. Algorithm Skill Comparison By Geography
7. Results for an Earlier Validation Period

## Introduction

The supplemental data here covers a few disparate sets of material. The first two sections provide additional technical detail on the methods of the paper that are not necessary to understand the approach, but are useful when replicating the results. We then share statistical breakdowns of the Loon observations which can also be derived by

---

processing the data set at (Candido, 2020), but we include for convenience here. The following three sections contain views of our results that we do not include in the main text, but are of interest to some readers of early drafts of this paper. Finally, we present an additional validation set that leads to similar conclusions as the results in the main text. We include this for completeness as some of our early discussions of this work used this validation set, rather than the newer validation set that allows us to compare against the European Centre for Medium-Range Weather Forecasts (ECMWF) ensemble system (ENS).

### **Computing the Conventional Analog Ensemble with MapReduce**

A barrier to operationalizing a global analog ensemble (AnEn) system is processing the corpus of analogs, which can easily grow to 100's of terabytes of data for three-dimensional global predictions over several years. The AnEn algorithm provides a natural partitioning as execution is independent for each  $\mathcal{P}$  (grid point and lead time). Every historical forecast and the current prediction contain a piece of data for every grid point. The challenge is to organize the data so that the calculations can be efficiently executed across many datacenter computers.

Our approach is to use the MapReduce paradigm (Dean & Ghemawat, 2004), which allows the computation to run on a distributed computing (cloud) infrastructure like Google's Flume (Chambers et al., 2010). The idea is to break the computation into two subsequent **Map** and **Reduce** phases, each of which operate many times in parallel on different portions of the data and output key-value pairs. Once written this way, the framework can handle scheduling the program's execution across many machines and moving the various subsets of data to the appropriate machines.

The above procedure is accomplished as follows. Let latitude, longitude, pressure, and forecast lead time be the tuple  $k$ , which is unique for each grid point. The first **Map** phase scans all historical forecast files and generates a key-value pair  $(k, t_f) \rightarrow x^i$  for each grid point and  $(k, t_f) \rightarrow y^i$  for each analysis point. The  $k$  corresponds to the location and forecast lead time for the particular  $x^i$  (past forecast) and  $t_f$  (calendar time being forecast). For every  $y^i$  we generate multiple key-value pairs corresponding to a  $k$  for *every* lead time the system will forecast.

Notice that the above rule gives each forecast-observation pair a unique key. Prior to the reduce phase all identical keys are grouped into one **Reduce** call by the MapReduce framework. The **Reduce** phase joins these  $x^i$  and  $y^i$  pairs into a single record and saves them as new key-value pairs  $k \rightarrow (x^i, y^i)$ . This first MapReduce gives us a historical corpus. This corpus could be built in advance of receiving a new forecast to post-process.

The second MapReduce groups the data by grid point and runs the AnEn algorithm. A **Map** phase on the forecast data file from the ECMWF generates key-value pairs  $k \rightarrow x^f$  for each grid point. The historical corpus key-value pairs are used directly. Note that  $x^f$  and every set of candidate analogs  $\{(x^i, y^i)\}$  for a grid point have the same key. The data is grouped by key and fed to the **Reduce** phase that has all the data needed to apply equations (1) and (3) to generate the forecast.

### **Details of Training a Distilled Analog Ensemble Model**

This section describes the low-level technical details of the distillation process.

Our training corpus is prepared by using a MapReduce similar to what is described in the previous section to process the set of forecast data files (both the 00Z and 12Z epochs) archived during the training period. Rather than running the AnEn algorithm

logic, i.e., equations (1) and (3), at each grid point given a new forecast, we instead save the candidate analogs (forecast-observation pairs) for a given grid point in a single record. These records are stored together on disk for retrieval by the training system, i.e., we have a set of records where each record corresponds to a unique latitude, longitude, pressure, and lead time grid point and contains all the viable forecast-observation pairs at this location at the appropriate lead time.

This data set is used to feed the training process of our deep neural network (DNN). We use a distributed architecture. We train our DNN based on the output of the AnEn, not directly from these forecast-observation pairs. Thus, we need to sample a hypothetical forecast to generate a training example of an input-output pair for the AnEn system. We use 10 datacenter worker processes that sample uniformly among grid points in records on disk, sample a hypothetical wind speed and wind heading forecast, and construct the input to the DNN (corresponding to this grid point and forecast) and the output (of the AnEn algorithm). In the results presented in this paper, we sample heading uniformly and wind speed from a beta distribution with  $\alpha = 1.2$ ,  $\beta = 3$ , and a coefficient of 100. Effectively this creates a weighted distribution of wind speeds which seems generally applicable to the pressure altitudes ranges of interest in the stratosphere.

