

Money Doesn't Grow on Trees, but Forecasts Do: Forecasting Extreme Precipitation with Random Forests

GREGORY R. HERMAN AND RUSS S. SCHUMACHER

Department of Atmospheric Science, Colorado State University, Fort Collins, Colorado

(Manuscript received 30 August 2017, in final form 29 March 2018)

ABSTRACT

Approximately 11 years of reforecasts from NOAA's Second-Generation Global Ensemble Forecast System Reforecast (GEFS/R) model are used to train a contiguous United States (CONUS)-wide gridded probabilistic prediction system for locally extreme precipitation. This system is developed primarily using the random forest (RF) algorithm. Locally extreme precipitation is quantified for 24-h precipitation accumulations in the framework of average recurrence intervals (ARIs), with two severity levels: 1- and 10-yr ARI exceedances. Forecasts are made from 0000 UTC forecast initializations for two 1200–1200 UTC periods: days 2 and 3, comprising, respectively, forecast hours 36–60 and 60–84. Separate models are trained for each of eight forecast regions and for each forecast lead time. GEFS/R predictors vary in space and time relative to the forecast point and include not only the quantitative precipitation forecast (QPF) output from the model, but also variables that characterize the meteorological regime, including winds, moisture, and instability. Numerous sensitivity experiments are performed to determine the effects of the inclusion or exclusion of different aspects of forecast information in the model predictors, the choice of statistical algorithm, and the effect of performing dimensionality reduction via principal component analysis as a preprocessing step. Overall, it is found that the machine learning (ML)-based forecasts add significant skill over exceedance forecasts produced from both the raw GEFS/R ensemble QPFs and from the European Centre for Medium-Range Weather Forecasts' (ECMWF) global ensemble across almost all regions of the CONUS. ML-based forecasts are found to be underconfident, while raw ensemble forecasts are highly overconfident.

1. Introduction

Locally extreme precipitation can cause a variety of costly, disruptive, and endangering impacts, including flooding, flash flooding, and landslides. In 2016 alone, these hazards combined caused more than 120 fatalities and \$10 billion in damage over the United States (NWS 2017b). The prediction of flash floods is a notoriously challenging forecast problem, requiring accurate prediction not only of heavy rainfall magnitudes, but also of the spatiotemporal distribution of that rainfall; of the hydrologic interactions among precipitation, terrain, and the land surface; and of antecedent precipitation and its effects on soil conditions. Forecasting precipitation processes responsible for most observed extreme rainfall over the contiguous United States (CONUS) is often considered among the most challenging problems in contemporary numerical weather

prediction (NWP; e.g., Fritsch and Carbone 2004; Novak et al. 2014). Given that the rainfall forecast alone presents such a considerable challenge, the additional hydrologic considerations in the flash flood forecast problem present an even more daunting task. While recent advances in heavy rainfall and flash flood forecasting have been made (e.g., Hapuarachchi et al. 2011; Novak et al. 2014; Barthold et al. 2015), forecasts still struggle in many situations (e.g., Delrieu et al. 2005; Lackmann 2013; Schumacher et al. 2013; Gochis et al. 2015; Nielsen and Schumacher 2016; among many others), and substantial progress remains to be made.

Contemporary operational dynamical forecast models often struggle to simulate accurately the physical processes responsible for extreme precipitation production. For example, models with parameterized convection often have a variety of persistent errors and biases associated with their depiction of convective systems, which are responsible for the majority of flooding rains over much of the CONUS (e.g., Schumacher and Johnson 2006; Stevenson and Schumacher 2014;

Corresponding author: Gregory R. Herman, gherman@atmos.colostate.edu

Herman and Schumacher 2016a). These include a tendency to underpredict total rainfall from convective systems (e.g., Schumacher and Johnson 2008; Herman and Schumacher 2016a); produce systems displaced too far to the north and west from where they are observed (e.g., Grams et al. 2006; Wang et al. 2009; Clark et al. 2010); initiate convection too early (e.g., Davis et al. 2003; Wilson and Roberts 2006; Clark et al. 2007); generate systems with too large an areal extent (e.g., Wilson and Roberts 2006); and propagate them incorrectly, too slowly, or not at all (e.g., Davis et al. 2003; Pinto et al. 2015). While convection-allowing models (CAMs) can better resolve the physical processes responsible for heavy rainfall generation (e.g., Kain et al. 2006; Weisman et al. 2008; Duda and Gallus 2013), they too can suffer from many of these biases (e.g., Kain et al. 2006; Lean et al. 2008; Kain et al. 2008; Weisman et al. 2008; Herman and Schumacher 2016a). Furthermore, although there is a plethora of CAM guidance out to the day-ahead time frame (out to 36 h to perhaps 48 h after initialization), due to current computational constraints, there is almost no operational CAM guidance running out to 2 days ahead and nothing operational that runs to 3 days ahead or beyond. Instead, global ensembles with parameterized convection serve as the primary source of forecast information and uncertainty quantification at these lead times. Nevertheless, there is considerable utility in skillful extreme precipitation forecasts at these longer lead times, since many mitigative actions may not be feasible to execute in a matter of hours but are easily accomplished with a day or more of warning. Statistical postprocessing of global ensemble output can potentially alleviate many of these dynamical model deficiencies and provide skillful extreme precipitation guidance at medium-range time scales. A specific focus on the day 2–3 period is warranted due to the increased existing operational emphasis on these lead times, compared with even longer ones, such as the excessive rainfall outlook produced by the Weather Prediction Center (Barthold et al. 2015), which forecasts locally excessive rainfall across the CONUS for days 1–3.

There is a long history of successful application of statistical postprocessing to dynamical model output (e.g., Klein et al. 1959; Glahn and Lowry 1972). Model output statistics (MOS; e.g., Glahn and Lowry 1972) is a simple, effective multivariate linear regression technique relating a set of dynamical model predictors to sensible weather predictands, such as minimum and maximum temperature, wind speeds, and precipitation probability. This basic technique has long demonstrated skill over both the underlying models and even human forecasters (e.g., Jacks et al. 1990; Vislocky and Fritsch 1997; Hamill et al. 2004; Baars and Mass 2005) but is

inherently limited by the linear assumptions underlying the method. Statistical postprocessing techniques have also been successfully applied to QPFs, from early linear approaches (e.g., Bermowitz 1975; Antolik 2000) to more contemporary techniques that can exploit more complex variable relationships, including neural networks (e.g., Hall et al. 1999), reforecast analogs (e.g., Hamill and Whitaker 2006; Hamill et al. 2015), logistic regression (LR; e.g., Applequist et al. 2002; Whan and Schmeits 2018, manuscript submitted to *Mon. Wea. Rev.*), random forests (RFs; e.g., Gagne et al. 2014; Ahijevych et al. 2016; Gagne et al. 2017; Whan and Schmeits 2018, manuscript submitted to *Mon. Wea. Rev.*), and other parametric techniques (e.g., Scheuerer and Hamill 2015; Whan and Schmeits 2018, manuscript submitted to *Mon. Wea. Rev.*). For other meteorological applications, other machine learning algorithms, such as support vector machines (e.g., Zeng and Qiao 2011; Herman and Schumacher 2016b) and boosting (e.g., Herman and Schumacher 2016b; Hong et al. 2016) have also been successfully applied. Related techniques have also been applied to forecasting related high-impact phenomena, such as severe hail (Brimelow et al. 2006; Gagne et al. 2015) and tornadoes (Alvarez 2014). One of the most powerful aspects of machine learning algorithms—and RFs in particular—is finding patterns and nonlinear interactions in the supplied training data (e.g., Breiman 2001). Depending on the extent and diversity of the data supplied in these experiments, trained RFs pose the theoretical capability of diagnosing and automatically correcting for various kinds of model biases, including context-dependent quantitative biases, such as QPF being systematically too high or too low, spatial displacement biases in the placement of extreme precipitation features, and, to some extent, temporal biases in the initiation or progression of extreme precipitation features.

This study makes a comprehensive investigation of using a global reforecast dataset to produce skillful and reliable probabilistic forecasts of locally extreme precipitation using the RF statistical postprocessing technique in the medium range. The following section provides further background and rigorously describes the data and methods used, algorithms employed, models trained, and experiments performed. Section 3 presents results of the sensitivity experiments conducted, while section 4 presents the final results of the trained models and provides two brief case studies illustrating the process. Section 5 summarizes the findings of this study, outlines complementary analysis of these models, identifies avenues for further research, and discusses the broader implications of the results on numerical weather prediction and postprocessing.

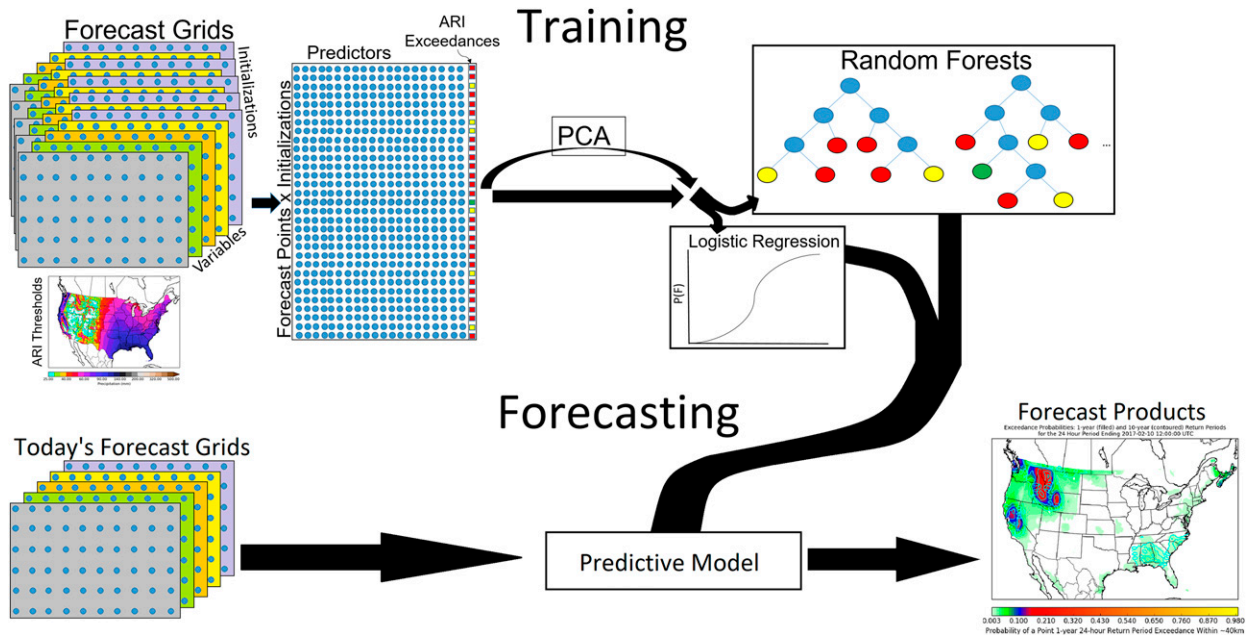


FIG. 1. Schematic representation of the forecast process for this study. GEFS/R forecasts are taken, assembled across fields, space, and time to form a training matrix, and past observations are used to associate a label with each forecast initialization–forecast day–forecast point triplet. The training matrix optionally undergoes preprocessing through PCA and then is input to one or more machine learning algorithms. From here, probabilistic ARI exceedance forecasts may be readily generated.

