| 1 | Machine Learning Enhancement of Storm Scale Ensemble                   |
|---|--|
| 2 | Probabilistic Quantitative Precipitation Forecasts                     |
| 3 | David John Gagne II*   |
|   | School of Meteorology, University of Oklahoma, Norman, Oklahoma        |
| 4 | Amy McGovern   |
|   | School of Computer Science, University of Oklahoma, Norman, Oklahoma   |
| 5 | Ming Xue   |
|   | Center for the Analysis and Prediction of Storms/School of Meteorology |
|   | University of Oklahoma, Norman, Oklahoma                               |
|   |  |
|   |  |

<sup>\*</sup>*Corresponding author address:* David John Gagne II, 120 David L. Boren Blvd., Suite 5900, Norman, OK 73072. Email: djgagne@ou.edu

#### ABSTRACT

Probabilistic quantitative precipitation forecasts challenge meteorologists due to the wide 7 variability of precipitation amounts over small areas and their dependence on conditions at 8 multiple spatial and temporal scales. Ensembles of convection-allowing numerical weather 9 prediction models offer a way to produce improved precipitation forecasts and estimates 10 of the forecast uncertainty. These models allow for the prediction of individual convective 11storms on the model grid, but they often displace the storms in space, time, and inten-12 sity, which results in added uncertainty. Machine learning methods can produce calibrated 13 probabilistic forecasts from the raw ensemble data that correct for systemic biases in the 14 ensemble precipitation forecast and incorporate additional uncertainty information from ag-15 gregations of the ensemble members and additional model variables. This study utilizes the 16 2010 Center for the Analysis and Prediction of Storms Storm Scale Ensemble Forecast system 17 and the National Severe Storms Laboratory National Mosaic and Multisensor Quantitative 18 Precipitation Estimate as input data for training logistic regressions and random forests to 19 produce a calibrated probabilistic quantitative precipitation forecast. The reliability and 20 discrimination of the forecasts are compared through verification statistics and a case study. 21

## <sup>22</sup> 1. Introduction

Most flooding fatalities occur due to flash floods, in which waters rise and fall rapidly 23 due to concentrated rainfall over a small area (Ashley and Ashley 2008). The first step 24 to anticipating flash floods is quantitative forecasting of the amount, location, and tim-25 ing of precipitation. These quantitative precipitation forecasts are challenging due to the 26 wide variability of precipitation amounts over small areas, the dependence of precipitation 27 amounts on processes at a wide range of scales, and the dependence of extreme precipitation 28 on any precipitation actually occurring (Bremnes 2004; Ebert 2001; Doswell et al. 1996). 29 Recent advances in numerical modeling and machine learning are working to address these 30 challenges. 31

Numerical Weather Prediction (NWP) models are now being run experimentally at 4 km 32 horizontal grid spacing, or storm scale, allowing for the formation of individual convective 33 cells without a convective parameterization scheme. These models better represent storm 34 processes and output hourly predictions, but they have the challenge of correctly placing 35 and timing the precipitation compared to models with coarser grid spacing and temporal 36 resolution. An ensemble of storm scale NWP models can provide improved estimates of un-37 certainty compared to coarser ensembles due to better sampling of the spatiotemporal errors 38 associated with individual storms (Clark et al. 2009). As each ensemble member produces 39 predictions of precipitation and precipitation ingredients, the question then becomes how 40 best to combine those predictions into the most accurate and useful consensus guidance. 41

The final product should highlight areas most likely to be impacted by heavy rain and
 also provide an uncertainty estimate for that impact. Probabilistic Quantitative Precipi-

tation Forecasts (PQPFs) incorporate both of these qualities. A good PQPF should have 44 reliable probabilities, such that a 40% chance of rain verifies 40% of the time over a large 45 sample (Murphy 1977). PQPFs should also discriminate between extreme and trace precipi-46 tation events consistently, so most extreme events occur with higher probabilities, and most 47 trace precipitation events are associated with low probabilities. Since individual ensemble 48 members may be biased in different situations, a simple count of the ensemble members that 49 exceed a threshold will often result in unreliable forecasts that discriminate poorly. Incor-50 porating trends from past forecasts and additional information from other model variables 51 can offset these biases and produce an enhanced PQPF. 52

Ensemble weather prediction (Toth and Kalnay 1993; Tracton and Kalnay 1993; Molteni 53 et al. 1996) has required various forms of statistical post-processing to produce accurate pre-54 cipitation forecasts and uncertainty estimates from the ensembles, but most previous studies 55 used coarser ensembles and longer forecast windows for calibration. The rank histogram 56 method (Hamill and Colucci 1997) showed that ensemble precipitation forecasts tended to 57 be underdispersive. Linear regression calibration methods have shown some skill improve-58 ments in Hamill and Colucci (1998). Eckel and Walters (1998). Krishnamurti et al. (1999). 59 and Ebert (2001). Hall et al. (1999), Koizumi (1999), and Yuan et al. (2007) applied neural 60 networks to precipitation forecasts and found increases in performance over linear regression. 61 Logistic regression, a transform of a linear regression to fit an S-shaped curve ranging from 0 62 to 1, has shown more promise, as in Applequist et al. (2002), which tested linear regression, 63 logistic regression, neural networks, and genetic algorithms on 24-hour PQPF and found 64 that logistic regression consistently outperformed the other methods. Hamill et al. (2004, 65 2008) also utilized the logistic regression with an extended training period for added skill. 66

Storm scale ensemble precipitation forecasts have been post-processed with smoothing 67 algorithms that produce reliable probabilities within longer forecast windows. Clark and 68 Coauthors (2011) applied smoothing algorithms at different spatial scales to the SSEF pre-69 cipitation forecasts and compared verification scores. Johnson and Wang (2012) compared 70 the skill of multiple calibration methods on neighborhood and object-based probabilistic 71 forecasts from the 2009 SSEF. Marsh et al. (2012) applied a Gaussian kernel density es-72 timation function to the NSSL 4 km Weather Research and Forecasting (WRF) to derive 73 probabilities from deterministic forecasts. These methods are helpful for predicting larger 74 scale events but smooth out the threats from extreme precipitation in individual convective 75 cells. Gagne II et al. (2012) took the first step in examining how multiple machine learning 76 approaches performed in producing probabilistic, deterministic, and quantile precipitation 77 forecasts over the central United States at individual grid points. 78

The purpose of this paper is to analyze PQPF predictions produced by multiple machine 79 learning techniques incorporating data from the Center for the Analysis and Prediction of 80 Storms (CAPS) 2010 Storm Scale Ensemble Forecast (SSEF) system (Xue et al. 2011; Kong 81 et al. 2011). In addition to the choice of algorithm, some variations in the algorithm setups 82 are also examined. The strengths and weaknesses of the machine learning algorithms are 83 shown through the analysis of verification statistics, variables chosen by the machine learning 84 models, and a case study. This paper expands on the work presented in Gagne II et al. (2012) 85 by including statistics for the eastern US, a larger training set more representative of the 86 precipitation probabilities, a different case study day with comparisons of multiple runs, 87 multiple precipitation thresholds, and more physical justification for the performance of the 88 machine learning algorithms. 89

### 90 2. Data

#### 91 a. Ensemble Data

The Center for the Analysis and Prediction of Storms (CAPS) 2010 Storm Scale Ensem-92 ble Forecast (SSEF) system (Xue et al. 2011; Kong et al. 2011) provides the input data for 93 the machine learning algorithms. The 2010 SSEF consists of 19 individual model members 94 from the WRF Advanced Research WRF (ARW), 5 members from the WRF Nonhydrostatic 95 Mesoscale Model (NMM), and 2 members from the CAPS Advanced Research Prediction 96 System (ARPS; Xue et al. 2000, 2001, 2003). Each member has a varied combination of 97 microphysics schemes, land surface models, and planetary boundary layer schemes. The 98 SSEF ran every weekday at 0000 UTC in support of the 2010 National Oceanic and Atmo-99 spheric Administration/Hazardous Weather Testbed Spring Experiment (Clark et al. 2012), 100 which ran from May 3 to June 18, for a total of 34 runs. The SSEF provides hourly model 101 output over the contiguous United States at 4 km horizontal grid spacing out to 30 hours. 102 Of the 26 members, the 14 members included initial condition perturbations derived from 103 the National Centers for Environmental Prediction Short Range Ensemble Forecast (SREF: 104 Du et al. 2006) members, and only they are included in our post-processing procedure. The 105 12 other members used the same control initial conditions although with different physics 106 options and models; designed to examine forecast sensitivities to physics parameterizations, 107 they do not contain the full set of initial conditions and model uncertainties and are therefore 108 excluded from our study. 109

#### 110 b. Verification Data

A radar-based verification dataset was used as the verifying observations for the SSEF. The National Mosaic and Multi-Sensor QPE (NMQ; Vasiloff et al. 2007) derives precipitation estimates from the reflectivity data of the NEXRAD radar network. The estimates are made on a grid with 1 km horizontal spacing over the continental US (CONUS). The original grid has been bi-linearly interpolated to the same grid as the SSEF.