Unlike many applications, we do not simply loop over the dataset on disk a fixed number of times or until training error stabilizes. This is because every time we touch a record we sample a new forecast and generate a new input-output pair. Saving these pairs on disk versus generating them online during training is an engineering trade-off, and we have chosen the latter approach.

These examples from the 10 worker processes are injected into a reservoir datacenter process, whose job is essentially to receive new examples, store them in a limited size buffer, and respond to requests (from the learning process) for samples. Rather than choosing a circular buffer or some other first-in, first-out structure, we use a flat data array of 1 million examples and, for each new example, sample an index in the array at random to replace. This means some examples will persist in the buffer longer, and some for a shorter period of time. The typical dwell time of an example in the buffer can be characterized probabilistically. The learning process repeatedly queries the reservoir for batches of training examples, which are selected uniformly at random from examples in the data array. A slowly changing flow of examples where each batch (on average) tends to be drawn from disparate parts of the function mapping being learned is conceptually similar to the replay buffer in deep reinforcement learning (Lin, 1992; Mnih et al., 2015).

We use the Tensorflow (Abadi et al., 2015) library to create and train our DNN. Our network has inputs of latitude, longitude, pressure altitude, forecast lead time, forecast direction, and forecast speed. We transform these into a graph layer that is normalized using the following code snippet where the array ‘domain’ represents the inputs described above.

```
nlat = tf.multiply(domain[:, 0], 1. / 90.0)
coslng = tf.cos(tf.multiply(domain[:, 1], np.pi / 180.0))
sinlng = tf.sin(tf.multiply(domain[:, 1], np.pi / 180.0))
npre = tf.multiply(tf.subtract(domain[:, 2], 4799.), 1. / (14432. - 4799.))
nlea = tf.multiply(tf.subtract(domain[:, 3], 43200.), 1. / (864000. - 43200.))
coshead = tf.cos(domain[:, 4])
```

```
sinhead = tf.sin(domain[:, 4])  
nspeed = tf.multiply(domain[:, 5], 1. / 100.)  
normalized_domain = tf.transpose(  
    tf.stack([nlat, coslng, sinlng, npre, nlea, coshead, sinhead, nspeed]))
```

We do this to avoid the discontinuity in longitude being present in our DNN, and to make our inputs have roughly the same order of magnitude (which is a domain trick to decrease training time).

At this point the network consists of 10 fully-connected hidden layers with ReLu activation functions. Each layer has a width of 50 elements. These plus an ultimate layer containing the ultimate post-processed speed and heading forecast (width 2, fully-connected, no activation function) comprise the trained layers of the DNN. We can also include additional network outputs such as forecast uncertainty (standard deviation of the forecast) or ensemble members. This is not discussed in this paper.

To train the network we use stochastic gradient descent with a learning rate of 0.0001 and batch size 100. We train until the root mean square error between the DNN forecasts and the AnEn mean forecasts (from the training examples) stabilizes. In the distilled AnEn used to generate the results in this paper, we trained the DNN with about 6 billion examples.

This network architecture was not tuned for efficiency, but instead chosen to demonstrate how a fairly standard and basic deep learning approach could be used to implement this algorithm.

## Statistics of the Loon Data Set

The following plots show the distribution of Loon’s approximately 10.5 million observations used for one of the comparisons between algorithms shown in the main text of the paper. This is the intersection of Loon’s dataset of observations of stratospheric winds from Loon (<http://www.loon.com>) high altitude balloons (Candido, 2020) and the region and time period for the validation paper used in our study.

Figure S1 shows the distribution of the data over pressure altitude and latitude.

Figure S2 shows the geographical distribution of the data.

### **Probabilistic Evaluation Metrics for Wind Speed**

In the main text of the paper we presented the CRPS, Spread Skill, and Rank Histogram plots for comparing the ensemble systems predictions on wind direction, an omitted plots for wind speed given a similar pattern on skill between the approaches. We include the figures for wind speed in Figure S3.

### **Confidence Intervals on Deterministic Evaluations**

Figures S4 and S5 show the same data as in Figure 2(a) in the main text, but include box plot views of the 90% bootstrap confidence intervals.

### **Algorithm Skill Comparison By Geography**

Figure S6 show the CRMSE averaged across all lead times grouped by geography. One can observe that the Distilled AnEn has higher skill (lower CRMSE) than the baseline ECMWF HRES generally across the stratosphere globally.

### **Results for an Earlier Validation Period**

Our original analysis of the methods included a comparison of AnEn mean against the ECMWF high-resolution deterministic forecast (HRES) and, for the probabilistic predictions, against a persistence ensemble (PeEn) over a year long validation period

from October, 2017 to September, 2018. However, to add a comparison to the ECMWF ENS in a revised version of the manuscript, we changed our validation period in the main text due to data availability.