2. Data and methods

There are several successive steps applied in creating the final forecasts evaluated in this study. A schematic overview of the forecast pipeline for the models trained in this study is depicted in Fig. 1. Many types of hydrometeorological information are first taken, then assembled in a methodical manner, further preprocessed for subsequent analysis, analyzed using a statistical machine learning algorithm, and finally, extreme precipitation forecast guidance is produced and evaluated. This section details each of these steps in the model development and evaluation process.

a. Datasets

Dynamical model data used for training the RF models in this study come from NOAA’s Second-Generation Global Ensemble Forecast System Reforecast (GEFS/R; Hamill et al. 2013) dataset. The GEFS/R is a global 11-member ensemble with parameterized convection and T254L42 resolution—which corresponds to an effective horizontal grid spacing of ~ 55 km at 40° latitude—initialized once daily at 0000 UTC back to December 1984. Perturbations are applied only to the initial conditions and are made using the ensemble transform with rescaling technique (Wei et al. 2008). The ensemble system used to generate these reforecasts is nearly static throughout its 30+ year period of coverage,

though updates to the operational data assimilation system over time have resulted in some changes in the bias characteristics of its forecasts over the period of record (Hamill 2017). Some forecast fields are preserved on the native Gaussian grid ($\sim 0.5^\circ$ spacing), while others are available only on a $1^\circ \times 1^\circ$ grid. Temporally, forecast fields are archived every 3 h out to 72 h past initialization and are available every 6 h beyond that. This study employs an almost 11-yr period of record to explore this forecast problem, using daily initializations from January 2003 through August 2013.

In creating probabilistic extreme precipitation forecast guidance, the predictand must first be concretely specified and a robust, consistent verification framework established. One of the many challenges in heavy rainfall and flash flood forecasting is the considerable difficulty in verifying events (e.g., Welles et al. 2007; Gourley et al. 2012; Barthold et al. 2015), as every approach has its deficiencies and limitations. It is attractive to consider the problem from a simple perspective of quantitative precipitation estimate (QPE) exceedances of some temporally static threshold. In particular, a fixed threshold (e.g., 50 mm h^{-1}) can be used as a proxy for flash flooding (e.g., Brooks and Stensrud 2000; Hitchens et al. 2013), as can exceedances of thresholds defined relative to the local precipitation climatology (e.g., Schumacher and Johnson 2006; Stevenson and Schumacher 2014; Herman and Schumacher 2016a),

such as average recurrence intervals (ARIs). An ARI defines a fixed frequency relative to the hydrometeorological climatology of the region; in particular, it corresponds to the expected duration, given the local climatology, between exceedances of a given threshold. For example, the 1-yr ARI for 24-h precipitation accumulations describes the accumulation amount for which one would expect the mean duration between exceedances of said amount to be 1 year. Past research has shown that a fixed-frequency ARI-based framework has better correspondence with heavy precipitation impacts than the use of any fixed threshold across the hydro-meteorologically diverse regions of the CONUS (e.g., Reed et al. 2007). From the perspective of forecast verification, defining extreme precipitation with respect to a fixed threshold exceedance raises challenges when applied uniformly across the CONUS. For example, skill differences observed between regions may simply be an artifact of a regionally varying event climatology rather than “true” regional differences in forecast skill (e.g., Hamill and Juras 2006). The ARI framework avoids this issue and provides reasonable correspondence with precipitation impacts while avoiding the additional complications, such as antecedent conditions, local hydrology, and urban effects (e.g., Herman and Schumacher 2016a), and is consequently used to quantify extreme rainfall for this study.

Specifically, forecast probabilities are issued for 24-h ARI exceedances at each GEFS/R archive grid point on its native Gaussian grid at all points across the CONUS, using a predictand with three categories: 1) no 1-yr ARI exceedance at any point within the gridpoint domain, 2) at least one 1-yr ARI exceedance, but no 10-yr ARI exceedances within the gridpoint domain, and 3) at least one 10-yr ARI exceedance within the gridpoint domain. For evaluation, probabilities from the middle and most severe categories are often aggregated to produce a 1-yr ARI exceedance probability. This approach has the advantage of retaining aspects of the anticipated event severity as would be retained in a regression context, but is largely lost when performing single-category classification. While there can be some additional complications, especially with respect to calibration, formulating the prediction problem as a single multicategory classification task rather than multiple distinct binary category models also ensures mathematical consistency of the exceedance probabilities within the generated probability mass functions in a way that the latter approach would not.

In aggregating multiple QPE-to-ARI threshold gridpoint comparisons in a single predictand, the forecasts issued correspond to neighborhood event probabilities, an increasingly popular method of

communicating probabilistic high-impact weather information in forecast operations (e.g., Barthold et al. 2015; NWS 2017a). Counting any one of several possible point exceedances as an “event” results in the event having a higher observed relative frequency relative to that of any of the individual point exceedances; the event frequency in this framework thus exceeds the purported frequency suggested by the ARI. However, the fixed-frequency property, and thus many of the aforementioned desirable properties of the framework, is approximately retained. For this study, focus is placed exclusively on two 24-h forecast periods: the 1200–1200 UTC period corresponding to forecast hours 36–60 from the GEFS/R forecast fields and the subsequent 24-h period encompassing forecast hours 60–84, denoted respectively as days 2 and 3. At these times, there is typically some knowledge to characterize the environmental conditions in which precipitation may form, but it is beyond the current range of operational CAM guidance.

Verification comes from the National Centers for Environmental Prediction (NCEP) stage IV precipitation analysis (Lin and Mitchell 2005) QPE product, created operationally since December 2001. Stage IV provides 24-h analyses over the CONUS on a ~4.75-km grid. It uses both rain gauge observations and radar-derived rainfall estimates to generate an analysis and is further quality controlled via NWS river forecast centers (RFCs) to ensure stray radar artifacts and other spurious anomalies do not appear in the final product. Despite some limitations (Herman and Schumacher 2016a; Nelson et al. 2016), its analysis quality, resolution, allowing better ability to capture precipitation extremes compared with other QPE products (e.g., Hou et al. 2014), and data record length make it preferable to analogous products.

The ARI thresholds associated with the 1- and 10-yr ARIs for 24-h precipitation accumulations are generated using the same methodology as Herman and Schumacher (2016a), where CONUS-wide thresholds are produced by stitching thresholds from several sources. NOAA’s Atlas 14 thresholds (Bonnin et al. 2004, 2006; Perica et al. 2011, 2013), an update from older work and currently under development, are used wherever they were available at the commencement of this study. For five northwestern states—Washington, Oregon, Idaho, Montana, and Wyoming—updated thresholds are not available, and derived NOAA Atlas 2 threshold estimates are used instead (Miller et al. 1973). Additionally, in Texas and the Northeast—New York, Vermont, New Hampshire, Maine, Massachusetts, Connecticut, and Rhode Island—Technical Paper 40

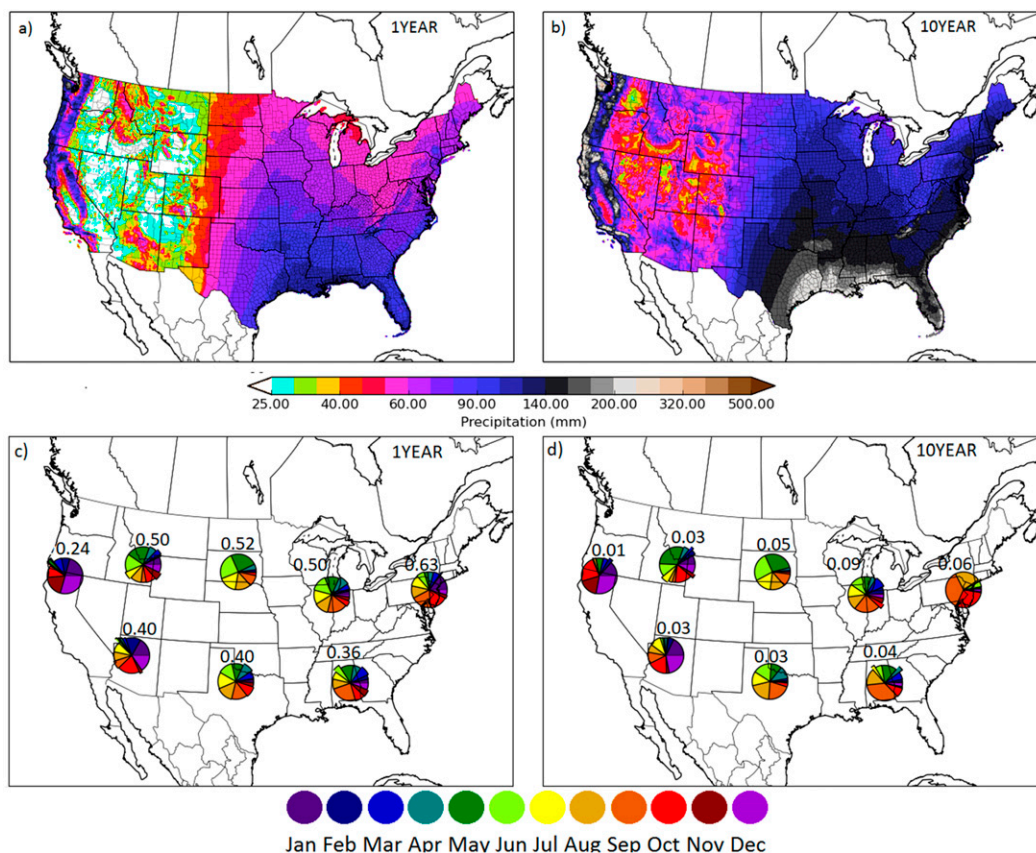


FIG. 2. ARI thresholds at the (a) 1- and (b) 10-yr ARI levels over the CONUS for a 24-h accumulation interval. Climatology of observed exceedances of the (c) 1-yr, 24-h and (d) 10-yr, 24-h ARI thresholds between January 2003 and August 2013 based on stage IV precipitation analysis. Pie charts indicate the monthly distribution of event occurrence within each study region, as shown in Fig. 3. Numbers above the pie charts indicate the mean number of exceedances per point per year within the region (a priori 1 and 0.1 for 1- and 10-yr ARIs, respectively).

(TP-40; Hershfield 1961) thresholds are used;¹ everywhere else uses the Atlas 14 threshold estimates. The 10-yr ARI thresholds (Fig. 2b) show a similar spatial pattern to the 1-yr ARI thresholds (Fig. 2a) but are substantially higher everywhere. More significantly, it is apparent that at both severity levels, there are large regional disparities in the threshold magnitudes. Over climatologically wet regions of the CONUS, such as the Pacific coastal mountains and immediately along the Gulf Coast, thresholds are as high as 100–150 and 250–300 mm for 1- and 10-yr ARIs, respectively. Over the central and eastern CONUS, thresholds tend to decrease smoothly with increasing latitude and distance from major bodies of water. Sharper variations are seen in areas of complex terrain over the western CONUS. In

the driest parts of the arid Southwest and Intermountain West, thresholds can be as low as 10–15 and 25–30 mm for the two ARI levels—a full order of magnitude difference from the largest thresholds at the same intensity level.

Forecast models in this study are trained separately for eight distinct yet cohesive and internally fairly hydrometeorologically homogeneous regions of the CONUS, using the delineation indicated in Fig. 3. Observed 1- and 10-yr ARI exceedance events that occurred during the period of record (Figs. 2c,d) highlight important regional differences in the seasonal climatology of ARI exceedances across the CONUS. In the Pacific coast (PCST) region, the vast majority of exceedances at both the 1- and 10-yr severity levels occur in the cool season, largely from atmospheric river events with large moisture transport impinging on coastal topography (e.g., Rutz et al. 2014; Herman and Schumacher 2016a). This seasonality holds to a lesser extent in the neighboring Southwest (SW) region, with some signal

¹ The northeastern states did receive updated Atlas 14 estimates in October 2015, but TP-40 thresholds were retained for consistency with prior work.