#### <sup>116</sup> c. Data Selection and Aggregation

The relative performance of any machine learning algorithm is conditioned on the dis-117 tribution of its training data. The sampling scheme for the SSEF is conditioned on the 118 constraints of 34 ensemble runs over a short, homogenous time period with 840,849 grid 119 points from each of the 30 time steps. The short training period and large number of grid 120 points preclude training a single model at each grid point, so a regional approach was used. 121 The SSEF domain was split into thirds (280,283 points per time step), and points were 122 selected with a uniform random sample from each subdomain in areas with quality radar 123 coverage. The gridded Radar Quality Index (RQI; Zhang et al. 2011) was evaluated at each 124 grid point to determine the trustworthiness of the verification data. Points with an RQI > 0125 were located within the useful range of a NEXRAD radar and included in the sampling. 126 For points with precipitation values less than 0.25 mm, 0.04% were sampled, and for points 127 with more precipitation, 0.4% were sampled. Grid 0 corresponds to the western third of 128 the CONUS, Grid 1 corresponds to the central third, and Grid 2 corresponds to the eastern 129 third. Fig. 1 shows that the north central section of the United States and the states south 130

of the Great Lakes were most heavily sampled. Grid 0 was excluded from post-processing
due to the low frequency of heavy precipitation for most of the region.

A comparison of the sampled rainfall distributions and the full rainfall distributions 133 for each subgrid are shown in Fig. 2 and Table 1. Undersampling of the 0 and trace 134 (below 0.25 mm) precipitation points was necessary because of the large number of no-135 precipitation events, which overwhelmed the signal from the actual precipitation events. 136 The random sampling of grid points helps reduce the chance of sampling multiple grid 137 points from the same storm without explicitly filtering subsets of the domain, which was 138 performed by Hamill et al. (2008). Filtering grid points surrounding each sampled point is 139 extremely computationally expensive because filtering requires multiple passes over the grid 140 while random sampling only requires 1 grid reordering. 141

Relevant model output variables, called predictors (Table 2), were also extracted from each ensemble member at each sampled point. These predictors captured additional information about the mesoscale and synoptic conditions in each model. The predictors from each ensemble member were then aggregated into 4 descriptive statistics for each predictor. The mean provided the most likely forecast value, and the standard deviation estimated the spread, and the minimum and maximum showed the extent of the forecasted values.

## $_{148}$ 3. Methods

#### <sup>149</sup> a. Machine Learning Methods

#### 150 1) LOGISTIC REGRESSION

One of the goals of this study is to compare the skill of more advanced machine learning methods with more traditional statistical methods. Logistic regression was used as the baseline statistical method. Logistic regressions are linear regression models in which a logit transformation is applied to the data so that the predicted values will range between 0 and 1. Two formulations of logistic regression were used: simple and multiple. The first formulation uses the ensemble mean precipitation forecast as expressed in Eqn. 1:

157 
$$p(R \ge t|x_1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1)}}$$
(1)

in which R is the precipitation amount, t is the threshold,  $x_1$  is the ensemble mean pre-158 cipitation forecast,  $\beta_0$  is the intercept term, and  $\beta_1$  is the weight given to the ensemble 159 mean precipitation forecast. Observed precipitation with an amount greater than or equal 160 to the chosen threshold is assigned a probability of 1, and precipitation amounts below the 161 threshold are assigned a probability of 0. The  $\beta$  coefficients are estimated by iteratively 162 maximizing the log-likelihood of the transformed training data. This formulation will adjust 163 the probability of the ensemble forecast based on systematic biases in the mean but will not 164 change the areal coverage of the precipitation forecasts. The second formulation incorpo-165 rates multiple variables with stepwise variable selection similar to standard Model Output 166

<sup>167</sup> Statistics (Glahn and Lowry 1972) approaches (Eqn. 2):

168 
$$p(R \ge t | x_1, x_2, \dots x_n) = \frac{1}{1 + \exp\left[-\left(\beta_0 + \sum_{n=1}^N \beta_n x_n\right)\right]}$$
(2)

in which  $x_1$  through  $x_n$  are the predictors chosen from Table 2 and  $\beta_1$  through  $\beta_n$  are the 169 weights for those predictors. The terms are determined through a forward selection process 170 that finds the set of up to 15 terms that minimize the Akaike Information Criterion (Akaike 171 1974), which rewards goodness of fit but penalizes for large numbers of terms. This approach 172 does produce the best fit regression model given the available parameters, but the searching 173 process can take extensive time given the large dimensionality of the training set. The qlm174 (generalized linear model) and step (stepwise variable selection) functions in the R statistical 175 package are used to generate the logistic regressions. A more detailed description of logistic 176 regression and the fitting process can be found in James et al. (2013). 177

#### 178 2) RANDOM FOREST

An non-parametric, nonlinear alternative to linear regression is the classification and regression decision tree (Breiman 1984). Decision trees recursively partition a multidimensional dataset into successively smaller subdomains by selecting variables and decision thresholds that maximize a dissimilarity metric. At each node in the tree, every predictor is evaluated with the dissimilarity metric, and the predictor and threshold with the highest metric value are selected as the splitting criteria for that node. After enough partitions, each subdomain is similar enough that the prediction can be approximated with a single value. The primary advantages of decision trees are that they can be human-readable and perform variable selection as part of the model growing process. The disadvantages lie in the brittleness of the trees. Trees can undergo significant structural shifts due to small variations in the training data, which results in large error variance.

Random forests (Breiman 2001) consist of an ensemble of classification and regression 190 trees (Breiman 1984) with two key modifications. First, the training data cases are bootstrap 191 resampled with replacement for each tree in the ensemble. Second, a random subset of the 192 predictors are selected for evaluation at each node. The final prediction from the forest is 193 the mean of the predicted probabilities from each tree. Random forests can produce both 194 probabilistic and regression predictions through this method. The random forest method 195 contains a few advantages that often lead to performance increases over traditional regres-196 sion methods. The averaging of the results from multiple trees produces a smoother range 197 of values than individual decision trees while also reducing the sensitivity of the model pre-198 dictions to minor differences in the training set (Strobl et al. 2008). The random selection of 199 predictors within the tree-building process allows for less optimal predictors to be included 200 in the model and increases the likelihood of the discovery of interaction effects among pre-201 dictors that would be missed by the stepwise selection method used in logistic regression 202 (Strobl et al. 2008). Random forests have been shown to improve predictive performance 203 on multiple problem domains in meteorology, including storm classification (Gagne II et al. 204 2009), aviation turbulence (Williams 2013), and wind energy forecasting (Kusiak and Verma 205 2011). For this project, we used the R random Forest library, which implements the original 206 approach (Breiman 2001). For the parameter settings, we chose to use 100 trees, a minimum 207 node size of 20, and the default values for all other parameters. A more detailed description 208

<sup>209</sup> of random forest can be found in James et al. (2013).