The PeEn is a simple way to generate an ensemble that consists of selecting the last  $N$  available ground-truth values to generate an  $N$ -member ensemble. It has been used in, e.g., Alessandrini, Delle Monache, Sperati, and Cervone (2015) and Cervone, Clemente-Harding, Alessandrini, and Monache (2017), as a probabilistic baseline forecast and can be interpreted as the probabilistic extension of a deterministic persistence forecast.

For the results shown in in Figures S7 and S8 the training dataset is the HRES forecasts produced from July, 2016, to September, 2017. We use this to choose weights used in the analog matching process. The validation period is over the HRES forecasts produced from October, 2017, to September, 2018. The data available in the AnEn matching includes all the forecasts in the training dataset plus any additional forecasts between the beginning of the validation time period but prior to the current forecast. This simulates operational use of an AnEn system. To evaluate the distilled AnEn we only use a DNN distilled from the training dataset.

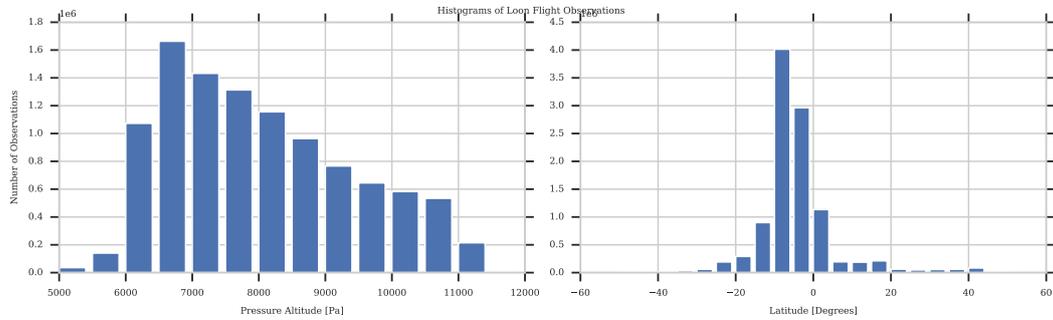
Please refer to the main text of the paper where the relevance of the metrics shown in the below figures are explained in greater detail, albeit for a different validation period.

## References

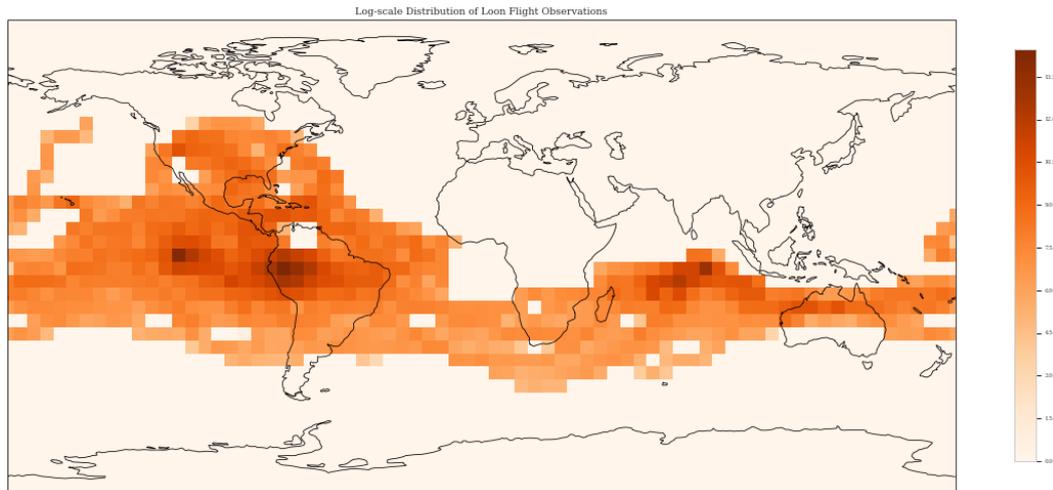
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2015). *TensorFlow: Large-scale machine learning on heterogeneous systems*. Retrieved from <https://www.tensorflow.org/> (Software available from tensorflow.org)