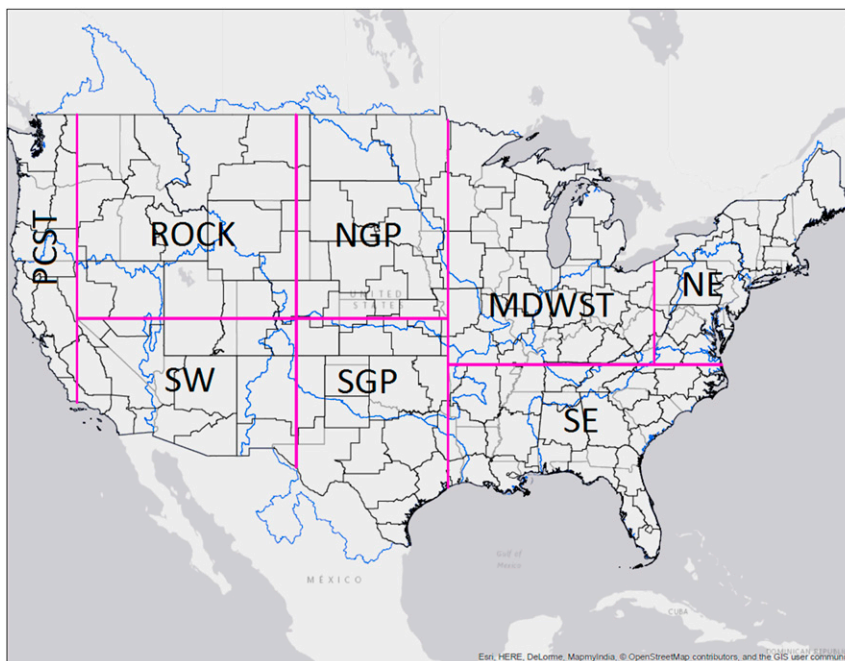


FIG. 3. Map depicting the regional partitioning of the CONUS used in this study and the labels ascribed to each region.

carrying over to the Rockies (ROCK) region as well. In the central and eastern regions, the majority of events occur during the warm season from more scattered convective-scale processes, particularly in the months of May, June, and July (e.g., Schumacher and Johnson 2006; Herman and Schumacher 2016a). Tropical cyclones can cause widespread and very significant rainfall and comprise a substantial portion of the extreme precipitation climatology, especially in the Northeast (NE) and Southeast (SE) regions. Because of the spatial extent of their impacts and immense rainfall totals they can produce, they form a much larger fraction of the climatology of 10-yr ARI exceedances (Fig. 3d) than 1-yr events (Fig. 3c). Additionally, the numbers are lower than would be expected; by the explicit exceedance frequencies associated with the thresholds, one would expect an average of one exceedance per point per year over the period of record for the 1-yr events (Fig. 3c) and 0.1 exceedances for 10-yr events (Fig. 3d). In reality, event counts are only approximately half of that. This is consistent with previous findings (e.g., Herman and Schumacher 2016a) and likely in part attributable to limitations in the stage IV product to capture extremes (e.g., Nelson et al. 2016). There is also quite a bit of region-to-region variability in event counts, particularly for 10-yr exceedances, much of which is attributable to statistical variability from having a short data record in relation to the event frequency.

b. Predictor assembly

Input predictors, or features, to the random forests can be partitioned into two categories: model predictors and background predictors. The former constitute the vast majority of inputs. Model predictors come from atmospheric fields forecast in the GEFS/R, which bear a known physical relationship with extreme precipitation. A core set of $f = 9$ fields used in this study are accumulated precipitation (APCP), convective available potential energy (CAPE), convective inhibition (CIN), precipitable water (PWAT), surface temperature (T2M) and specific humidity (Q2M), surface zonal (U10) and meridional winds (V10), and mean sea level pressure (MSLP). Sensitivity experiments explore the use of additional upper-air atmospheric fields; a full list of fields used in this study, their associated symbols used in this manuscript, and the grids on which they are each archived is included in Table 1. The spatiotemporal variations in these fields are considered as well. Spatially, predictors are structured in a forecast-point relative sense. In the control model, GEFS/R forecast values up to $r = 4$ grid boxes ($\sim 2^\circ$) latitudinally or longitudinally displaced in any direction relative to the forecast point are considered. Temporally, simulated fields are considered at each archive time during the forecast interval, which corresponds to every 3 h during the day 2 period and every 6 h during the day 3 period, for a total of $t = 9$

TABLE 1. Summary of dynamical model fields examined in this study, including the abbreviated symbol to which each variable is referred throughout the paper, a description of each variable, the predictor group with which the field is associated in the manuscript text, and the highest resolution for which the field can be obtained from the GEFS/R.

Symbol	Description	Predictor Group	Grid
APCP	Precipitation accumulation in past (3) 6 h	Core	Native Gaussian
CAPE	Surface-based convective available potential energy	Core	Native Gaussian
CIN	Surface-based convective inhibition	Core	Native Gaussian
MSLP	Mean sea level pressure	Core	Native Gaussian
PWAT	Total precipitable water	Core	Native Gaussian
Q2M	Specific humidity 2 m above ground	Core	Native Gaussian
T2M	Air temperature 2 m above ground	Core	Native Gaussian
U10	Zonal component of 10-m wind	Core	Native Gaussian
V10	Meridional component of 10-m wind	Core	Native Gaussian
Q300	Specific humidity at 300 hPa	Upper-Air Extra	1° × 1°
Q500	Specific humidity at 500 hPa	Upper-Air Core	1° × 1°
Q700	Specific humidity at 700 hPa	Upper-Air Extra	1° × 1°
Q850	Specific humidity at 850 hPa	Upper-Air Core	1° × 1°
T250	Temperature at 250 hPa	Upper-Air Extra	1° × 1°
T500	Temperature at 500 hPa	Upper-Air Core	1° × 1°
T700	Temperature at 700 hPa	Upper-Air Extra	1° × 1°
T850	Temperature at 850 hPa	Upper-Air Core	1° × 1°
U250	Zonal component of 250-hPa wind	Upper-Air Extra	1° × 1°
U500	Zonal component of 500-hPa wind	Upper-Air Core	1° × 1°
U700	Zonal component of 700-hPa wind	Upper-Air Extra	1° × 1°
U850	Zonal component of 850-hPa wind	Upper-Air Core	1° × 1°
V250	Meridional component of 250-hPa wind	Upper-Air Extra	1° × 1°
V500	Meridional component of 500-hPa wind	Upper-Air Core	1° × 1°
V700	Meridional component of 700-hPa wind	Upper-Air Extra	1° × 1°
V850	Meridional component of 850-hPa wind	Upper-Air Core	1° × 1°
W850	Vertical velocity (omega) at 850 hPa	Upper-Air Core	1° × 1°

and 5 forecast periods for the day 2 and 3 periods, respectively. All told, this yields $tf(2r + 1)^2$ model predictors, which yields, respectively, $M = 6561$ and 3645 model predictors for the day 2 and day 3 control models. The other category of predictors—background predictors (Table 2)—are those that are solely associated with the forecast point and have no relation to the present meteorology. These include the location of the

point, as well as the ARI characteristics of the point and in the surrounding area.

c. Dimensionality reduction

There are a large number of model predictors, and they are also highly correlated—spatially, temporally, and across variables. With millions of training examples and thousands of features, the forecast problem can

TABLE 2. List of background predictors used in this study and their associated symbols and descriptions.

Symbol	Description
ARI1_LOCAL_MEDIAN	Median of 1-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI1_LOCAL_MIN	Minimum of 1-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI1_LOCAL_MAX	Maximum of 1-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI10_LOCAL_MEDIAN	Median of 10-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI10_LOCAL_MIN	Minimum of 10-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI10_LOCAL_MAX	Maximum of 10-yr ARIs whose closest GEFS/R grid point is the forecast point.
ARI1_REGIONAL_MEDIAN	Median of 1-yr ARIs that lie within the domain from which model predictors are drawn.
ARI1_REGIONAL_MIN	Minimum of 1-yr ARIs that lie within the domain from which model predictors are drawn.
ARI1_REGIONAL_MAX	Maximum of 1-yr ARIs that lie within the domain from which model predictors are drawn.
ARI10_REGIONAL_MEDIAN	Median of 10-yr ARIs that lie within the domain from which model predictors are drawn.
ARI10_REGIONAL_MIN	Minimum of 10-yr ARIs that lie within the domain from which model predictors are drawn.
ARI10_REGIONAL_MAX	Maximum of 10-yr ARIs that lie within the domain from which model predictors are drawn.
LAT	Latitude of forecast point
LON	Longitude of forecast point

become computationally intractable. Further, having many highly correlated features can readily result in model overfitting—making predictions based on noise affecting an individual native feature rather than the underlying signal—a phenomenon commonly termed the “curse of dimensionality” (e.g., [Friedman 1997](#)). There are numerous ways these concerns can be addressed; broadly speaking, the most common approaches are either feature selection or feature extraction. In feature selection, a subset of initial predictors is chosen that collectively bears the strongest predictive relationship with the predictand, whereas in feature extraction, a smaller set of new predictors is derived from the original set. Both of these procedures can be performed subjectively through manual means or objectively through automated means. In this case, all of the input predictors are believed to have a physical relationship with extreme precipitation, and choosing only the most predictive fields (e.g., model QPF) and discarding the rest risks removing valuable predictive information not contained in the retained predictor set. The primary issue with the input predictors in this case is not that many may not have any physical bearing on the predictand, but rather that each predictor represents a value at a different point of a continuous field or a different property at the same point, and all are thus necessarily highly correlated to one another. Furthermore, while one could conceivably extract features using field averages or some other predetermined method, this may not be optimal. For example, it may be better to weight values closer to the forecast point more heavily, while still retaining some information from the far-field predictors. Given the uncertainty in optimally constructing features by manual means, it is more convenient and repeatable to instead extract features objectively. Though it has some limitations (e.g., [Shlens 2014](#)), principal components analysis (PCA; [Ross et al. 2008](#); [Pedregosa et al. 2011](#)) is a robust and frequently utilized approach for dimensionality reduction. This creates a small set of uncorrelated predictors that explains the signal in the forecast data and gives insight into the regional modes of atmospheric variability as depicted in the GEFS/R model [explored in more depth in [Herman and Schumacher \(2018\)](#)], while leaving the noise in lower-order principal components (PCs), acting in principle to both alleviate overfitting and manage computational requirements.

d. Machine learning algorithms and sensitivity experiments

The primary statistical algorithm used in this study is random forests ([Breiman 2001](#)). RFs are in essence an ensemble of decision trees, where traditionally each tree

individually makes a deterministic prediction about the outcome of the predictand; the relative frequencies of each possible predictand outcome in the ensemble of trees are then used to make a probabilistic forecast. Much further detail on tree and RF construction and mechanics can be found in [appendix A](#), as well as in [McGovern et al. \(2017\)](#) and other sources. There are also several parameters that can be tuned to the particular forecast problem in order to maximize model performance. Fourfold cross validation is used for model development in this study, whereby each model configuration examined is trained four times—once each on three-quarters of the training data—and then evaluated on the final withheld quarter. To avoid issues of sample independence and approximately mimic information that would be available in an operational context, 974 consecutive initializations are used for each quarter of training data. All parameter settings and sensitivity experiments are evaluated in this framework. The set of RF parameters tuned is described in [appendix A](#), and the results are presented in [appendix B](#).