In addition to gains in performance, the random forest can also be used to rank the 210 importance of each input variable (Breiman 2001). Variable importance is computed by 211 first calculating the accuracy of each tree in the forest on classifying the cases that were not 212 selected for training, known as the out-of-bag cases. Within the out-of-bag cases, the values 213 of each variable are randomly rearranged, or permuted, and those cases are then re-evaluated 214 by each tree. The mean variable importance score is then the difference in prediction accuracy 215 on the out-of-bag cases averaged over all trees. Variable importance scores can vary randomly 216 among forests trained on the same dataset, so the variable importance scores from each of 217 the 34 forests trained for cross-validation were averaged together for a more robust ranking. 218

#### 219 b. Evaluation Methods

223

Two scores were used to assess the probabilistic forecasts. The Brier Skill Score (BSS) (BSS; Brier 1950) is one method used to evaluate probabilistic forecasts. The Brier Skill Score can be decomposed into three terms (Murphy 1973), as shown in Eq. 3:

$$BSS = \frac{\frac{1}{N}\sum_{k=1}^{K}n_k(\overline{o}_k - \overline{o})^2 - \frac{1}{N}\sum_{k=1}^{K}n_k(p_k - \overline{o}_k)^2}{\overline{o}(1 - \overline{o})}$$
(3)

<sup>224</sup> N is the number of forecasts, K is the number of probability bins,  $n_k$  is the number of <sup>225</sup> forecasts in each probability bin,  $\overline{o}_k$  is the observed relative frequency for each bin,  $\overline{o}$  is the <sup>226</sup> climatological frequency, and  $p_k$  is the forecast probability for a particular bin k. The first <sup>227</sup> term in the numerator describes the resolution of the forecast probability, which should be

maximized and increases as the observed relative frequency differs more from climatology. 228 The second term in the numerator describes the reliability of the forecast probability, which 229 should be minimized and decreases with smaller differences between the forecast probability 230 and observed relative frequency. The denominator term is the uncertainty, which is based 231 on the climatological probability of precipitation and cannot be reduced through calibration. 232 BSS increases with skill from 0. The components of the BSS can be displayed graphically 233 with an attributes diagram (Wilks 2011), in which the observed relative frequency of binned 234 probability forecasts are plotted against lines showing perfect reliability, no skill where the 235 reliability and resolution are equal, and no resolution where the observed relative frequency 236 and probability equal climatology. 237

The Area Under the Relative Operating Characteristic (ROC) curve, or AUC (Mason 238 1982) evaluates how well a probabilistic forecast correctly identifies heavy precipitation 239 events compared to identifying the events based on random chance. To calculate AUC, 240 first, the decision probability threshold is varied from 0 to 1 in equal steps. At each step, a 241 contingency table is constructed by splitting the probabilities into two categories with the 242 decision threshold (Table 3). Using Table 3, the following scores can be computed (Table 4). 243 Probability of detection (POD) is the ratio of hits to the total number of observed events. 244 The False Alarm Ratio (FAR) accounts for the number of false alarms compared to total yes 245 forecasts. The Probability of False Detection (POFD) is the ratio of false alarms to total no 246 forecasts. The Equitable Threat Score (ETS) is the skill score officially used to assess the 247 skill of 24 hour precipitation forecasts, and is sensitive to the climatological base rate of the 248 validation data. The Peirce Skill Score (PSS) is another skill score that is insensitive to the 249 climatological base rate but is sensitive to hedging by over forecasting rare events (Doswell 250

et al. 1990). The bias ratio compares determines if false alarms or misses are more prevalent. From the contingency table, the POD and POFD are calculated and plotted against each other, forming a ROC curve. The AUC is the area between the right and bottom sides of the plot and the curve itself. AUC with values above 0.5 has positive skill. AUC only determines how well the forecast discriminates between two categories, so it does not take the reliability of the forecast into account.

The ROC curve also can be used to determine an "optimal" decision threshold to convert 257 a probabilistic forecast to a deterministic forecast. As the decision probability increases, the 258 POD decreases from 1 to 0 while 1-POFD increases from 0 to 1. At some decision threshold, 259 POD and 1-POFD should be equal. At this optimal threshold (OT), the user has an equal 260 chance of correctly detecting rain events and no-rain events. This point occurs where the 261 ROC curve is farthest from the positive diagonal "no-skill" line. In addition, the vertical 262 distance between the ROC curve and the positive diagonal is the Peirce Skill Score (PSS; 263 Peirce 1884; Hansen and Kuipers 1965) since PSS=POD-POFD (Manzato 2007). Since the 264 PSS can be calculated from the contingency table at each decision threshold, finding the 265 maximum PSS also finds the OT. The performance of the threshold choice can be validated 266 by treating the predictions as a binary classification problem and using the binary verification 267 statistics derived from the contingency table (Table 4). 268

#### 269 c. Experimental Procedure

Each model was trained and evaluated using a leave-one-day-out cross validation procedure. For the 34 daily model runs in the training set, each machine learning model was

trained on 33 days and tested on 1. Separate models were trained for each 6 hour forecast 272 period and each subgrid to better capture the trends in the weather and ensemble dispersive-273 ness. Models were primarily trained at two thresholds,  $0.25 \text{ mm h}^{-1}$  and  $6.35 \text{ mm h}^{-1}$ . The 274  $0.25 \text{ mm h}^{-1}$  threshold models were used to determine if rain was to occur at a point or not, 275 and the 6.35 mm  $h^{-1}$  models were trained and evaluated only on the conditions that either 276 rain occurred or the ensemble mean precipitation forecast predicted rain. Deterministic rain 277 forecasts were derived from the  $0.25 \text{ mm h}^{-1}$  models with the ROC curve optimal threshold 278 technique to evaluate any biases in the areal forecasts. Additional models were trained at 279 the 2.54 and 12.70 mm  $h^{-1}$  thresholds to evaluate the skill in predicting lighter and heavier 280 precipitation, but only the 6.35 mm  $h^{-1}$  models were physically evaluated with variable im-281 portance and the case study. The 0.25, 2.54, 6.35, and 12.70 mm  $h^{-1}$  were chosen because 282 they correspond to the 0.01, 0.1, 0.25, and 0.5 in  $h^{-1}$  thresholds, respectively, that are used 283 for determining trace and heavy precipitation amounts. The probabilities shown in the case 284 study are the joint probabilities of precipitation greater than or equal to 6.35 and 0.25 mm 285  $h^{-1}$ . They were calculated by multiplying the conditional probability of precipitation greater 286 than or equal to  $6.35 \text{ mm h}^{-1}$  with the probability of precipitation greater than or equal to 287  $0.25 \text{ mm h}^{-1}$ . 288

## 289 4. Results

#### <sup>290</sup> a. Deterministic Rain Model Evaluation

Evaluation of the deterministic forecasts of precipitation location shows that the multiple 291 logistic regression and random forest do add skill compared to the uncalibrated ensemble 292 probability and the simple logistic regression. Fig. 3 shows only slight variations in the opti-293 mal threshold with forecast hour and that the thresholds for each machine learning algorithm 294 are similar. The low threshold for the raw ensemble indicates that the best precipitation cov-295 erage is found when any ensemble member forecasts rain. PSS and ETS show similar hourly 296 trends, but PSS is higher than ETS since ETS tends to maximize at a higher probability 297 threshold. The multi-predictor machine learning algorithms (random forest and multiple 298 logistic regression) provide the most improvement for the first 15 hours of the forecast with 299 only a slight improvement for the afternoon and evening. The improvement comes from 300 an increased POD without also increasing the FAR. Similar results are seen in grid 2 in 301 terms of thresholds and score values, and the multi-predictor algorithms are able to provide 302 consistent improvement over all forecast hour (Fig. 4). Similar results to these are found in 303 the SSEF verification from Kong et al. (2011) although the ETS is lower in that paper due 304 to that study verifying against all grid points, including points with no radar coverage. 305

#### 307

#### 1) EVALUATION OVER ALL FORECAST HOURS

The attributes diagrams for the raw ensemble and machine learning algorithms show 308 how the forecasts were altered to improve their calibration and skill. Fig. 5 contains the 309 attributes diagrams for the ensemble and algorithms applied to grids 1 and 2 at the 6.35 mm 310  $h^{-1}$  threshold for all forecast hours. The raw ensemble in both domains tends to be over-311 confident with probabilistic forecasts above climatology and under confident with forecasts 312 below climatology, resulting in a negative BSS. In spatial terms, the ensemble probabilities 313 are too high for areas where it forecasts heavy rain, and it is missing some areas where rain 314 actually occurred. The simple logistic regression improves the ensemble forecast by rescaling 315 the probabilities based on the ensemble mean rain amount. In grid 1 and grid 2 this gener-316 ally results in all of the probabilities being scaled closer to the climatological probability to 317 address the general overconfidence of the ensemble. This scaling brings the observed relative 318 frequencies above climatology into the positive skill zone, but they are still overconfident. 319

The multiple logistic regression and random forest incorporate additional ensemble vari-320 ables to add more information to the probability estimates. This does result in a large 321 improvement in BSS over the simple logistic regression. The multiple logistic regression 322 performs the best in grid 1 because it produces nearly perfect reliability between 0 and 50%323 while maintaining positive skill from 60 to 90%. The random forest also places its forecasts 324 close to the perfect reliability line but slightly further away than the stepwise logistic re-325 gression, and the forecast probabilities do not exceed 70%, so the random forest has a lower 326 maximum resolution than the multiple logistic regression. 327

In grid 2, the BSS for all three models decreases. The observed relative frequency for the random forest lies just above the no skill line and then increases above 50%, which may be a due to chance from the low sample size in the higher percentage range. The multiple logistic regression performs better with the lower probability forecasts but has negative skill with the higher probabilities. The weaker performance may be due to having fewer heavy rain events in grid 2 during the study time period.