- Alessandrini, S., Delle Monache, L., Sperati, S., & Cervone, G. (2015). An analog ensemble for short-term probabilistic solar power forecast. *Applied Energy*, *157*, 95–110. Retrieved 2016-11-21, from <http://www.sciencedirect.com/science/article/pii/S0306261915009368>
- Candido, S. (2020, April). *Loon stratospheric sensor data*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.3763022> doi: 10.5281/zenodo.3763022
- Cervone, G., Clemente-Harding, L., Alessandrini, S., & Monache, L. D. (2017). Short-term photovoltaic power forecasting using artificial neural networks and an analog ensemble. *Renewable Energy*, *108*, 274 - 286. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0960148117301386> doi: <https://doi.org/10.1016/j.renene.2017.02.052>
- Chambers, C., Raniwala, A., Perry, F., Adams, S., Henry, R., Bradshaw, R., & Nathan. (2010). Flumejava: Easy, efficient data-parallel pipelines. In *Acm sigplan conference on programming language design and implementation (pldi)* (p. 363-375). 2 Penn Plaza, Suite 701 New York, NY 10121-0701. Retrieved from <http://dl.acm.org/citation.cfm?id=1806638>
- Dean, J., & Ghemawat, S. (2004). Mapreduce: Simplified data processing on large clusters. In *Osd'04: Sixth symposium on operating system design and implementation* (pp. 137–150). San Francisco, CA.
- Lin, L.-J. (1992). Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, *8*(3-4), 293–321.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... others (2015). Human-level control through deep reinforcement learning. *Nature*,

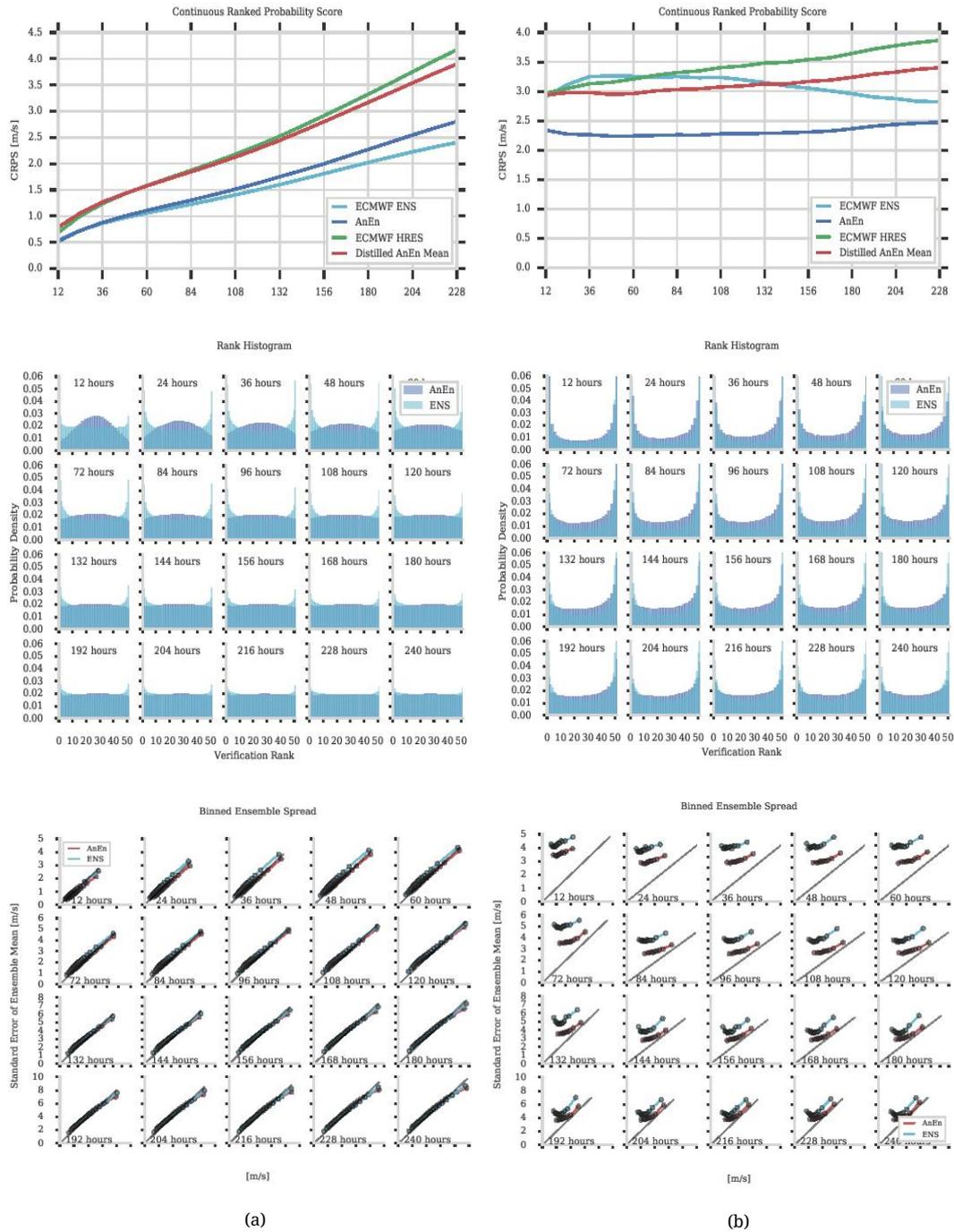
518(7540), 529.