In this study, there are a great number of dynamical model data considered as input information on which the RF can base a prediction. A suite of sensitivity experiments are conducted, as summarized in [Table 3](#), in order to investigate which aspects of forecast information contribute most to forecast skill. Experiments include exploring the following:

- Sensitivity to the inclusion of horizontal variations in atmospheric fields by varying the previously described predictor radius parameter R from 0 to 4.
- Sensitivity to the inclusion of additional upper-air atmospheric fields by comparing the inclusion and exclusion of two sets of fields, as noted explicitly in [Table 1](#). The first incorporates temperature, specific humidity, zonal and meridional winds at 850 and 500 hPa, and 850-hPa vertical velocity in the so-called Upper-Air Core predictor group, while an additional experiment further includes those same fields at 700 and 250 hPa.
- Sensitivity of predictor temporal resolution. Predictor density is 3-hourly for day 2 guidance and 6-hourly for day 3 guidance; models are additionally trained with predictors at 12-hourly temporal density for both lead times and 6-hourly temporal density for the day 2 forecast model and compared against the control versions.
- Sensitivity to the type and the extent of use of ensemble information, a question that has implications for how operational centers allocate their computational resources. Using forecast information from only the GEFS/R’s control member in model training

TABLE 3. Summary of the models trained in this study and the corresponding names designated to the models. An “×” indicates the process is performed or the information is used; a lack of one indicates the opposite. MEDIAN corresponds to the ensemble median, CTRL corresponds to the ensemble control member’s fields, and CNFDB uses the median in addition to the second-from-lowest and second-from-highest member values for each field. Horizontal radius is listed in grid boxes from forecast point; time step denotes the number of hours between GEFS/R forecast field predictors. Slashes indicate the first number applies to the day 2 version of the model, while the latter number applies to the day 3 version. Letters enclosed by parentheses indicate subversions of models, with one parameter changed to the value adjacent to the letter. Asterisks indicate a model applies only to day 2 and not day 3. Otherwise, models apply to all eight forecast regions and have both day 2 and day 3 versions. Those models with boldface names are incorporated into the weighted blend of the final model configuration.

Model name	CTL_NPCA	CTL_PCA	UAC_PCA	UAF_PCA	CORE_CNFDB	CORE_CTRL	CORE_LSPACE	CORE_LTIME	CTL_LR
Algorithm	RF	RF	RF	RF	RF	RF	RF	RF	LR
PCA preprocessed		×	×	×	×				×
Uses core fields	×	×	×	×	×	×	×	×	×
Uses UAC fields			×	×					
Uses UAE fields				×					
Ensemble information	MEDIAN	MEDIAN	MEDIAN	MEDIAN	CNFDB	CTRL	MEDIAN	MEDIAN	MEDIAN
Horizontal radius	4	4	4	4	4	4	0 (a), 1 (b), 2 (c), 3 (d)	4	4
Time step	3/6	3/6	3/6	3/6	3/6	3/6	3/6	12 (a), 6 (b*)	3/6

(CTRL) is compared with using the ensemble median from the full ensemble (MEDIAN) and then further with the use of the ensemble second-lowest and second-highest values for each atmospheric field in conjunction with the median (CNFDB) to evaluate the impact of this dimension of forecast information. This follows the findings of Herman and Schumacher (2016b), who found relatively little sensitivity in performance with respect to how ensemble information is used, but using the near-minimum, median, and near-maximum values outperformed using the mean and spread.

- Sensitivity to predictor preprocessing methodology. Models are trained with and without the aforementioned PCA preprocessing step, and an assessment of the effect of this preprocessing step on model skill is made by comparing the two.
- The effect of region size on forecast skill, hypothesizing that models trained for larger regions may exhibit higher skill due to more available training data. This is performed by aggregating the ROCK and SW regions into a new WEST one, combining the southern Great Plains (SGP), northern Great Plains (NGP), and Midwest (MDWST) regions into a CENTRAL region, and collecting SE and NE regions into a single EAST region, while leaving PCST—with its unique extreme precipitation climatology—untouched.
- Sensitivity of model performance as a function of model algorithm, specifically by comparing with logistic regression, a common and comparatively simpler alternative to statistically deriving forecast probabilities. Further discussion of LR and other machine learning alternatives to the RF algorithm is included in appendix A.

e. Model evaluation

Based on the parameter tuning and sensitivity experiment results, final model configurations are selected. The final model is run over a completely withheld 4-yr evaluation period spanning September 2013–August 2017. The forecasts generated from the final model are compared with those from the full ensemble of raw GEFS/R QPFs, as well as the full 50-member ECMWF global ensemble, accessed from TIGGE (Molteni et al. 1996; Bougeault et al. 2010). The comparison with the former provides an assessment of what improvement, if any, these models yield, compared with the raw guidance from which their forecasts are derived when evaluated in a real-time setting. The latter, meanwhile, provides an assessment for how these forecasts compare with state-of-the-science operational ensemble guidance available at these lead times. To make these comparisons, the QPF

from each ensemble member of the two ensembles is regridded onto the ~ 4.75 -km stage IV HRAP grid on which the Atlas thresholds lie using a first-order conservative scheme (Ramshaw 1985). These regridded QPFs are then compared with the 1- and 10-yr ARI thresholds to create deterministic exceedance forecasts with respect to the two thresholds for each ensemble member. These binary grids are then upscaled to the GEFS/R grid using the same procedure as the verification upscaling: any exceedance in the downscaled grid corresponds to an exceedance at the nearest GEFS/R point in the upscaled grid. Since the predictand categories are necessarily mutually exclusive, the 1-yr ARI exceedance grids are modified so that any member forecasting a 10-yr ARI exceedance at a point is not forecasting an exceedance of between 1 and 10 years at that same point and time period. The prevailing operational method of generating forecast probabilities from a dynamical ensemble—democratic voting, whereby the fraction of ensemble members forecasting the event is used as the forecast probability (e.g., Buizza et al. 1999; Eckel 2003)—is applied to each ensemble to generate the exceedance probabilities for the reference forecasts.

Skill, both in the final assessment of model performance as well as in all aforementioned sensitivity experiments, is quantified by means of the rank probability skill score (RPSS) with a climatological reference:

$$\text{RPSS} = 1.0 - \frac{\sum_{d=1}^D \left\{ \sum_{p=1}^P \left[\sum_{m=1}^K \left(\sum_{j=1}^m P_{jpd} - O_{jpd} \right)^2 \right] \right\}}{\sum_{d=1}^D \left\{ \sum_{p=1}^P \left[\sum_{m=1}^K \left(\sum_{j=1}^m P_{\text{clim},j} - O_{jpd} \right)^2 \right] \right\}}, \quad (1)$$

with D forecast days; P forecast points; K predictand categories; P_{jpd} and O_{jpd} correspond, respectively, to the forecast probability and observance of predictand category j on day d and at point p ; and P_{clim} corresponds to the climatological frequency of occurrence, as defined by the respective ARIs of the predictand. A score of 1.0 indicates a perfect forecast, and a score of 0.0 indicates model performance equivalent to forecasting climatology. Final assessment also includes analysis of reliability, both subjectively through reliability diagrams and quantitatively via the Murphy (1973) decomposition of the Brier score (BS), for category j^* :

$$\text{BS}_{j^*} = \sum_{n=1}^N (P_{Nj^*} - O_{Nj^*})^2 = \frac{1}{N} \sum_{c=1}^C N_{cj^*} (P_{cj^*} - \overline{O_{cj^*}})^2 - \frac{1}{N} \sum_{c=1}^C N_{cj^*} (\overline{O_{j^*}} - \overline{O_{cj^*}})^2 + \overline{O_{j^*}}(1 - \overline{O_{j^*}}), \quad (2)$$

where there are $N = DP$ total forecasts, broken into C discrete probability bins with N_c forecasts being issued for each bin c ; $\overline{O_{j^*}}$ denotes the climatological (based on the period of record) frequency of observing event category j^* ; $\overline{O_{cj^*}}$ denotes the proportion of forecasts in probability bin c observing event category j^* , where j^* is the aggregation of event categories of at least j in the RPSS framework, and $\overline{O_{j^*}}(1 - \overline{O_{j^*}})$, the so-called “uncertainty” term, also represents the BS of a climatological forecast. Converting to a Brier skill score (BSS) framework by dividing out by this term,

$$\text{BSS}_{j^*} = 1.0 - \frac{\text{BS}_{j^*}}{\text{BS}_{\text{clim},j^*}} = \frac{\frac{1}{N} \sum_{c=1}^C N_{cj^*} (\overline{O_{j^*}} - \overline{O_{cj^*}})^2}{\underbrace{\overline{O_{j^*}}(1 - \overline{O_{j^*}})}_{\text{Resolution}}} - \frac{\frac{1}{N} \sum_{c=1}^C N_{cj^*} (P_{cj^*} - \overline{O_{cj^*}})^2}{\underbrace{\overline{O_{j^*}}(1 - \overline{O_{j^*}})}_{\text{Reliability}}}. \quad (3)$$

This analysis is conducted for both the 1- and 10-yr thresholds.

Skill calculations and comparisons are made for the host of sensitivity experiments and for each region, lead time, and model configuration. For each comparison, statistical significance is assessed by bootstrapping to obtain identical sets of cases for each of the two forecast sets being compared. Skill scores are derived from the subsample of each forecast set, and a skill difference is computed. This process is repeated 1000 times to generate a distribution of skill differences, and statistical significance is ascertained with respect to whether the 0.5th- and 99.5th-percentile skill score difference values from the bootstrap trials overlap zero. This 99% confidence bound is used in contrast to 90% or 95% bounds to compensate for concerns arising from conducting statistical significance analysis on numerous different comparisons. While some uncertainty analysis is included in the figures presented, much of the statistical significance difference results discussed in the text are omitted for the sake of concision.

3. Results: Sensitivity experiments

Examining forecast skill as a function of time step between atmospheric field predictors (i.e., the CORE_LTIME models of Table 3; Fig. 4a) yields two striking findings concerning 1) the large variations in forecast skill across regions and 2) the evidently low sensitivity of forecast skill to time step length within any given region.

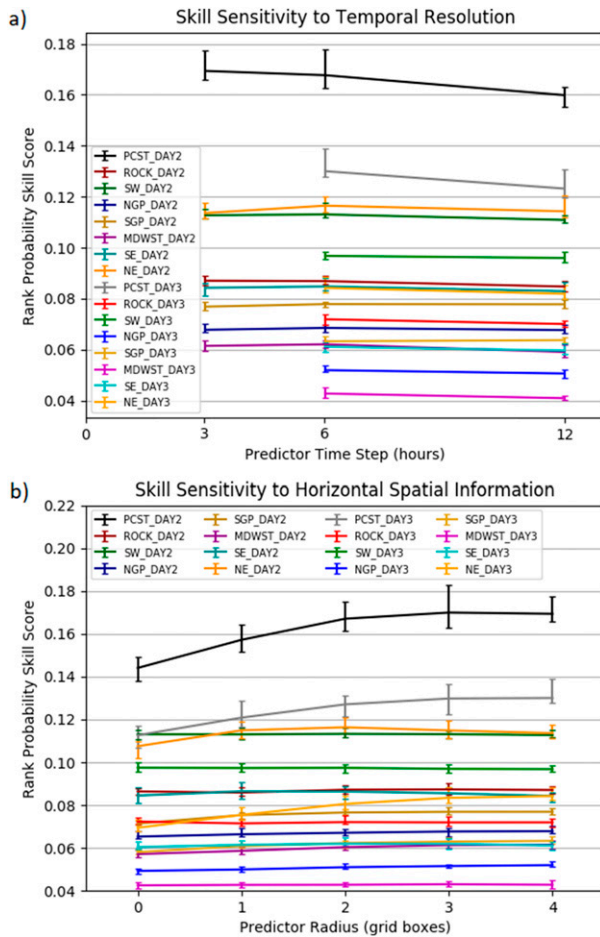


FIG. 4. Sensitivity experiment RPSS results for (a) the CORE_LTIME models, as a function of the time step between incorporation of new atmospheric field forecast values, and (b) the CORE_LSPACE models, as a function of the radius of predictor information incorporated, each including both day 2 and day 3 versions of the model and for each region studied. Lines correspond to a particular day–region pair as indicated in the respective panel legends. Error bars in both panels correspond to 90% confidence bounds obtained by bootstrapping.

For the 3-h time step, predictors are gathered from a total of nine forecast times; with the 6-h step, five forecast times are used, and with the 12-h time step, a total of three forecast times are used. The 12-h time step therefore has one-third of the total predictors as the model with the 3-h time step, but still yields nearly identical forecast skill results. In most regions and forecast periods, there is a slight degradation in performance going from the 6- to 12-h time step, but the difference is not generally statistically significant by a 99% bootstrap skill score difference test (not shown). The one exception to this is in the PCST region, which has much higher skill overall than the other regions for both

forecast periods and exhibits somewhat higher sensitivity to the predictor time step than the other regions, particularly in going from 6 to 12 h, with RPSS differences of approximately 6%.