The ROC curves show that the multiple-predictor machine learning algorithms enhance 334 the discrimination abilities of the ensemble. The raw ensemble in both subgrids has slightly 335 positive skill in terms of AUC (Fig. 6). At the optimal PSS threshold, the raw ensemble 336 detects only 67% of heavy rain events (POD) in grid 1, and 86% of its positive forecasts are 337 false alarms (FAR). In grid 2, the detection ability is worse with a 66% POD and a 91%338 FAR. The Bias scores greater than 1 indicate a larger proportion of false alarms than misses. 339 Since the simple logistic regression rescales the predictions over a smaller probability range, 340 it has a slightly lower AUC than the raw ensemble. Its optimal threshold is higher than the 341 raw ensemble, so it has a correspondingly lower POD, POFD, and FAR. It also has a smaller 342 bias score at the optimal threshold. The multiple logistic regression and random forest have 343 a much larger AUC and POD and similar FAR compared to the raw ensemble. The multiple 344 logistic regression and random forest have very similar scores with only slight differences in 345 the POFD, FAR, and bias. The same relationships hold in grid 2 except that there is a 346 slightly larger difference in the scores of the multiple logistic regression and random forest. 347 The overall scores are also slightly lower for grid 2, which is likely due to the lower frequency 348 of convective precipitation events in the training data. 349

#### 350 2) EVALUATION BY FORECAST HOUR

Comparisons of BSS and AUC by hour and by sub-grid show additional trends in the 351 probability of precipitation forecasts. The raw ensemble consistently has the worst BSS (Fig. 352 7). For all models, the best performance occurs at forecast hour 1 then decreases sharply at 353 hour 2 before stabilizing. This initial decrease is likely due to the radar data assimilation 354 placing the storms in the same place initially and then having the individual storms diverge 355 from the predicted storm motions. There is a slight increase in performance between hours 356 6 and 12 for grid 2. This increase may be due to the diurnal cycle of convection resulting in 357 larger storm clusters during this time period, or it may be due to the model fully spinning 358 up. There is another major decrease in performance in grid 1 between hours 12 and 24. 359 This time period is when the greatest uncertainty exists due to convective initiation and 360 the tendency for initial convection to be isolated. The multiple logistic regression and the 361 random forest do not have any statistically significant ( $\alpha < 0.05$ ) differences in their AUC 362 and BSS. The biggest departure from the ensemble mean in terms of AUC occurred in the 363 18 to 24 hour range for both grids, which corresponds with peak convective activity. The 364 simple logistic regression did still improve on the ensemble forecast in terms of reliability 365 but not to the same extent as the multiple-predictor models, and it made the AUC worse. 366 The temporal trends in the BSS match those found in the 2009 SSEF by Johnson and Wang 367 (2012). The multiple-predictor methods appeared to produce similar increases in BSS to the 368 single predictor neighborhood and object-based methods. 369

Machine learning models were trained at 3 precipitation thresholds (2.54, 6.35, and 12.70) 371  $mm h^{-1}$ ) in order to evaluate how skill varies with precipitation intensity. Fig. 8 shows the 372 random forest BSS by hour for the three precipitation thresholds. All three follow the same 373 diurnal patterns, but the 2.54 mm  $h^{-1}$  forecasts have consistently much higher skill than 374 the other thresholds. While the 2.54 and 6.35 mm  $h^{-1}$  show skill for all forecast hours, the 375  $12.70 \text{ mm h}^{-1}$  forecasts do not show any skill from 14 to 26 hours. The decreasing skill with 376 threshold size is likely due to the smaller heavy precipitation areas and spatial and timing 377 errors in placement of convection in the ensembles. Similar results were found for the other 378 machine learning methods (not shown). 379

#### 380 d. Variable Importance

Variable importance scores were calculated and averaged over each random forest and 381 sub-grid to determine if the random forests were choosing relevant variables and how the 382 choice of variables was affected by region. Variable importance is indicative of how random-383 izing the values of each variable affects the random forest performance. It accounts for how 384 often a variable is used in the model, the depth of the variable in the tree, and the num-385 ber of cases that transit through the branch containing that variable, but the importance 386 score cannot be decomposed into those factors. The top 8 variable importance scores for the 387 random forests trained on each 6-hour period in Grid 1 are shown in Table 5. In the first 388 6 hours, all of the top 5 variables are aggregations of the hour precipitation or precipitable 389 water with the rest being hour max upward wind. In this time period, the ensemble members 390

are very similar, so there is great overlap among the precipitation regions. By hours 7-12, the max upward wind becomes more dominant in the rankings, although the precipitation max and standard deviation are still found among the top 5 variables. Vertical velocities become the most common feature in hours 13 through 18 with only precipitable water and specific humidity contributing moisture information. For hours 19 through 30, the standard deviation of Surface-Based CAPE appears, which is likely associated with the presence of nearby boundaries.

The grid 2 variable importances highlight the greater importance of moisture and lower importance of vertical velocities in the Eastern United States. The predicted precipitation is again only important in hours 1 through 12, but the precipitable water and specific humidity show high importance through the entire forecast period. Upward and Downward winds are in the rankings for each time period, but they tend to be on toward the bottom of the top 8. Surface-Based CAPE is also important in the latter hours of the forecasts, but the mean and max are selected instead of the standard deviation.

#### 405 e. Case Study: 13 May 2010

The case of 13 May 2010 illustrates the spatial characteristics, strengths, and weaknesses of the precipitation forecasts from the SSEF and the machine learning methods. Since the SSEF is run through forecast hour 30, the last six forecast hours of one run overlap with the first six hours of the next run. This overlap allows for the comparison of two runs on the same observations and illustrates the effects of lead time on the forecast probabilities. Fig. 9 shows the distribution of the observed 1-hour precipitation at 02 UTC on May 13. The bulk of the

precipitation originates from a broken line of discrete supercell thunderstorms positioned 412 along a stationary front stretching from northern Oklahoma into northern Missouri with 413 additional storms in Iowa and ahead of the dry line in Oklahoma and Texas. Additional 414 precipitation is falling in Virginia from a mesoscale convective system. A comparison of 415 the 2-hour and 26-hour dew point and temperature SSEF forecasts (Fig. 10) shows major 416 differences in the placement and strength of the fronts and dry line. The cold front is further 417 north and more diffuse in the 26 hour forecast while the dry line is further east, moving from 418 the central Texas panhandle to western Oklahoma. The differences in the placement of the 419 surface boundaries also affect the placement of the precipitation areas in both forecasts. 420