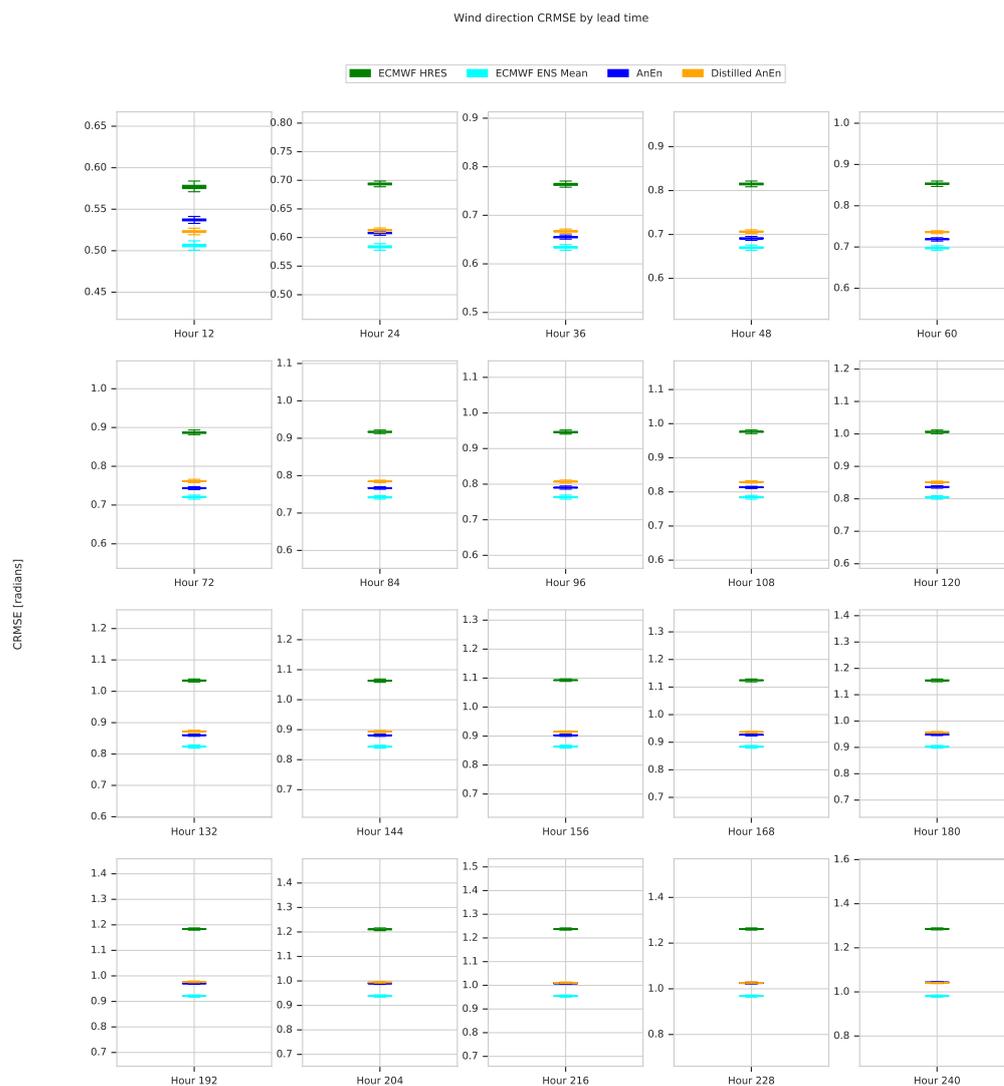
**Figure S1.** Distribution of Loon’s measurements as a function of pressure altitude and latitude.