Similar to the temporal resolution findings, there is a general lack of sensitivity as a function of predictor spatial extent (Fig. 4b). This finding comes in stark contrast to that of Herman and Schumacher (2016b), who found great sensitivity of predictor spatial extent in forecasting airport flight rule conditions. Albeit weak, a slight improvement in skill for most forecast period–region combinations can be noted with increasing predictor radius, often to the extent that the skill difference between 0 and 4 gridbox radii is statistically significant (not shown). Two regions in particular, the NE and PCST, exhibit by far the most sensitivity to predictor spatial extent, with differences of roughly 0.02 observed over the evaluated interval. Also of note is that a radius of 4 grid boxes—the highest number evaluated—did not always yield the best performance results; most notably, the day 2 model for the NE region maximized skill at a radius of 2, with a slight deterioration of forecast skill with increasing radius thereafter. In those regions where the GEFS/R cannot explicitly resolve the processes responsible for producing extreme precipitation, the RF is ultimately making forecasts more on environmental factors; these do not vary drastically in time or space, and thus a single number or small set of numbers at or immediately surrounding the forecast point is sufficient to characterize the basic properties of the environment. This is all that the RF is really using for much of its predictions [see Herman and Schumacher (2018) for more detail]. However, in regions impacted more readily by larger-scale systems where the dynamical model can more directly simulate the precipitation processes, such as PCST and the NE, the spatial variations in atmospheric fields carry more signal rather than noise and thus contribute more predictive value.

Like varying spatial and temporal density, there is relatively little sensitivity to the inclusion of more atmospheric fields (Fig. 5a). Slight but consistent improvement is observed in adding the core upper-air fields as predictors, but adding further levels beyond the core group was found to not improve predictive skill and actually resulted in a decrease in skill for the PCST, NE, ROCK, and SE regions—those that are most affected by larger-scale precipitation systems. Though still rather small, somewhat more distinct sensitivity to type of ensemble information included (Fig. 5b) can be seen here across all regions, with improvements seen using predictor information from the GEFS/R ensemble median versus using only the control member, and slight

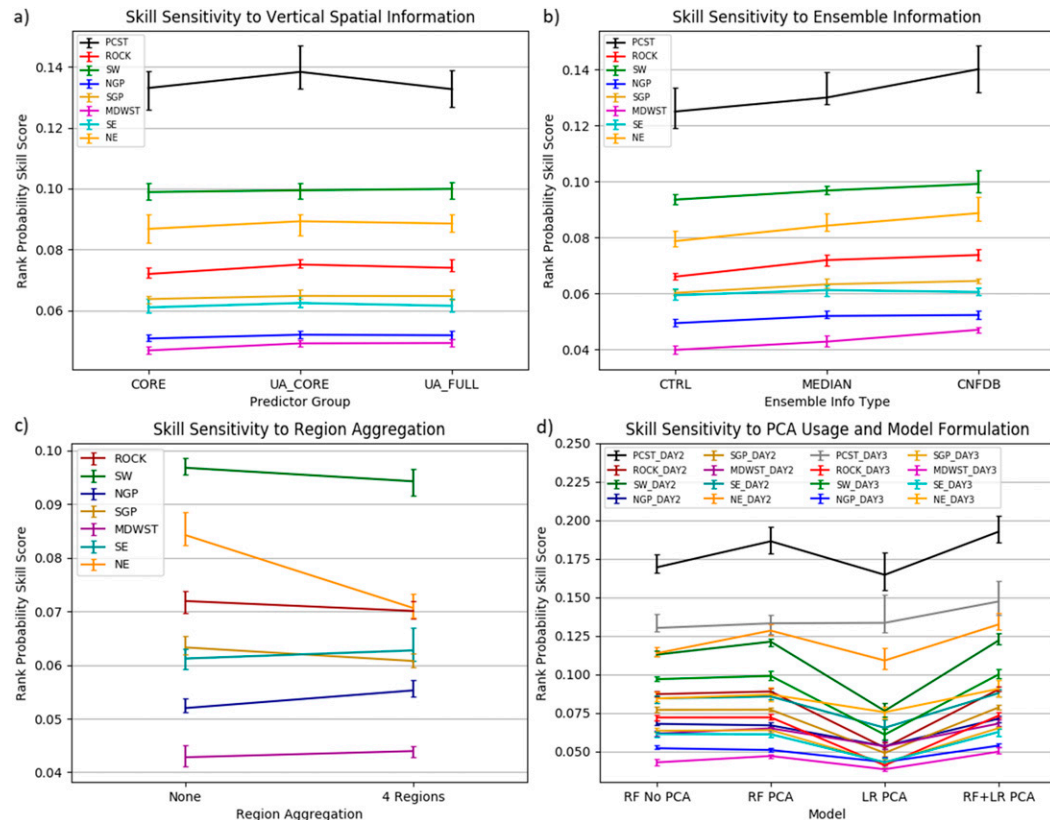


FIG. 5. Sensitivity experiment RPSS results. (a) A function of the atmospheric fields included as input to the RF algorithm for day 3 forecast and broken out by region. From left to right, the columns correspond to results using just the Core atmospheric field group, both the Core and Upper-Air Core groups, and the Core, Upper-Air Core, and Upper-Air Extra groups. For more information on which fields are included in each predictor group, consult Table 1. (b) A function of the type of GEFS/R information used as input predictors to the RF algorithm for day 3 forecasts and broken out by region. From left to right, the columns correspond to results using just the forecast fields from the GEFS/R control member, the ensemble median forecast values from the full ensemble, and the ensemble median, second-from-minimum, and second-from-maximum forecast values from the full ensemble. (c) A function of region aggregation, with the left column using the eight regions depicted in Fig. 3 and the right column using training data that aggregates data from seven of the eight original regions into three regions, as described in the text. (d) A function of model algorithm for different forecast days and regions as indicated in the figure legend. From left to right, columns correspond to results of the CTL_NPCA model, CTL_PCA model, CTL_LR model, and a weighted combination of models as described in the text. For all panels, error bars correspond to 90% confidence bounds obtained by bootstrapping.

further improvement using the ensemble second-from-minimum and second-from-maximum in addition to the ensemble median. The largest differences in magnitude are again for the PCST region, but in this experiment, clear and statistically significant (not shown) improvements are also seen for low-skill, convectively active regions such as MDWST.

Aggregating regions (Fig. 5c) results in a slight degradation in forecast skill. In principle, it is possible for a decision tree to automatically forecast for specific regions by splitting first on the latitude and longitude predictors, and then further partitioning based on meteorological variables thereafter. However, these

findings demonstrate that there is some—albeit limited—utility in manually partitioning training data with distinct hydrometeorological relationships, rather than relying on the machine learning algorithm to discern the distinction automatically. Comparing the impact of applying PCA preprocessing to the RF (Fig. 5d, leftmost two columns) shows that performing PCA tends to either improve performance, as is the case for the PCST, NE, SW, and MDWST regions, or make little difference, as seen in the ROCK, NGP, SGP, and SE regions. The positive differences tend to be larger in magnitude, both in relative and absolute senses, for day 2 model versions, compared with day 3. Forecasts produced through LR

tend to be substantially worse than those generated by RFs (Fig. 5d, center columns). However, the exact magnitude to which this is the case varies by region; substantial differences in skill are seen between RF and LR forecasts for the SW, ROCK, and SGP regions, while there is almost no skill difference between the day 3 forecasts in the PCST region. This may suggest the linear assumptions inherent to the LR algorithm perform better in larger-scale systems than in the more convectively active ones in which the responsible processes are highly nonlinear, but this causality is not entirely clear. Finally, a weighted average of RF and LR forecasts outperforms its component members for all regions and forecast periods. The extent of overperformance is strongly tied to the skill difference between the RF and LR models; when the skill difference is small, the value of the weighted average is comparatively large to when the RF performs much better than LR (cf. Fig. 5d, PCST and SW lines). Since these weighted averages performed the best in cross validation, a weighted average using each of the CTL_NPC, CTL_PCA, and CTL_LR models was chosen for the final model configuration.

4. Results: Final model performance

For both the final ML models and the forecasts from the raw QPFs of both the GEFS/R and ECMWF (Fig. 6), a usually statistically significant deterioration in forecast skill from day 2 to day 3 is evident in each CONUS region over the 4-yr test period. Forecast skill is significantly higher in regions with extreme precipitation associated partially or primarily with synoptic-scale precipitation episodes, such as PCST, SW, and ROCK, rather than smaller-scale convective systems that characterize extreme precipitation, as in the NGP, SGP, and MDWST regions. At an extreme, the NGP and SGP GEFS/R raw QPFs exhibit no skill in predicting ARI exceedances at these lead times. Especially for the ML models, the bigger day 2 versus day 3 skill differences are also seen where the skill is higher, again suggesting the direct forecasting of the precipitation as opposed to forecasts more reflecting the forecast environment, either dynamically via parameterized convection in the case of raw QPFs or directly in the case of the ML model forecasts. Furthermore, the ML models exhibit a larger skill deterioration between days 2 and 3 than either of the raw ensemble forecast sets.

Comparing the forecast systems, the ECMWF forecasts consistently and statistically significantly outperform the GEFS/R forecasts at all lead times, except in the SE region (Fig. 6). Encouragingly, the ML model forecasts are statistically significantly more skillful for all

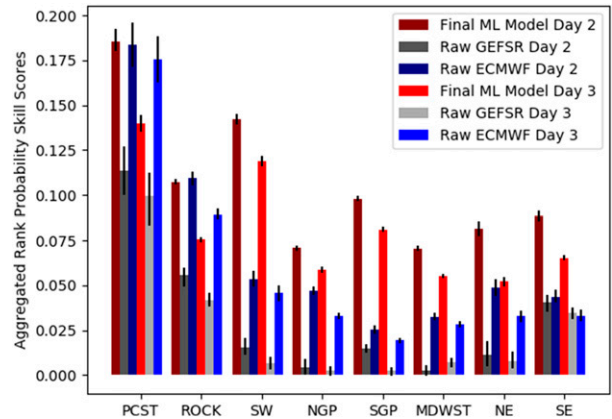


FIG. 6. Final RPSS results obtained over the 4-yr test period spanning September 2013–August 2017, broken out by region. Red bars correspond to the results of the final forecast models trained in this study, while gray bars depict results from the raw GEFS/R QPF probabilities derived from the full ensemble. Dark bars illustrate day 2 performance results, while lighter colors show results for day 3. Error bars correspond to 90% confidence bounds obtained by bootstrapping.

eight regions and both lead times, compared with the GEFS/R forecasts from which they are based. The post-processing is thus clearly accomplishing its purpose of improving forecast skill. But it is also apparent that the GEFS/R is not a state-of-the-science model for extreme QPF prediction, given its lower skill compared with the ECMWF. The real test of the ML model then is how it compares with current best operational guidance for these lead times, represented here with the ECMWF ensemble. The comparison (Fig. 6) is generally quite favorable, with the day 3 ML forecasts outperforming even the day 2 ECMWF forecasts across all regions, except ROCK and PCST. In the nonwestern regions, the extent of overperformance is quite considerable when comparing equal lead times, with skill score improvements of factors of 2 to 3 seen in many comparisons. In the ROCK and PCST regions, the ML and ECMWF forecasts performed about equally at day 2, and ECMWF performed slightly better at day 3. Overall, the ML models demonstrated the ability to consistently outperform current operational model guidance, especially in convectively active regions where there is no operational guidance that can dynamically resolve the physical processes producing extreme precipitation at these lead times.