Comparisons of the raw ensemble probabilities and the simple logistic regression (Fig. 11) 421 show the effects of calibrating the probabilities only on precipitation. All of the precipitation 422 probabilities shown are generated by multiplying the conditional PQPF of greater than 423  $6.35 \text{ mm h}^{-1}$  with the probability of precipitation, which removed spurious low conditional 424 probabilities of precipitation. In the 2 hour forecast, there is little dispersion among the 425 ensemble members, so there is very high confidence in the raw ensemble. In this case though, 426 the high confidence areas are displaced west from the actual heavy precipitation regions. 427 resulting in some missed precipitation areas. The simple logistic regression corrects for 428 this overconfidence by lowering the SSEF probabilities. The area with a greater than 10%429 probability is noticeably smaller and approximates the area occupied by the observed rain 430 Since it rescales the probabilities, it does not translate the predicted rain areas. areas. 431 resulting in more misses. At 26 hours, the raw ensemble probabilities are lower and spread 432 over a wider area, indicating greater ensemble dispersion. The probabilities in Oklahoma 433 and Kansas are shifted further east due to the positioning of the surface boundaries. In this 434

case, it results in the observed precipitation being captured better than they were in the 2 435 hour forecast with fewer misses but more false alarms. The simple logistic regression again 436 reduces the probabilities and the area covered by them so that it does not have significantly 437 more false alarms than at 2 hours, but it misses some precipitation that was captured by the 438 raw ensemble. Of the machine learning methods, the simple logistic is best for accurately 439 depicting areal coverage, but it overcorrected in terms of downscaling the raw probabilities. 440 The results from the multiple-predictor methods show the advantages of incorporating 441 additional predictors. In Fig. 12, the multiple logistic regression is compared with the 442 random forest at 2 and 26 hours. In the 2 hour forecasts, both the multiple logistic regression 443 and random forest expand their probabilities over a wider area and shift them slightly east 444 compared to the raw ensemble, enabling them to capture the heavy precipitation fully. 445 The random forest also captures the precipitation areas north of the cold front better than 446 the logistic regression, which may weight SBCAPE too highly and not handle cold sector 447 precipitation as well. At 26 hours, the areas increase and the probabilities are generally 448 lower. Both algorithms capture nearly all precipitation areas, but they also produce more 449 false alarms. Precipitation did extend through central Texas into Mexico ahead of the dry 450 line, but none of that precipitation was heavy. The multiple logistic regression highlighted 451 that precipitation as well as ahead of the dry line. Unlike the 2 hour forecast, it did cover the 452 precipitation north of the front in Iowa. The random forest covered all of the precipitation 453 areas in the Plains and seemed to handle the east coast precipitation coverage better. 454

The extended probability areas of the random forest and multiple logistic regression can be explained by examining some of the variables deemed important by the random forest. In Fig. 13, the Hour Max Upward Wind, SBCAPE, Precipitable Water, and 850 mb Specific

Humidity are shown. The Hour Max Upward Wind was the most important variable for 458 the 25-30 hour random forest (Table 5), and its spatial distribution closely matches that 459 of the random forest, including the positive areas in southwest Texas. It also supports the 460 precipitation north of the cold front in Iowa. The SBCAPE was highest in central Texas 461 in the area where both machine learning models placed probabilities but no heavy rain 462 occurred. The machine learning algorithms were able to account for situations in which 463 high CAPE and no precipitation did not occur together by not assigning probabilities in the 464 Gulf of Mexico. Precipitable water and 850 mb specific humidity were also high over a wide 465 area and are generally correlated well with precipitation but occur over a much larger area 466 than the actual precipitation. By weighting these variables together, the multi-predictor 467 algorithms were able to grow the forecasted area to include points at which the mix of 468 ingredients was favorable for precipitation. 469

## 470 5. Discussion

The results of the machine-learning post-processing of storm scale ensemble precipita-471 tion forecasts displayed not only improvements to the forecasts but also limitations of the 472 ensemble, the algorithms, and grid-point-based framework. First, the post-processing model 473 performance is constrained by the information available from the ensemble. If most of the 474 ensemble members are predicting precipitation in the wrong place or not at all, and the 475 environmental conditions are also displaced, then the machine learning algorithm will not 476 be able to provide much additional skill. Second, the machine learning model will only make 477 predictions based on the range of cases it has previously seen. For higher rain thresholds, the 478

algorithms will need more independent samples in order to make skilled predictions. Third, 479 the grid-point framework does not fully account for the spatial information and error in the 480 ensemble. The spatial error could be incorporated further by smoothing the verification grid 481 with a neighborhood filter (Ebert 2009) or warping the ensemble forecast to more closely 482 match the observed precipitation (Gilleland et al. 2010) and then training. Machine learn-483 ing algorithms could also be used to improve the calibration from object-based frameworks 484 (Johnson and Wang 2012) and could incorporate information from the additional predictors 485 within the bounds of the objects. 486

## 487 6. Conclusions

Multiple machine learning algorithms were applied to the 2010 CAPS Storm Scale Ensem-488 ble Forecast (SSEF) system in order to improve the calibration and skill of its probabilistic 489 heavy precipitation forecasts. Two types of machine learning methods were compared over 490 a period from May 3 through June 18 with both verification statistics and a case study. 491 Verification statistics showed that all of the machine learning methods improved the cali-492 bration of the SSEF precipitation forecasts but only the multiple-predictor methods were 493 able to calibrate the models better and discriminate more skillfully between light and heavy 494 precipitation cases. Hourly performance varied with diurnal storm cycles and the increasing 495 dispersiveness of the ensemble members. Comparisons of the rankings of predictors indi-496 cated that the ensemble predicted precipitation was only important for the first 12 hours 497 of the model runs. After that period, the upward wind and atmospheric moisture variables 498 became better indicators of the placement of precipitation. The case study showed that the 499

multiple-predictor machine learning methods could shift the probability maxima to better 500 match the actual precipitation areas, but they would also produce more false alarm areas 501 in the process. For shorter-term forecasts, the false corrections were made without a sig-502 nificant increase in the false alarm area. Calibrating the probabilities with only ensemble 503 rainfall predictions results in predicted areas that are too small and still displaced from 504 the observed precipitation. The multiple-predictor machine learning algorithms did prove 505 especially beneficial in that situation. Ultimately, machine learning techniques can provide 506 an enhancement to precipitation forecasts by consistently maximizing the potential of the 507 available information. 508

#### 509 Acknowledgments.

Special thanks go to my Master's committee members Fanyou Kong and Michael Rich-510 man. Zac Flamig provided assistance with the NMQ data. CAPS SSEF forecasts were 511 supported by a grant (NWSPO-2010-201696) from the NOAA Collaborative Science, Tech-512 nology, and Applied Research (CSTAR) Program and the forecasts were produced at the 513 National Institute for Computational Science (http://www.nics.tennessee.edu/). Scientists 514 at CAPS, including Fanyou Kong, Kevin Thomas, Yunheng Wang, and Keith Brewster con-515 tributed to the design and production of the CAPS ensemble forecasts. This study was 516 funded by the NSF Graduate Research Fellowship under Grant 2011099434 and by NSF 517 grant AGS-0802888. The editor, Michael Baldwin, and the two anonymous reviewers pro-518 vided very helpful feedback that strengthened the quality of the paper. 519

## REFERENCES

- Akaike, H., 1974: A new look at the statistical model identification. *IEEE Transactions on* 522
- Automatic Control, 19, 716–723. 523
- Applequist, S., G. E. Gahrs, R. L. Pfeffer, and X. Niu, 2002: Comparison of methodologies 524 for probabilistic quantitative precipitation forecasting. Wea. Forecasting, 17, 783–799.
- Ashley, S. T. and W. S. Ashley, 2008: Flood fatalities in the united states. J. Appl. Meteorol., 526 **47**, 805–818. 527
- Breiman, L., 1984: Classification and regression trees. Wadsworth International Group, 358 528 pp. 529
- Breiman, L., 2001: Random forests. Mach. Learn., 45, 5 32. 530
- Bremnes, J. B., 2004: Probabilistic forecasts of precipitation in terms of quantiles using 531 NWP model output. Mon. Wea. Rev., 132, 338–347. 532
- Brier, G. W., 1950: Verification of forecasts expressed in terms of probability. Mon. Wea. 533 *Rev.*, **78**, 1–3. 534
- Clark, A. J. and Coauthors, 2011: Probabilistic precipitation forecast skill as a function of 535 ensemble size and spatial scale in a convection-allowing ensemble. Mon. Wea. Rev., 139, 536 1410 - 1418.537