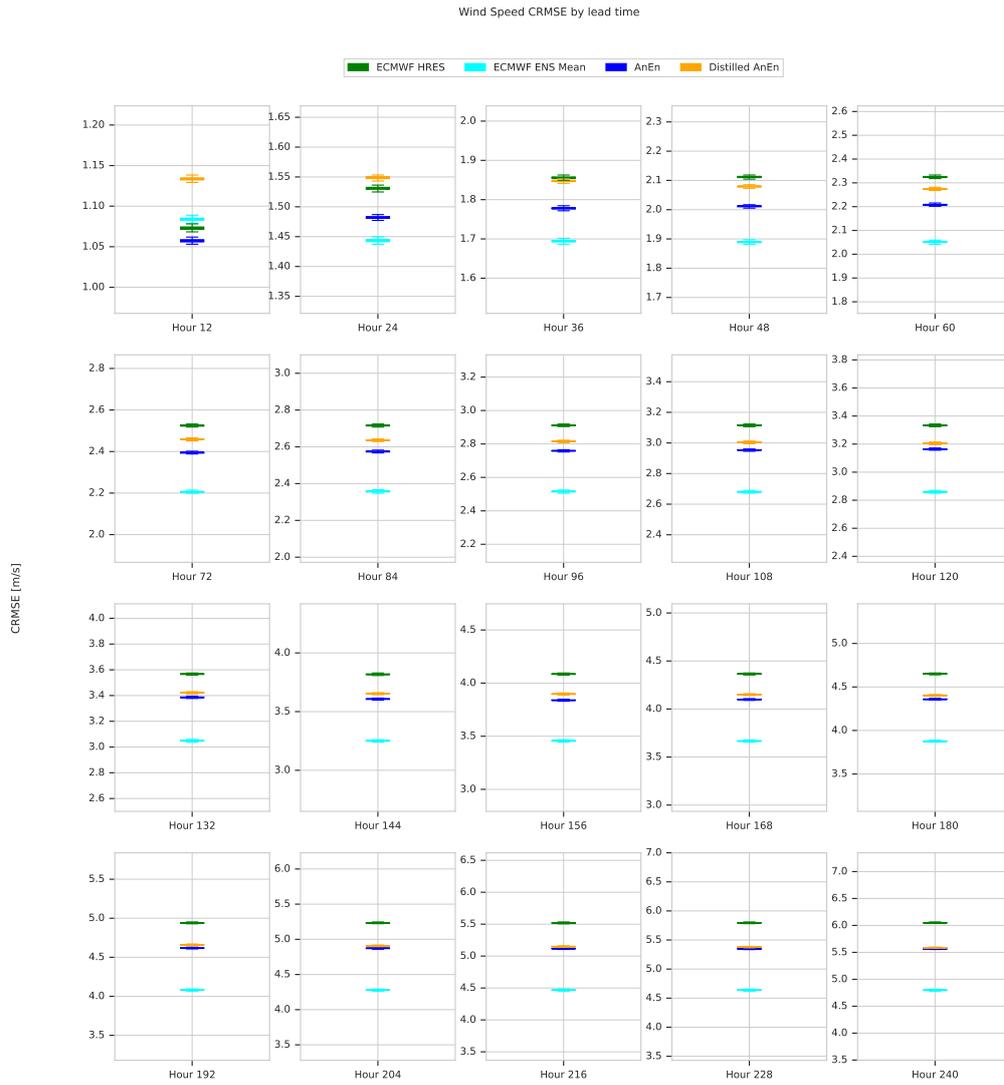
**Figure S2.** Geographical distribution of Loon’s measurements.



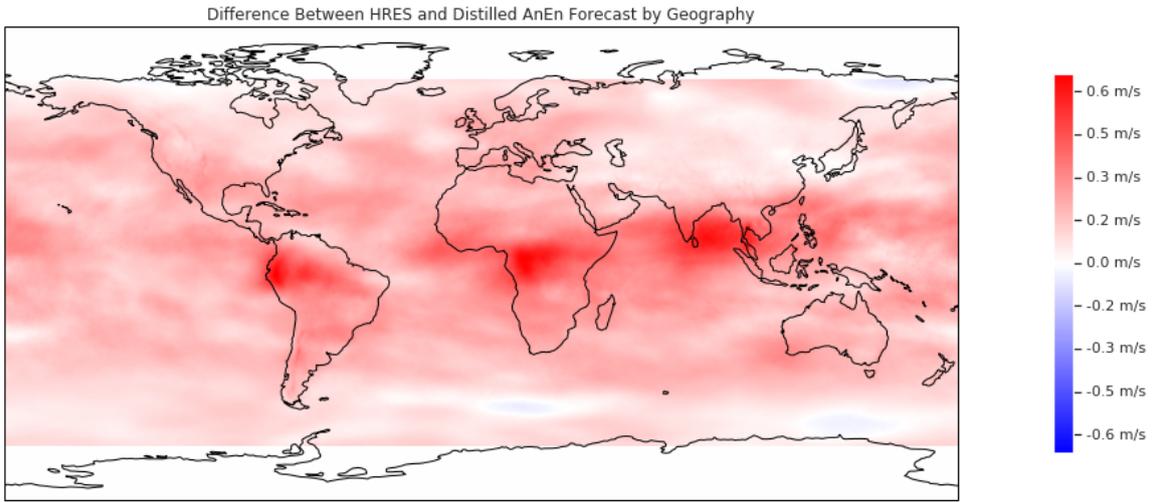
**Figure S3.** Probabilistic forecast evaluation metrics comparing the AnEn forecast of wind speed to forecasts produced by a ENS. Results with HRES analysis as ground truth are shown on the left (a), while results against Loon’s measurements are on the right (b). From top to bottom, the metrics shown are CRPS, rank histogram, and binned-spread skill.