Reliability diagrams of day 2 raw GEFS/R and ECMWF forecasts (Fig. 7) reveal highly overconfident probabilistic exceedance forecasts for all regions, both severity levels, and both ensembles, as evidenced by the shallow slope relative to the one-to-one line in each panel. The raw GEFS/R forecasts (Figs. 7a,b) are relatively sharp, with more than 0.01% of forecasts falling

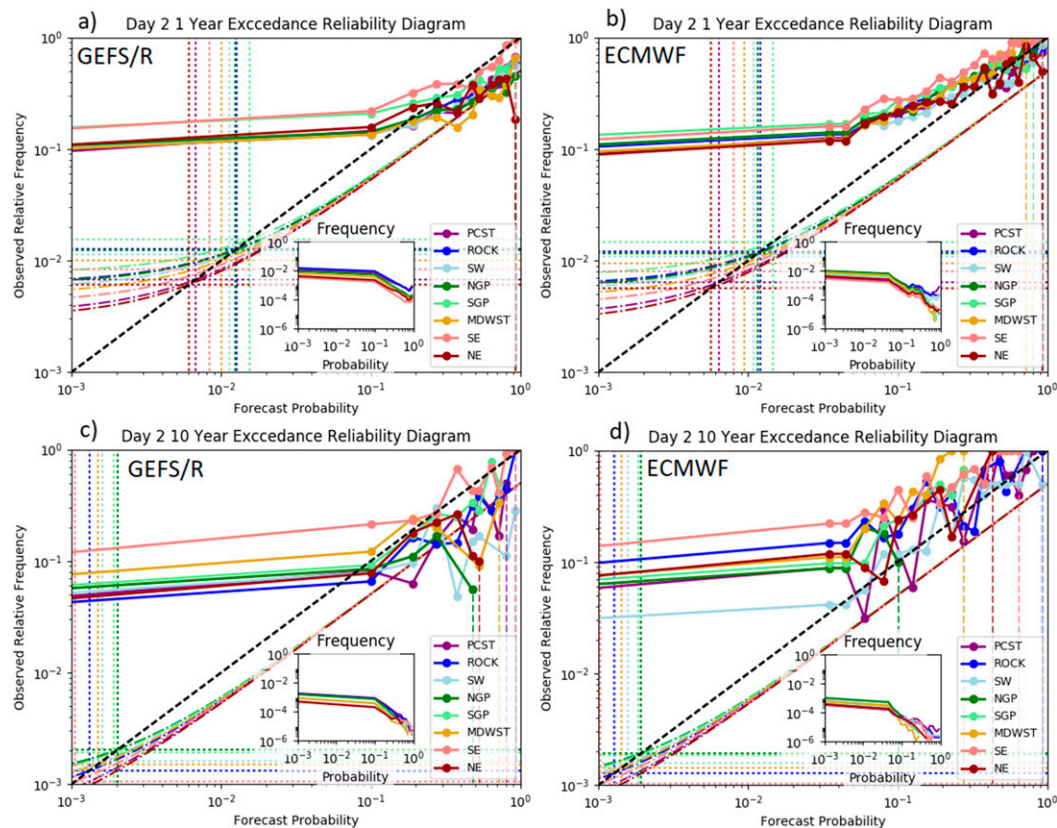


FIG. 7. Reliability diagrams for day 2 forecasts generated from raw QPFs of the full GEFS/R and ECMWF ensembles. Colored opaque lines with circular points indicate observed relative frequency as a function of forecast probability; the dashed black line is the one-to-one line, indicating perfect reliability. Colors correspond to the performance of the forecasts over different regions, as indicated in the legend in the lower-right of each panel. Inset panels indicate the total proportion of forecasts falling in each forecast probability bin, using the logarithmic scale on the left-hand side of each panel; lines are again colored by region in accordance with the legend. The 1-yr exceedance forecast from (a) GEFS/R and (b) ECMWF and the 10-yr exceedance forecast from (c) GEFS/R and (d) ECMWF. All axes are logarithmic as labeled. Colored dotted lines indicate the climatological event probability for each region for the ARI level of the corresponding panel, while the dash-dotted lines indicate no skill lines for the color-corresponding region. The curves continue off the left end of each panel toward the ORF of forecasts in the zero forecast probability bin.

into each probability bin above 10% and a vast majority of zero probability forecasts (not shown). For all regions, there are cases where every ensemble member has simultaneously predicted a 1-yr exceedance (Fig. 7a), but the same is not true for 10-yr exceedance predictions in the northeastern regions: NE, NGP, and MDWST (Fig. 7b). The ECMWF (Figs. 7b,d) is also overconfident, but we see that it is also negatively biased for all cases. Its degree of overconfidence is dampened, compared with the GEFS/R, and it is not as sharp, with fewer occurrences of very high forecast probabilities, except in the westernmost regions of ROCK and PCST (Fig. 7b, inset panels). With 50 members rather than 11, there is also substantially more resolution across the probability spectrum in the ECMWF forecasts. By the very nature of how these forecasts are generated, quite a bit of

sharpness is inherent at the cost of reliability, since it is not possible for probabilities near the climatological event frequency to be issued for either raw ensemble, but particularly for the GEFS/R.

The day 2 reliability diagrams for 1-yr exceedance forecasts from the different components of the final model—CTL_NPCA, CTL_PCA, and CTL_LR—are shown in Fig. 8. The CTL_NPCA (Fig. 8a) shows markedly different characteristics than either of the raw ensembles. In particular, all of the regions exhibit an underconfidence signal, with low probability events below about 2% for 1-yr events (Fig. 8a) occurring with observed relative frequencies below the forecast probabilities. The relative event frequencies are conversely appreciably higher than the forecast probabilities would indicate for probabilities above 5%. Among the regions,

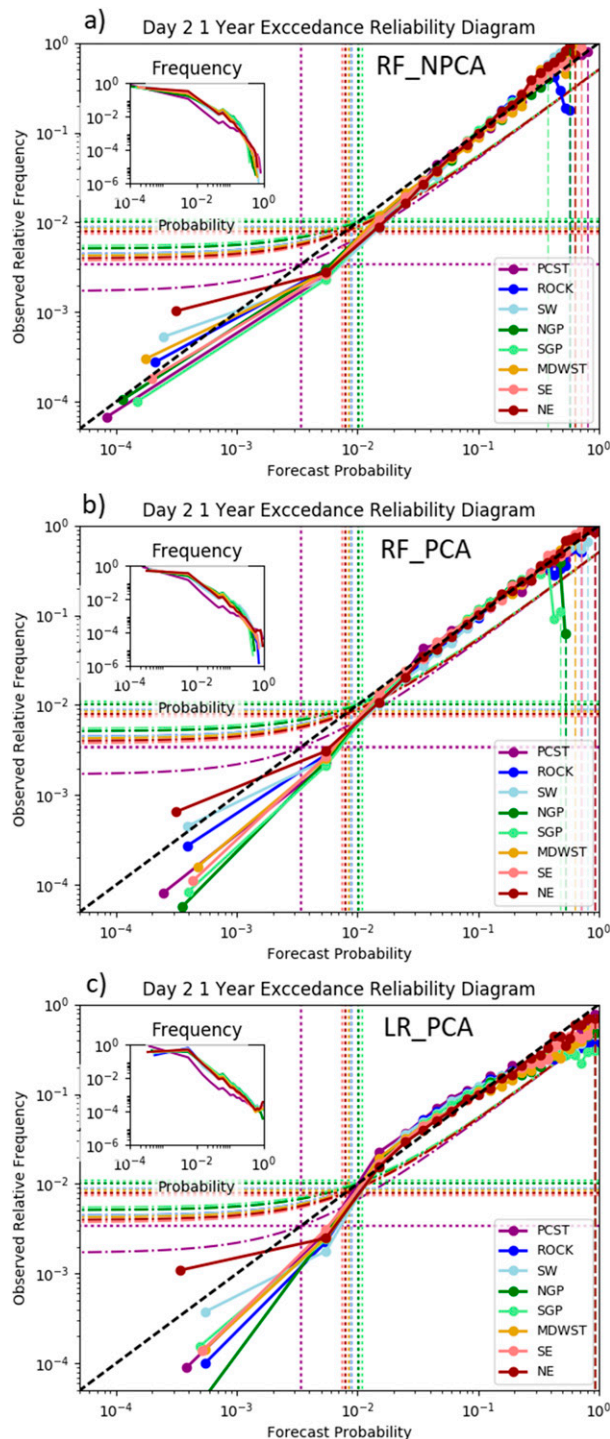


FIG. 8. Reliability diagrams for day 2 forecasts of 1-yr ARI exceedances for different statistical algorithms. Panel characteristics as in Fig. 7, but note that axes have been modified to include more of the low-probability tail due to increased resolution in the plotted forecast sets. Forecasts from the (a) CTL_NPCCA model, (b) CTL_PCA model, and (c) CTL_LR model. Bin right edges correspond to forecast probabilities of $0, 1 \times 10^{-10}, 1 \times 10^{-7}, 1 \times 10^{-4}, 1 \times 10^{-3}, 0.01, 0.02, 0.03, 0.04, 0.05, 0.07, 0.09, 0.11, 0.14, 0.17, 0.21, 0.25, 0.30, 0.35, 0.40, 0.45, 0.50, 0.55, 0.60,$

the PCST probabilities are the most negatively biased, while NE probabilities are the most positively biased. Overall, reliability is much better than for either raw ensemble, but this comes at the expense of sharpness. Fewer than 1 in 10 000 forecasts are above about 20% (Fig. 8a, inset panels), and maximum probabilities are in the 30%–80% range, depending on the lead time and region, compared with 100% for all lead times and regions in the raw ensembles. The CTL_PCA model (Fig. 8b) exhibits very similar reliability characteristics to the CTL_NPCCA model, including the underconfidence, reduced sharpness compared with the raw ensembles, and different regional probability bias characteristics. It tends to be more negatively biased than CTL_NPCCA at low and high probabilities (cf. Figs. 8a,b), correctly so at high probabilities and undesirably so at low ones. The CTL_LR model (Fig. 8c) exhibits some similarities and some differences with the RF-based models. PCST forecasts are consistently the most negatively biased, followed by ROCK and the SE, with NE region forecasts being the least negatively biased. However, unlike the RF-based forecasts, the LR model issues a larger number of high probabilities; for example, forecasts in the highest probability bin were issued for most regions (Fig. 8c). At the highest probabilities, the forecasts revert to being positively biased, as they are for events with probabilities issued in the 0.01%–1% range. At very low probabilities, LR-based forecasts are substantially more negatively biased than for RF-based forecasts, leading to considerable overconfidence overall when considering that the vast majority of forecasts issued occur on this low-probability end of the spectrum. While LR (and regression in general) is effective at removing bias in a global sense, since a single regression equation must necessarily apply globally to all forecasts, it inherently cannot perform more localized, context-dependent forms of bias correction, leading to forecast probability-dependent model biases.

The final ML model reliability (Fig. 9) unsurprisingly reflects a blend of the component members, retaining some of the underconfidence of the RF-based models while adding a bit of sharpness from the CTL_LR model in regions where it verified skillfully enough in cross validation (e.g., PCST; Fig. 5d) to garner much weight. The probability distribution for 1-yr exceedance events

←
0.675, 0.75, 0.85, and 1.0, except that the first five probability bins have been aggregated into a single frequency-weighted probability bin for plotting on the figure.