521

525

- <sup>538</sup> Clark, A. J., W. A. Gallus, M. Xue, and F. Kong, 2009: A comparison of precipitation
  <sup>539</sup> forecast skill between small convection-allowing and large convection-parameterizing en<sup>540</sup> sembles. *Wea. Forecasting*, 24, 1121–1140.
- <sup>541</sup> Clark, A. J., S. J. Weiss, J. S. Kain, and Coauthors, 2012: An overview of the 2010 hazardous
  <sup>542</sup> weather testbed experimental forecast program spring experiment. *Bull. Amer. Meteor.*<sup>543</sup> Soc., 93, 55–74.
- <sup>544</sup> Doswell, C. A., H. E. Brooks, and R. A. Maddox, 1996: Flash flood forecasting: An <sup>545</sup> ingredients-based methodology. *Weather and Forecasting*, **11**, 560–581.
- <sup>546</sup> Doswell, C. A., R. Davies-Jones, and D. L. Keller, 1990: On summary measures of skill in <sup>547</sup> rare event forecasting based on contingency tables. *Wea. Forecasting*, **5**, 576–585.
- <sup>548</sup> Du, J., J. McQueen, G. DiMego, Z. Toth, D. Jovic, B. Zhou, and H. Chuang, 2006:
  <sup>549</sup> New dimension of NCEP short-range ensemble forecasting (SREF) system: Inclusion
  <sup>550</sup> of wrf members. *Preprint, WMO Expert Team Meeting on Ensemble Prediction System*,
  <sup>551</sup> Exeter, United Kingdom, WMO, 5, URL http://www.emc.ncep.noaa.gov/mmb/SREF/
  <sup>552</sup> reference.html.
- Ebert, E., 2001: Ability of a poor man's ensemble to predict the probability and distribution
  of precipitation. *Mon. Wea. Rev.*, **129**, 2461–2480.
- Ebert, E. E., 2009: Neighborhood verification: A strategy for rewarding close forecasts.
  Weather and Forecasting, 24 (6), 1498–1510.
- <sup>557</sup> Eckel, F. A. and M. K. Walters, 1998: Calibrated probabilistic quantitative precipitation
  <sup>558</sup> forecasts based on the MRF ensemble. *Wea. Forecasting*, **13**, 1132–1147.

27

- Gagne II, D. J., A. McGovern, and J. Brotzge, 2009: Classification of convective areas using
  decision trees. *Journal of Atmospheric and Oceanic Technology*, 26 (7), 1341–1353.
- Gagne II, D. J., A. McGovern, and M. Xue, 2012: Machine learning enhancement of storm
   scale ensemble precipitation forecasts. *Proceedings of the Conference on Intelligent Data Understanding (CIDU)*, Boulder, CO, IEEE-CIS, 39–46.
- Gilleland, E., et al., 2010: Spatial forecast verification: image warping. Tech. Rep.
   NCAR/TN-482+STR, NCAR.
- Glahn, H. R. and D. A. Lowry, 1972: The use of model output statistics (MOS) in objective
  weather forecasts. J. Appl. Meteor., 11, 1203–1211.
- Hall, T., H. E. Brooks, and C. A. Doswell, 1999: Precipitation forecasting using a neural
   network. Wea. Forecasting, 14, 338–345.
- Hamill, T. M. and S. J. Colucci, 1997: Verification of Eta-RSM short-range ensemble forecasts. Mon. Wea. Rev., 125, 1312–1327.
- Hamill, T. M. and S. J. Colucci, 1998: Evaluation of Eta-RSM ensemble probabilistic precipitation forecasts. *Mon. Wea. Rev.*, **126**, 711–724.
- Hamill, T. M., R. Hagedorn, and J. S. Whitaker, 2008: Probabilistic forecast calibration
  using ECMWF and GFS ensemble reforecasts. Part II: Precipitation. *Mon. Wea. Rev.*,
  136, 2620–2632.
- 577 Hamill, T. M., J. S. Whitaker, and X. Wei, 2004: Ensemble reforecasting: Improving

- medium-range forecast skill using retrospective forecasts. Mon. Wea. Rev., 132, 1434–
  1447.
- Hansen, A. W. and W. J. A. Kuipers, 1965: On the relationship between the frequency of
  rain and various meteorological parameters. *Meded. Verhand.*, 81, 2–15.
- James, G., D. Witten, T. Hastie, and R. Tibshirani, 2013: An Introduction to Statistical Learning with Applications in R. Springer, 430 pp.
- Johnson, A. and X. Wang, 2012: Verification and calibration of neighborhood and objectbased probabilistic precipitation forecasts from a multimodel convection-allowing ensemble. *Mon. Wea. Rev.*, **140**, 3054–3077.
- Koizumi, K., 1999: An objective method to modify numerical model forecasts with newly
  given weather data using an artificial neural network. *Wea. Forecasting*, 14, 109–118.
- Kong, F., et al., 2011: Evaluation of CAPS multi-model storm-scale ensemble forecast for
- the NOAA HWT 2010 spring experiment. 24th Conf. Wea. Forecasting/20th Conf. Num.
- <sup>591</sup> Wea. Pred., Seattle, WA, Amer. Meteor. Soc., Paper 457.
- 592 Krishnamurti, T. N., C. M. Kishtawal, T. E. LaRow, D. R. Bachiochi, Z. Zhang, C. E.
- <sup>593</sup> Williford, S. Gadgil, and S. Surendran, 1999: Improved weather and seasonal climate <sup>594</sup> forecasts from multimodel superensemble. *Science*, **285**, 1548–1550.
- Kusiak, A. and A. Verma, 2011: Prediction of status patterns of wind turbines: A datamining approach. J. Sol. Energy Eng., 133, doi:10.1115/1.4003188.

- <sup>597</sup> Manzato, A., 2007: A note on the maximum Peirce skill score. *Wea. Forecasting*, **22**, 1148– <sup>598</sup> 1154.
- Marsh, P. T., J. S. Kain, V. Lakshmanan, A. J. Clark, N. Hitchens, and J. Hardy, 2012:
  A method for calibrating deterministic forecasts of rare events. *Wea. Forecasting*, 27, 531–538.
- Mason, I., 1982: A model for assessment of weather forecasts. Aust. Meteor. Mag., 30,
  291–303.
- Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF ensemble
  prediction system: Methodology and validation. *Quart. J. Roy. Meteor. Soc.*, **122**, 73–
  119.
- Murphy, A. H., 1973: A new vector partition of the probability score. J. Appl. Meteor., 12,
  595–600.
- Murphy, A. H., 1977: The value of climatological, categorical, and probabilistic forecasts in the cost-loss ratio situation. *Mon. Wea. Rev.*, **105**, 803–816.
- Peirce, C. S., 1884: The numerical measure of the success of predictions. *Science*, 4, 453–454.
- 612 Strobl, C., A.-L. Boulesteix, T. Kneib, T. Augustin, and A. Zeileis, 2008: Conditional
- variable importance for random forests. BMC Bioinformatics, 9 (1), 307, doi:10.1186/
- <sup>614</sup> 1471-2105-9-307, URL http://www.biomedcentral.com/1471-2105/9/307.
- <sup>615</sup> Toth, Z. and E. Kalnay, 1993: Ensemble forecasting at NMC: The generation of perturba-<sup>616</sup> tions. *Bull. Amer. Meteor. Soc.*, **74**, 2317–2330.

- Tracton, M. S. and E. Kalnay, 1993: Operational ensemble prediction at the National Me teorological Center: Practical aspects. *Wea. Forecasting*, 8, 379–398.
- <sup>619</sup> Vasiloff, S., et al., 2007: Improving QPE and very short term QPF: An initiative for a <sup>620</sup> community-wide integrated approach. *Bull. Amer. Meteor. Soc.*, **88**, 1899–1911.
- Wilks, D. S., 2011: Statistical Methods in the Atmospheric Sciences. 3d ed., Academic Press,
  676 pp.
- Williams, J. K., 2013: Using random forests to diagnose aviation turbulence. Mach. Learn.,
  doi:10.1007/s10994-013-5346-7.
- <sup>625</sup> Xue, M., K. K. Droegemeier, and V. Wong, 2000: The Advanced Regional Prediction System
- (ARPS) A multiscale nonhydrostatic atmospheric simulation and prediction model. Part
  I: Model dynamics and verification. *Meteor. Atmos. Phys.*, **75**, 161–193.
- Xue, M., D. Wang, J. Gao, K. Brewster, and K. K. Droegemeier, 2003: The Advanced
  Regional Prediction System (ARPS), storm-scale numerical weather prediction and data
  assimilation. *Meteor. Atmos. Phys.*, 82, 139–170.
- <sup>631</sup> Xue, M., et al., 2001: The Advanced Regional Prediction System (ARPS) A multiscale
  <sup>632</sup> nonhydrostatic atmospheric simulation and prediction model. Part II: Model physics and
  <sup>633</sup> applications. *Meteor. Atmos. Phys.*, **76**, 134–165.
- <sup>634</sup> Xue, M., et al., 2011: CAPS realtime storm scale ensemble and high resolution forecasts for
- the NOAA hazardous weather testbed 2010 spring experiment. 24th Conf. Wea. Forecast-
- <sup>636</sup> ing/20th Conf. Num. Wea. Pred., Seattle, WA, Amer. Meteor. Soc., PAPER 9A.2.