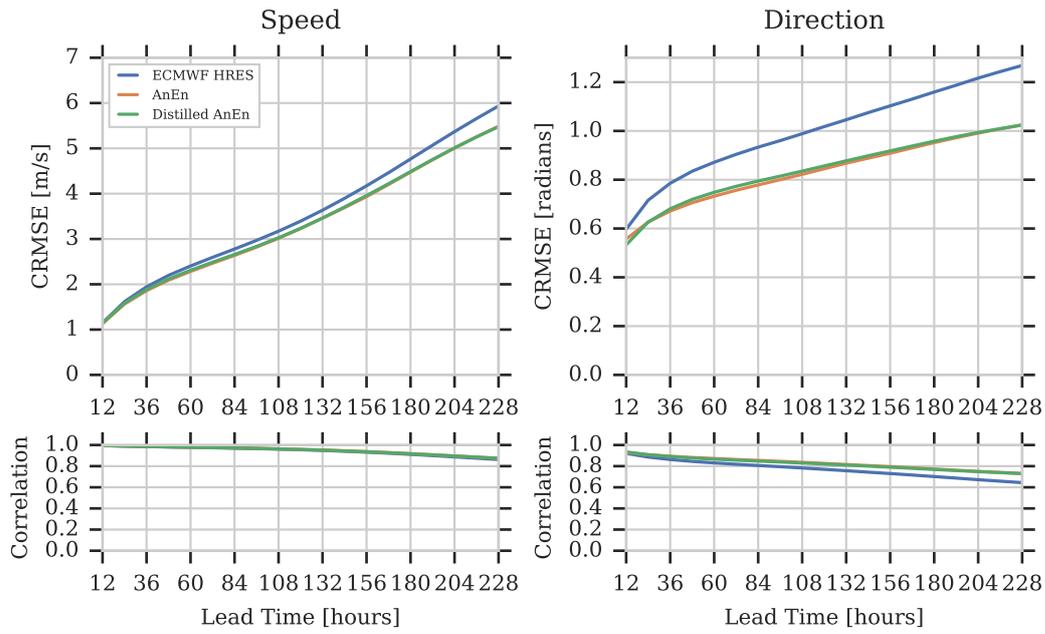
**Figure S4.** CRMSE for wind direction predictions including boxplots showing the bootstrap 90% confidence intervals.



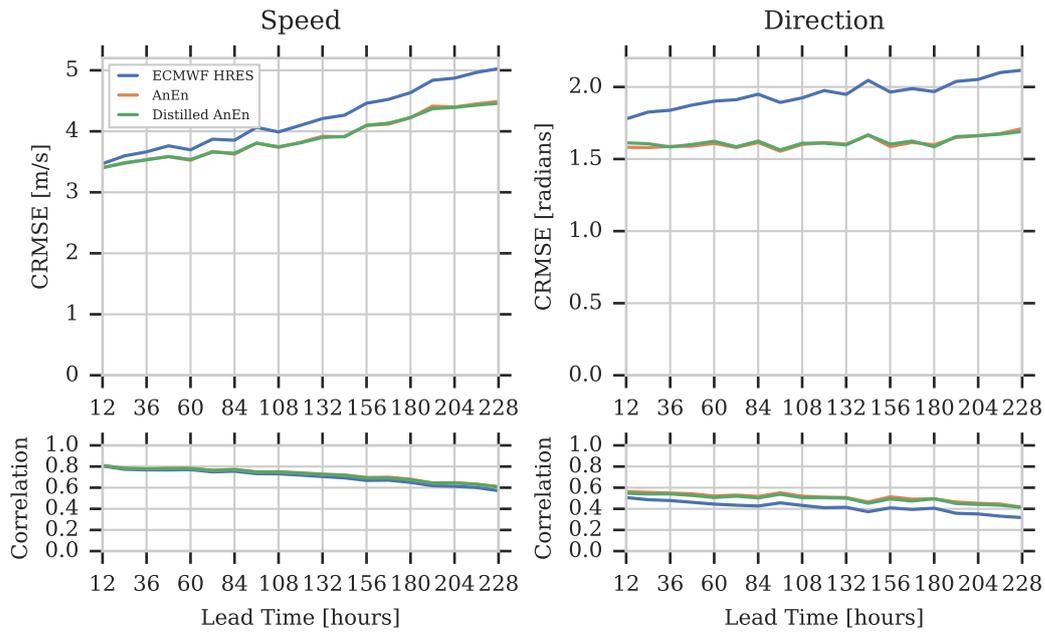
**Figure S5.** CRMSE for wind speed predictions including boxplots showing the bootstrap 90% confidence intervals.