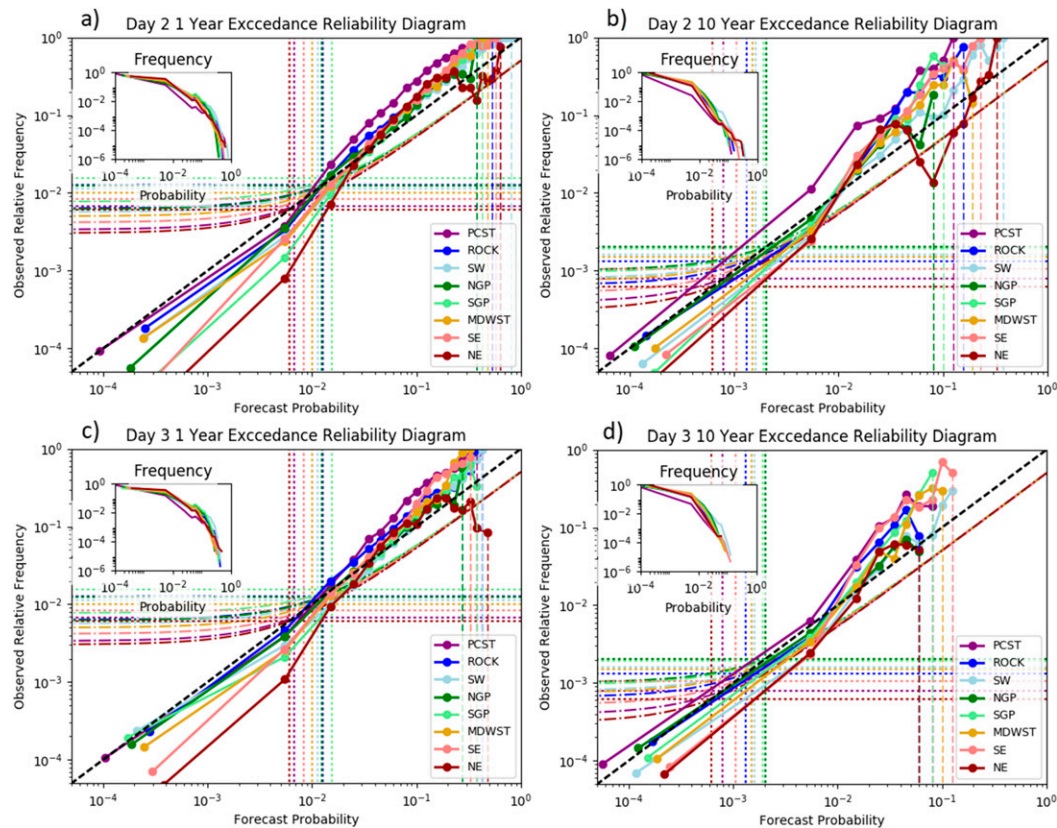


FIG. 9. Reliability diagrams for the final forecast model, with panel attributes as in Fig. 8. Day 2 forecast results for (a) 1- and (b) 10-yr ARI exceedance forecasts and day 3 (c) 1- and (d) 10-yr ARI exceedance forecasts.

is not markedly different between the day 2 and day 3 forecasts (cf. Figs. 9a,c), but the relatively higher probabilities issued for 10-yr exceedances on day 2 do not occur at the day 3 lead times (cf. Figs. 9b,d). This is consistent with increasing confidence in very extreme events with decreasing lead time—something seen very pronounced in the final model, but to a much lesser extent in the raw ensemble forecasts.

The relationship between the reliability analysis and skill via the Brier score decomposition (Murphy 1973) quantitatively solidifies many of the general observations discerned by inspection of the reliability diagrams. Though sharper than competing forecasts, the raw GEFS/R forecasts consistently exhibit the worst resolution component contribution to forecast skill for all regions and severity levels, both for day 2 forecasts (Figs. 10a,c) and day 3 forecasts (Figs. 11a,c) due to an inability to actually distinguish events from nonevents by resolving the responsible physical mechanisms. The final ML models exhibit better resolution term skill contributions than the ECMWF ensemble forecasts, with the exception of the ROCK and NGP regions for 1-yr events (Figs. 10a, 11a). Between the component

models, resolution term skill tended to best for CTL_NPCA forecasts over the test period, particularly at the 10-yr severity level (e.g., Fig. 10c), but the extent of the difference tended to be relatively small, and there were numerous instances where PCA-based models exhibited more resolution. The weighted average consistently exhibited higher resolution than any of the component members. With respect to the reliability contribution to skill (Figs. 10b,d for day 2; Figs. 11b,d for day 3), ECMWF forecasts were—perhaps surprisingly, given the lack of explicit calibration—the most reliable forecast set for all regions and lead times, while in many cases, the ML models had a more negative contribution to the total skill than the raw GEFS/R, likely resulting from the underconfidence. The resolution term is, at largest, 1 and at least 0 in this decomposition, while the reliability term is, at most, 0. The magnitude of the resolution terms is consistently several factors larger than the reliability term for all forecast sets, and the differences in that term generally have a larger absolute impact on the overall Brier skill scores.

While by no means a comprehensive characterization of the system, a sample of real cases over the test period

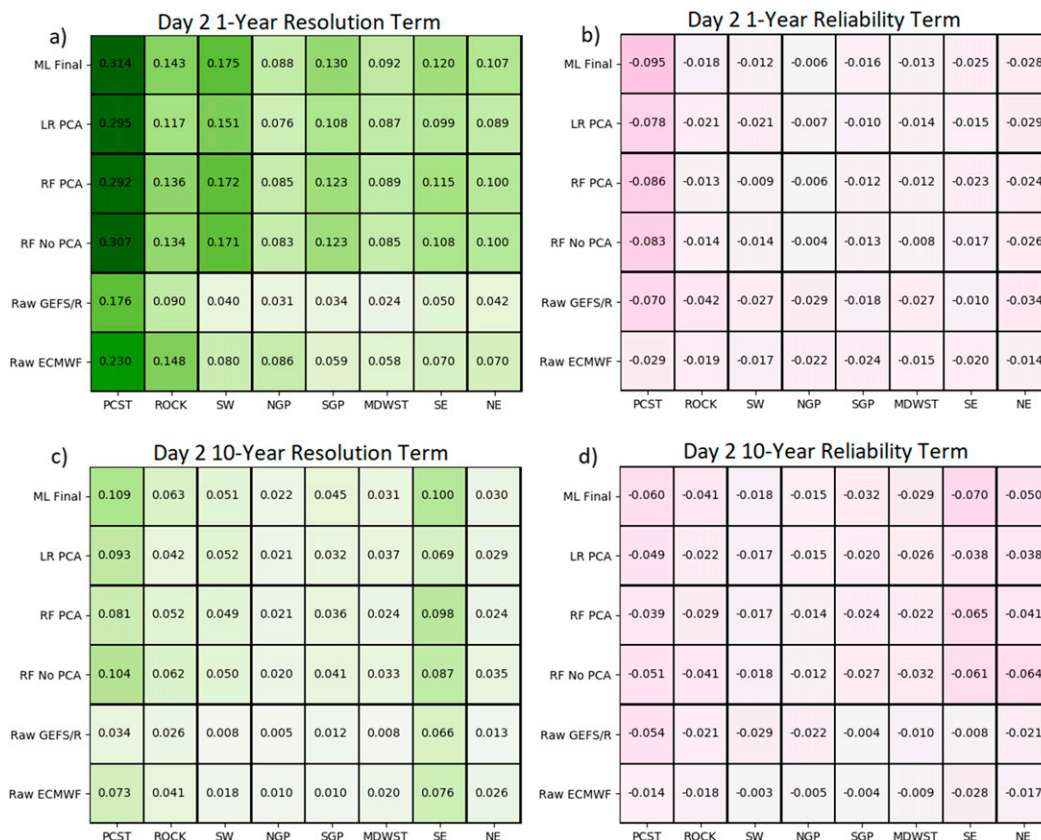


FIG. 10. Modified Murphy (1973) decomposition results, following Eq. (3) in text. (a) Equation (3) resolution term for all models and regions for day 2 forecasts at the 1-yr severity level and (b) the reliability term results for the same forecasts and severity level. (c),(d) As in (a),(b), but for 10-yr ARI exceedance forecasts. Numeric values indicate the value of the corresponding term of the table, as indicated by the model label (row) and region (column).

is presented to illustrate some of the strengths and weaknesses of the system. On the evening of 19 May and the morning of 20 May 2015, a vigorous mesoscale convective system developed over southern Oklahoma and northern Texas, producing very heavy rainfall that contributed to historic flooding in the region during May 2015 (e.g., Wolter et al. 2016). Stage IV analysis (Fig. 12a) reveals that the 24-h precipitation totals exceeded 1-yr ARI thresholds within much of an E–W band encompassing the region, with embedded areas of 10-yr exceedances along the state border region (Fig. 12b). While the ECMWF ensemble forecasts indicate some possibility of extreme precipitation in that region during this time frame at day 3 (Fig. 12d), the probabilities are displaced too far to the south and west, and the probabilities of 10-yr exceedances are very low. There is some improvement in positioning with the day 2 forecast (Fig. 12c), but it remains too far west and with probabilities still quite low, particularly at the 10-yr ARI level. Raw GEFS/R forecasts at day 3 (Fig. 12f) indicate quite high risk for a 1-yr exceedance over a fairly narrow

area, better positioned than the ECMWF ensemble at the same lead time but still too far to the west. Outside of this area, the GEFS/R indicates almost no risk of an extreme rainfall event and also indicates no risk of a 10-yr exceedance anywhere in the domain. The day 2 forecast (Fig. 12e) looks similar to the day 3 outlook, except that the probabilities are reduced somewhat in the target area, which also has incorrectly displaced farther to the south and west. The ML model depicts a much different picture. It exudes much less confidence, with lower maximum probabilities, compared with either raw ensemble, but nonzero exceedance probabilities of both 1- and 10-yr exceedances across much of the domain for both days 3 (Fig. 12h) and 2 (Fig. 12g). Importantly, the model elevated probabilities compared with the raw guidance in the place that extreme precipitation was actually observed (to the east of where it was forecast in the GEFS/R). In fact, at day 2 (Fig. 12g), the probability maximum is located right where the heaviest precipitation actually occurred, displaced well to the north and east of where it was forecast in the

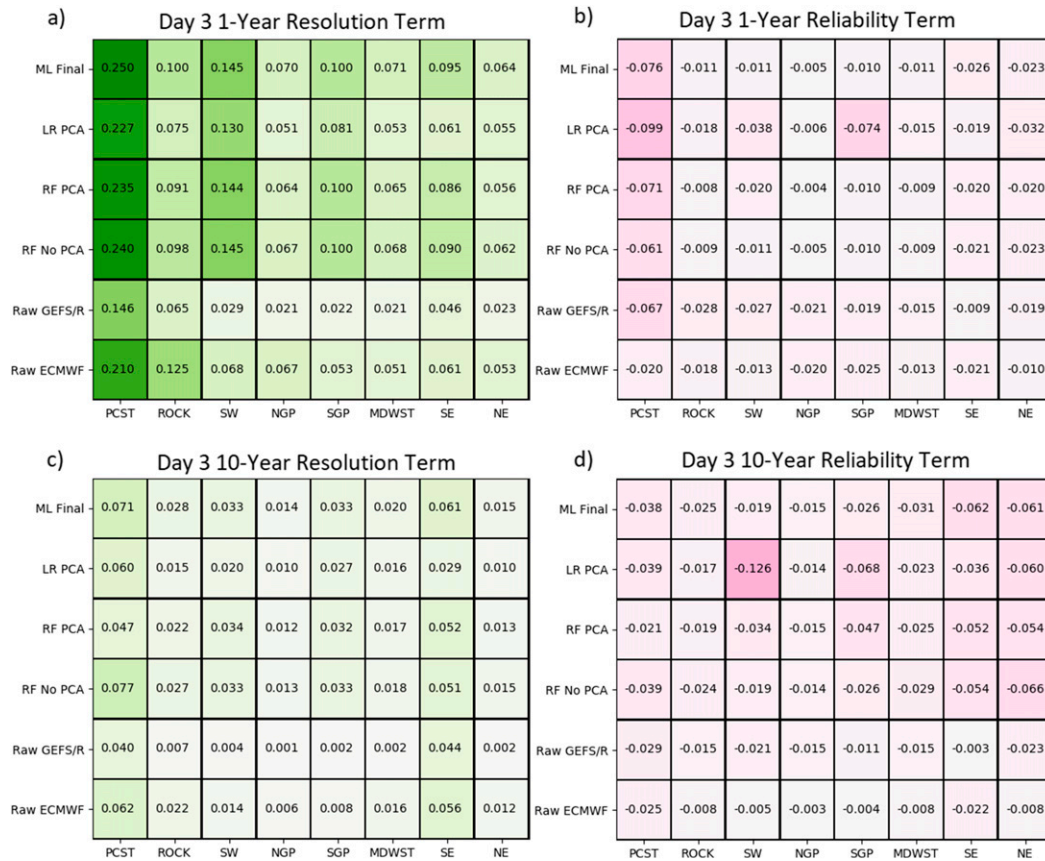


FIG. 11. As in Fig. 10, but for day 3 forecasts.