- Yuan, H., X. Gao, S. L. Mullen, S. Sorooshian, J. Du, and H. H. Juang, 2007: Calibration
   of probabilistic quantitative precipitation forecasts with an artificial neural network. Wea.
   *Forecasting*, 22, 1287–1303.
- <sup>640</sup> Zhang, J., Y. Qi, K. Howard, C. Langston, and B. Kaney, 2011: Radar quality index (RQI)—
- a combined measure of beam blockage and vpr effects in a national network. Proc. Eighth
- Int. Symp. on Weather Radar and Hydrology, Exeter, United Kingdom, Royal Meteoro logical Society.

## 644 List of Tables

|   | thresholds.   | 34  |
|---|---|---|
| 2 | The names and descriptions of model predictors sampled from the SSEF runs.    |   |
|   | CAPE is Convective Available Potential Energy, and CIN is Convective Inhi-    |   |
|   | bition.   | 35  |
| 3 | An example binary contingency table for whether or not rain is forecast.      | 36  |
| 4 | The scores calculated from the binary contingency table for use with the      |   |
|   | optimal threshold predictions.  | 37  |
| 5 | Top 8 variable importance scores for the random forests trained over 6-hourly |   |
|   | periods in Grid 1.  | 38  |
| 6 | Top 8 variable importance scores for the random forests trained over 6-hourly |   |
|   | periods in Grid 2.  | 39  |
|   | 2<br>3<br>4<br>5<br>6   | <ul> <li>thresholds.</li> <li>2 The names and descriptions of model predictors sampled from the SSEF runs.<br/>CAPE is Convective Available Potential Energy, and CIN is Convective Inhibition.</li> <li>3 An example binary contingency table for whether or not rain is forecast.</li> <li>4 The scores calculated from the binary contingency table for use with the optimal threshold predictions.</li> <li>5 Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 1.</li> <li>6 Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 2.</li> </ul> |

TABLE 1. Frequencies of total and sampled grid points for the different precipitation thresholds.  $\mbox{Grid} \ | \ \ 0{-}0.25 \ \mbox{mm} \ \ 0.25{-}6 \ \mbox{mm} \ \ 6 \ \mbox{mm}$ 

| Grid           | 0-0.25  mm        | 0.25-6  mm | 6  mm     |
|----------------|-------------------|------------|-----------|
| Grid 1 Total   | 218,052,271       | 13,062,623 | 1,884,725 |
| Grid 1 Sampled | 86,704            | $51,\!802$ | $7,\!490$ |
| Grid 2 Total   | $194,\!343,\!572$ | 10,923,189 | 1,137,694 |
| Grid 2 Sampled | 77,210            | 4,3218     | $4,\!497$ |

TABLE 2. The names and descriptions of model predictors sampled from the SSEF runs. CAPE is Convective Available Potential Energy, and CIN is Convective Inhibition.

| Variable                     | Levels                          |
|------------------------------|---------------------------------|
| Hour Precipitation           | Surface                         |
| Surface-Based CAPE           | Surface                         |
| Surface-Based CIN            | Surface                         |
| Dewpoint                     | 2 m                             |
| Pressure                     | Mean Sea Level                  |
| Composite Radar Reflectivity | Column Maximum                  |
| Precipitable Water           | Column Sum                      |
| Height                       | 700  mb                         |
| U-Wind                       | 700 mb                          |
| V-Wind                       | 700  mb, 500  mb                |
| Specific Humidity            | 850  mb, 700  mb, 500  mb       |
| Temperature                  | 2  m, 850  mb, 700  mb, 500  mb |
| Hour Max Reflectivity        | Column and Time Maximum         |
| Hour Max Upward Wind         | Column and Time Maximum         |
| Hour Max Downward Wind       | Column and Time Maximum         |

|          |     | Observed  |                   |  |  |
|----------|-----|-----------|-------------------|--|--|
|          |     | Yes       | No                |  |  |
| Forecast | Yes | a (hit)   | b (false alarm)   |  |  |
| Torecast | No  | c  (miss) | d (true negative) |  |  |

TABLE 3. An example binary contingency table for whether or not rain is forecast.

| Score        | Formula   |
|--------------|---|
| POD          | $\frac{a}{a+c}$                                 |
| FAR          | $\frac{b}{a+b}$                                 |
| POFD         | $\frac{b}{b+d}$                                 |
| ETS          | $\frac{a - a_{random}}{a + b + c - a_{random}}$ |
| $a_{random}$ | $\frac{(a+c)(a+b)}{a+b+c+d}$                    |
| PSS          | $\frac{a}{a+c} - \frac{b}{b+d}$                 |
| Bias         | $\frac{a+b}{a+c}$                               |

TABLE 4. The scores calculated from the binary contingency table for use with the optimal threshold predictions.

| Variable                      | Mean   | Variable                      | Mean   |
|-------------------------------|--------|-------------------------------|--------|
| Grid 1 Hours 1-6              |        | Grid 1 Hours 7-12             |        |
| Hour Precipitation SD         | 0.0118 | Hour Max Upward Wind Mean     | 0.0129 |
| Precipitable Water Mean       | 0.0107 | Hour Max Upward Wind Max      | 0.0126 |
| Hour Precipitation Mean       | 0.0105 | Hour Max Upward Wind SD       | 0.0113 |
| Precipitable Water Max        | 0.0103 | Hour Precipitation SD         | 0.0105 |
| Hour Precipitation Max        | 0.0096 | Hour Precipitation Max        | 0.0088 |
| Hour Max Upward Wind Mean     | 0.0088 | Hour Max Downward Wind SD     | 0.0087 |
| Hour Max Upward Wind Max      | 0.0088 | Hour Max Downward Wind Min    | 0.0084 |
| Hour Max Upward Wind SD       | 0.0081 | Specific Humidity 700 mb Max  | 0.0076 |
| Grid 1 Hours 13-18            |        | Grid 1 Hours 19-24            |        |
| Hour Max Upward Wind Mean     | 0.0096 | Hour Max Downward Wind Min    | 0.0114 |
| Hour Max Downward Wind SD     | 0.0088 | Surface-Based CAPE SD         | 0.0112 |
| Hour Max Downward Wind Mean   | 0.0086 | Hour Max Downward Wind Mean   | 0.0110 |
| Hour Max Downward Wind Min    | 0.0084 | Hour Max Downward Wind SD     | 0.0105 |
| Hour Max Upward Wind SD       | 0.0071 | Specific Humidity 850 mb Mean | 0.0092 |
| Precipitable Water Max        | 0.0070 | Temperature 700 mb Min        | 0.0089 |
| Specific Humidity 700 mb Max  | 0.0069 | Hour Max Upward Wind Mean     | 0.0088 |
| Hour Max Upward Wind Max      | 0.0062 | Hour Max Upward Wind Max      | 0.0084 |
| Grid 1 Hours 25-30            |        |                               |        |
| Hour Max Upward Wind Mean     | 0.0123 |                               |        |
| Hour Max Downward Wind Mean   | 0.0119 |                               |        |
| Hour Max Upward Wind SD       | 0.0112 |                               |        |
| Hour Max Upward Wind Max      | 0.0108 |                               |        |
| Hour Max Downward Wind SD     | 0.0105 |                               |        |
| Hour Max Downward Wind Min    | 0.0104 |                               |        |
| Specific Humidity 850 mb Mean | 0.0078 |                               |        |
| Surface-Based CAPE SD         | 0.0077 |                               |        |

TABLE 5. Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 1.