**Figure S6.** Geographical distribution of CRMSE for the distilled AnEn prediction of wind speed with HRES analysis as ground truth.

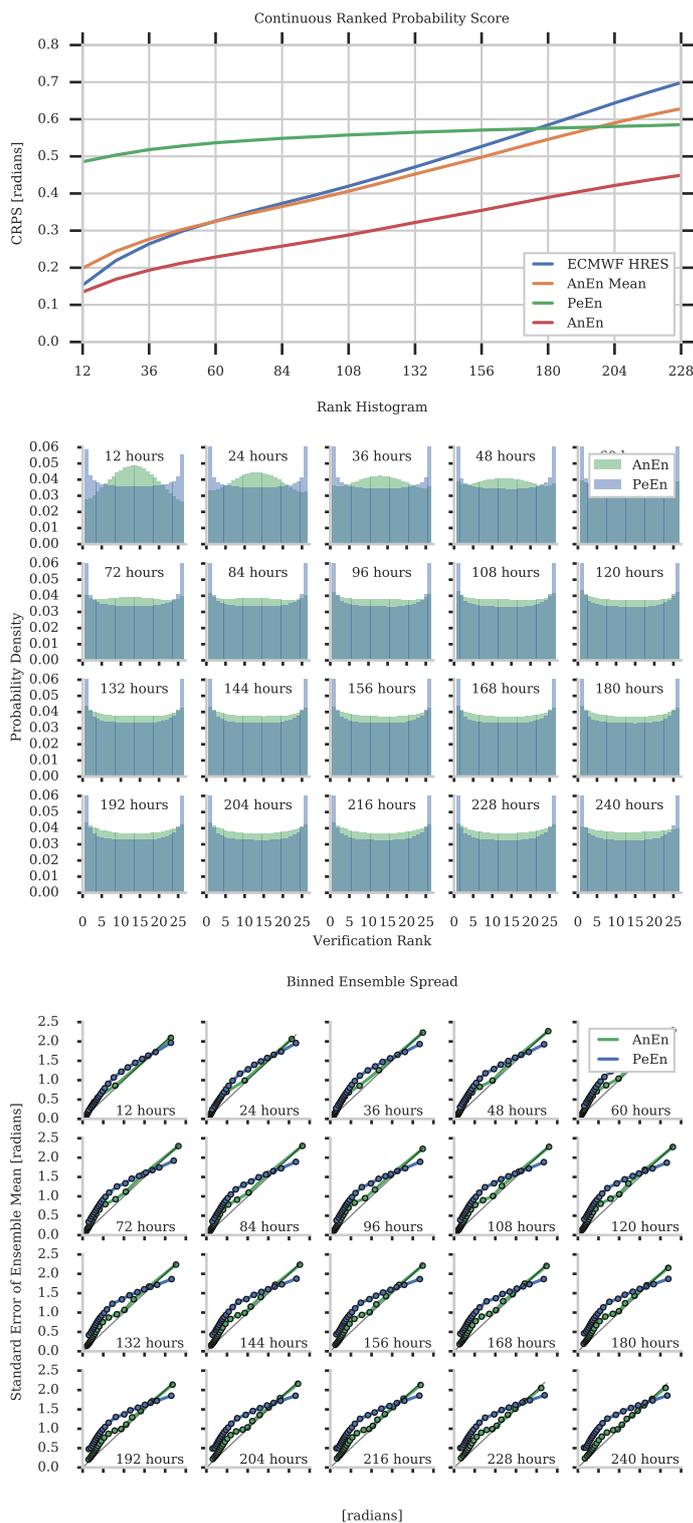


(a)



(b)

**Figure S7.** A deterministic wind speed and direction forecast skill comparison between the HRES, AnEn, and Distilled AnEn over all lead times is shown using as ground truth (a) HRES analysis and (b) Loon observations of stratospheric winds.



**Figure S8.** Probabilistic forecast evaluation metrics comparing the AnEn forecast of wind direction to forecasts produced by HRES, AnEn mean, and PeEn using HRES analysis as the ground truth. From top to bottom, the metrics shown are CRPS, rank histogram, and binned-spread skill.