| Variable                      | Mean   | Variable                      | Mean   |
|-------------------------------|--------|-------------------------------|--------|
| Grid 2 Hours 1-6              |        | Grid 2 Hours 7-12             |        |
| Hour Precipitation Mean       | 0.0116 | Precipitable Water Mean       | 0.0111 |
| Precipitable Water Mean       | 0.0114 | Precipitable Water Min        | 0.0085 |
| Precipitable Water Min        | 0.0104 | Precipitable Water Max        | 0.0079 |
| Precipitable Water Max        | 0.0097 | Hour Max Upward Wind Mean     | 0.0077 |
| Hour Max Upward Wind Mean     | 0.0094 | Hour Max Upward Wind SD       | 0.0066 |
| Hour Max Reflectivity Mean    | 0.0091 | Hour Precipitation Mean       | 0.0065 |
| Hour Max Upward Wind Max      | 0.0082 | Hour Precipitation SD         | 0.0062 |
| Hour Precipitation Max        | 0.0080 | Hour Max Upward Wind Max      | 0.0062 |
| Grid 2 Hours 13-18            |        | Grid 2 Hours 19-24            |        |
| Precipitable Water Mean       | 0.0077 | Specific Humidity 850 mb Mean | 0.0162 |
| Hour Max Upward Wind Mean     | 0.0075 | Specific Humidity 850 mb Max  | 0.0128 |
| Specific Humidity 850 mb Mean | 0.0065 | Hour Max Upward Wind Mean     | 0.0110 |
| Precipitable Water Min        | 0.0065 | Surface-Based CAPE Max        | 0.0106 |
| Precipitable Water Max        | 0.0056 | Surface-Based CAPE Mean       | 0.0103 |
| Temperature 700 mb Mean       | 0.0056 | Hour Max Upward Wind SD       | 0.0101 |
| Hour Max Upward Wind Max      | 0.0055 | Hour Max Downward Wind Mean   | 0.0098 |
| Temperature 500 mb Min        | 0.0054 | Hour Max Upward Wind Max      | 0.0090 |
| Grid 2 Hours 25-30            |        |                               |        |
| Specific Humidity 850 mb Mean | 0.0114 |                               |        |
| Hour Max Upward Wind Mean     | 0.0094 |                               |        |
| Precipitable Water Mean       | 0.0084 |                               |        |
| Specific Humidity 850 mb Max  | 0.0083 |                               |        |
| Hour Precipitation SD         | 0.0072 |                               |        |
| Hour Max Upward Wind Max      | 0.0071 |                               |        |
| Hour Max Downward Wind Mean   | 0.0070 |                               |        |
| Surface-Based CAPE Max        | 0.0069 |                               |        |

TABLE 6. Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 2.

# 657 List of Figures

| 658 | 1 | Map of the number of grid points sampled within a 400 $\rm km^2$ box spatially      |    |
|-----|---|---|----|
| 659 |   | averaged with a 300 km radius median filter. The domain sub grids are also          |    |
| 660 |   | shown and labeled.  | 42 |
| 661 | 2 | Histogram comparing the relative frequencies of the full precipitation distri-      |    |
| 662 |   | bution for the each subgrid to the sampled rainfall distributions.                  | 43 |
| 663 | 3 | For a precipitation threshold of 0.25 mm $h^{-1}$ , the optimal probability thresh- |    |
| 664 |   | olds (OT) and binary verification statistics calculated at the OT for the SSEF      |    |
| 665 |   | and each machine learning model in grid 1.  | 44 |
| 666 | 4 | For a precipitation threshold of 0.25 mm $h^{-1}$ , the optimal probability thresh- |    |
| 667 |   | olds (OT) and binary verification statistics calculated at the OT for the SSEF      |    |
| 668 |   | and each machine learning model in grid 2.  | 45 |
| 669 | 5 | Attributes diagrams for the raw ensemble and each machine learning model.           |    |
| 670 |   | The solid line with circles is the reliability curve, which indicates the observed  |    |
| 671 |   | relative frequency of rain events above the 6.35 mm threshold for each prob-        |    |
| 672 |   | ability bin. The dot-dashed line with triangles shows the relative frequency        |    |
| 673 |   | of all forecasts that fall in each probability bin. If a point on the reliability   |    |
| 674 |   | curve falls in the gray area, it contributes positively to the BSS. The horizon-    |    |
| 675 |   | tal and vertical dashed lines are located at the climatological frequency for a     |    |
| 676 |   | particular sub grid. The diagonal dashed line indicates perfect reliability.        | 46 |
| 677 | 6 | Relative Operating Characteristic (ROC) curves for the raw ensemble and             |    |
| 678 |   | each machine learning model. The diagonal lines indicate PSS.                       | 47 |

| 679 | 7  | BSS and AUC comparisons by hour.   | 48 |
|-----|----|--|----|
| 680 | 8  | Evaluation of random forests trained with 3 precipitation thresholds at each           |    |
| 681 |    | forecast hour. Brier Skill Score (BSS) is the evaluation metric.                       | 49 |
| 682 | 9  | Observed 1-hour precipitation on 13 May 2010 at 02 UTC.                                | 50 |
| 683 | 10 | Filled contours of 2 m dew point and 2 m temperature for the 2-hour and                |    |
| 684 |    | 26-hour SSEF forecasts valid on 13 May 2010 at 0200 UTC.                               | 51 |
| 685 | 11 | The 2-hour (left) and 26-hour (right) forecasts of the SSEF ensemble proba-            |    |
| 686 |    | bility (top) and the simple logistic regression (bottom). The green contours           |    |
| 687 |    | indicate the observed areas of 1-hour precipitation greater than 6.35 mm $\rm h^{-1}.$ | 52 |
| 688 | 12 | The 2-hour (left) and 26-hour (right) forecasts of the multiple logistic regres-       |    |
| 689 |    | sion (top) and the random forest (bottom). The green contours indicate the             |    |
| 690 |    | observed areas of 1-hour precipitation greater than 6.35 mm $h^{-1}$ .                 | 53 |
| 691 | 13 | Maps of the SSEF 26-hour forecasts of predictors considered important by the           |    |
| 692 |    | random forest model.   | 54 |



FIG. 1. Map of the number of grid points sampled within a  $400 \text{ km}^2$  box spatially averaged with a 300 km radius median filter. The domain sub grids are also shown and labeled.



FIG. 2. Histogram comparing the relative frequencies of the full precipitation distribution for the each subgrid to the sampled rainfall distributions.



FIG. 3. For a precipitation threshold of 0.25 mm  $h^{-1}$ , the optimal probability thresholds (OT) and binary verification statistics calculated at the OT for the SSEF and each machine learning model in grid 1.



FIG. 4. For a precipitation threshold of 0.25 mm  $h^{-1}$ , the optimal probability thresholds (OT) and binary verification statistics calculated at the OT for the SSEF and each machine learning model in grid 2.



FIG. 5. Attributes diagrams for the raw ensemble and each machine learning model. The solid line with circles is the reliability curve, which indicates the observed relative frequency of rain events above the 6.35 mm threshold for each probability bin. The dot-dashed line with triangles shows the relative frequency of all forecasts that fall in each probability bin. If a point on the reliability curve falls in the gray area, it contributes positively to the BSS. The horizontal and vertical dashed lines are located at the climatological frequency for a particular sub grid. The diagonal dashed line indicates perfect reliability.



FIG. 6. Relative Operating Characteristic (ROC) curves for the raw ensemble and each machine learning model. The diagonal lines indicate PSS.



FIG. 7. BSS and AUC comparisons by hour.



FIG. 8. Evaluation of random forests trained with 3 precipitation thresholds at each forecast hour. Brier Skill Score (BSS) is the evaluation metric.



Observed Precipitation Valid at 13 May 2010 0200 UTC

FIG. 9. Observed 1-hour precipitation on 13 May 2010 at 02 UTC.



FIG. 10. Filled contours of 2 m dew point and 2 m temperature for the 2-hour and 26-hour SSEF forecasts valid on 13 May 2010 at 0200 UTC.



FIG. 11. The 2-hour (left) and 26-hour (right) forecasts of the SSEF ensemble probability (top) and the simple logistic regression (bottom). The green contours indicate the observed areas of 1-hour precipitation greater than 6.35 mm  $h^{-1}$ .



FIG. 12. The 2-hour (left) and 26-hour (right) forecasts of the multiple logistic regression (top) and the random forest (bottom). The green contours indicate the observed areas of 1-hour precipitation greater than 6.35 mm  $h^{-1}$ .



FIG. 13. Maps of the SSEF 26-hour forecasts of predictors considered important by the random forest model.