GEFS/R (Fig. 12e). Additionally, while still low, the 10-yr event probabilities are much higher over the verifying area when compared with either raw ensemble, with maximum day 2 probabilities of around 30% and 3% for 1- and 10-yr exceedances, respectively. Finally, in contrast to the raw guidance, the ML model became increasingly confident in an event occurring with decreasing lead time (cf. Figs. 12g,h).

A different mesoscale precipitation produced extreme precipitation over southwestern Wisconsin, southeastern Minnesota, and northeastern Iowa during the evening and overnight hours of 21 and 22 September 2016, respectively. Based on ST4 QPE (Fig. 13a), much of the area experienced 1-yr ARI exceedances for the 24-h period ending 1200 UTC 22 September 2016, and within the 1-yr exceedance area, there were many embedded cells that produced 10-yr ARI exceedances (Fig. 13b). ECMWF forecasts at day 3 indicated risk of extreme rainfall, even at the 10-yr severity level (Fig. 13d), but the location was poor, with exceedance probabilities high in eastern Minnesota and northern Wisconsin where extreme rainfall was not observed, and very low probabilities in northeastern Iowa and

southeastern Wisconsin where it was. Both the positioning and risk of very extreme precipitation improved for the day 2 forecast issuance (Fig. 13c), but probabilities still remained too far to the north. The GEFS/R at day 3 (Fig. 13f) indicated very little risk of extreme precipitation in the area, with just one member correctly predicting a 1-yr exceedance in southeastern Minnesota. The risk of an event occurring within the domain increased for the day 2 issuance, but the locations got worse, with maximum risk indicated in eastern Nebraska, western Iowa, and northeastern Wisconsin and the only 10-yr prediction occurring in the latter location. Somewhat like the raw GEFS/R, the ML model had only some indication of extreme precipitation risk at day 3 (Fig. 13h). However, it both had the higher probabilities (near 10% in both cases) distributed over a much larger area and indicated some risk of a 10-yr event, with probability maxima near 1.5%. Additionally, it had the maximum probability axis nearly collocated with where heaviest precipitation occurred: well to the south of the ECMWF probabilities, albeit still slightly too far to the north. The day 2 forecast issuance (Fig. 13g) was largely similar. The two main changes are a correctly increased

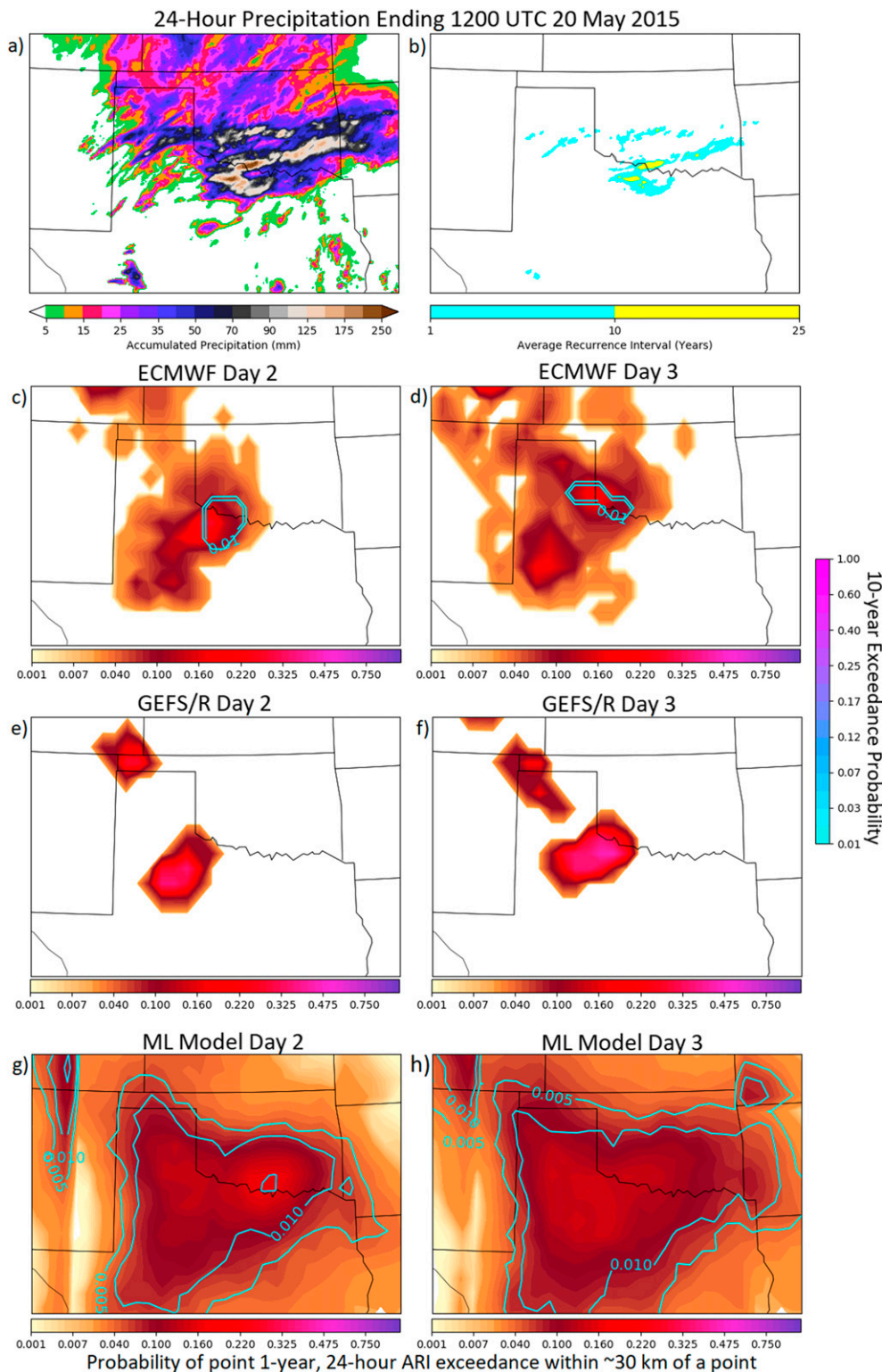


FIG. 12. Case study depicting forecasts from the final ML model and both reference ensembles for the 24-h period ending 1200 UTC 20 May 2015. (a) 24-h stage IV QPE ending at 1200 UTC 20 May 2015 and (b) corresponding ARI exceedances of 1- and 10-yr thresholds. (c) ECMWF ensemble neighborhood ARI exceedance probabilities in the filled (1 yr) and unfilled (10 yr) contours for the 36–60-h forecast initialized 0000 UTC 18 May 2015 and (d) for the 60–84 h-forecast initialized 0000 UTC 17 May 2015. (e),(f) As in (c),(d), but for forecasts from the raw GEFS/R QPFs. (g),(h) As in (c),(d), but for 36–60- and 60–84-h forecasts, except for from the final version of the ML model trained in this study. Contours for 10-yr events are 0.005, 0.01, 0.03, 0.05, 0.075, 0.10, 0.125, 0.15, 0.175, 0.20, 0.25, 0.3, 0.4, 0.5, 0.6, 0.8, and 1.0.

risk in the area where the event actually verified and an incorrectly increased risk of heavy precipitation in eastern Nebraska, where the raw GEFS/R had heavy precipitation on day 2 (Fig. 13e).

5. Discussion and conclusions

An ML model based on RFs and LR is used to generate CONUS-wide probabilistic forecasts for the exceedance of 1- and 10-yr ARI thresholds for 24-h precipitation accumulations during the day 2 and day 3 periods. Approximately 11 years of GEFS/R forecasts, in particular the ensemble median, are used to train these models, and forecasts are made using numerous simulated atmospheric fields (Table 1) varying in both space and time, in addition to a variety of geographic and climatological forecast predictors (Table 2). Separate models are trained for each of the two 24-h periods and for each of eight different regions of the CONUS, as depicted in Fig. 3. A variety of sensitivity experiments are performed, as outlined in Table 3, to ascertain the utility of different aspects of forecast information in predicting locally extreme precipitation. Finally, the final forecast models were evaluated and compared with forecasts based only on the ensemble of raw QPFs from the GEFS/R and ECMWF. The ML models trained in this study demonstrably outperformed the raw GEFS/R forecasts for all regions and forecast lead times (Fig. 6), often more than doubling the forecast skill and adding substantially more than 24-h lead time improvement in forecast skill. With the exception of the PCST and ROCK regions, the same held for comparison of the ML model forecasts with ECMWF ensemble forecasts as well. Both raw ensembles tended to be negatively biased and highly overconfident in predicting extreme QPFs (Fig. 7), particularly at the 10-yr ARI for central CONUS regions; this was reversed in the final ML model forecasts, which were more reliable at higher probabilities but generally underconfident (Fig. 9).

In general, unlike past studies (e.g., Herman and Schumacher 2016b), in most regions, the temporal resolution and extent of spatially displaced predictors from the forecast point considered had little to no impact on forecast skill (Fig. 4), in addition to the use of upper-level information and additional ensemble information (Fig. 5). These results are suggestive of two findings. First, most of the relevant information about predictors displaced spatiotemporally from the forecast point, other atmospheric fields, or other ensemble member information can be derived with at least moderate accuracy using just the information from the ensemble median from a core set of fields collocated and concurrent with

the forecast point and time. Thus, these additional predictors contain only limited independent forecast information, at least for this coarse dynamical model and this underdispersive ensemble configuration. It also suggests that for the most part, the predictive ability is coming primarily through a characterization of the overall environment, which can be reasonably summarized with only a subset of predictors, rather than the simulated spatiotemporal variability and full 3D characterization of the atmospheric evolution in the underlying dynamical model. This finding comes in contrast to similar studies of other forecast problems using the GEFS/R, such as the Herman and Schumacher (2016b) study that investigated the use of the GEFS/R to create ML-based probabilistic forecasts of cloud ceiling and visibility at different airports and found considerable value in the inclusion of spatially displaced predictors. However, there is at least one major exception: none of this really held for the PCST region. Here, more complex models with more predictors notably improved forecast skill. This is perhaps in part because the physical processes associated with extreme precipitation are much better resolved in the GEFS/R in this region, compared with the others, and so the added information adds usable forecast utility beyond simply duplicatively characterizing the atmospheric environment for the forecast. The largest skill difference of the sensitivity experiments came for most regions in changing algorithmic assumptions and processes (Fig. 5d); the simpler linear assumptions of LR tended to degrade forecast skill, compared with the more limited assumptions underlying the RF models.

The results of this study reveal that the application of more sophisticated statistical methods and ML algorithms such as RFs can demonstrably improve forecasts of extreme precipitation and potentially other rare, high-impact weather events in the medium range when compared with the methods and techniques that are most prevalent in forecast operations today. One unique aspect here is the scope of this model; while most past studies that employed these techniques for numerical weather prediction have focused on a small domain or just a sampling of points, the models trained here demonstrate an ability to generate skillful, reliable forecasts year-round for all of the CONUS and a range of lead times. There are many forecast problems that remain to be explored, but the results of this study and others strongly suggest that further development and application of these data-intensive statistical techniques could substantially improve our forecasts over the current state of the art, even compared with using more sophisticated dynamical models. To that end, implementation of this methodology for operational use to assist Weather Prediction Center forecasters with the

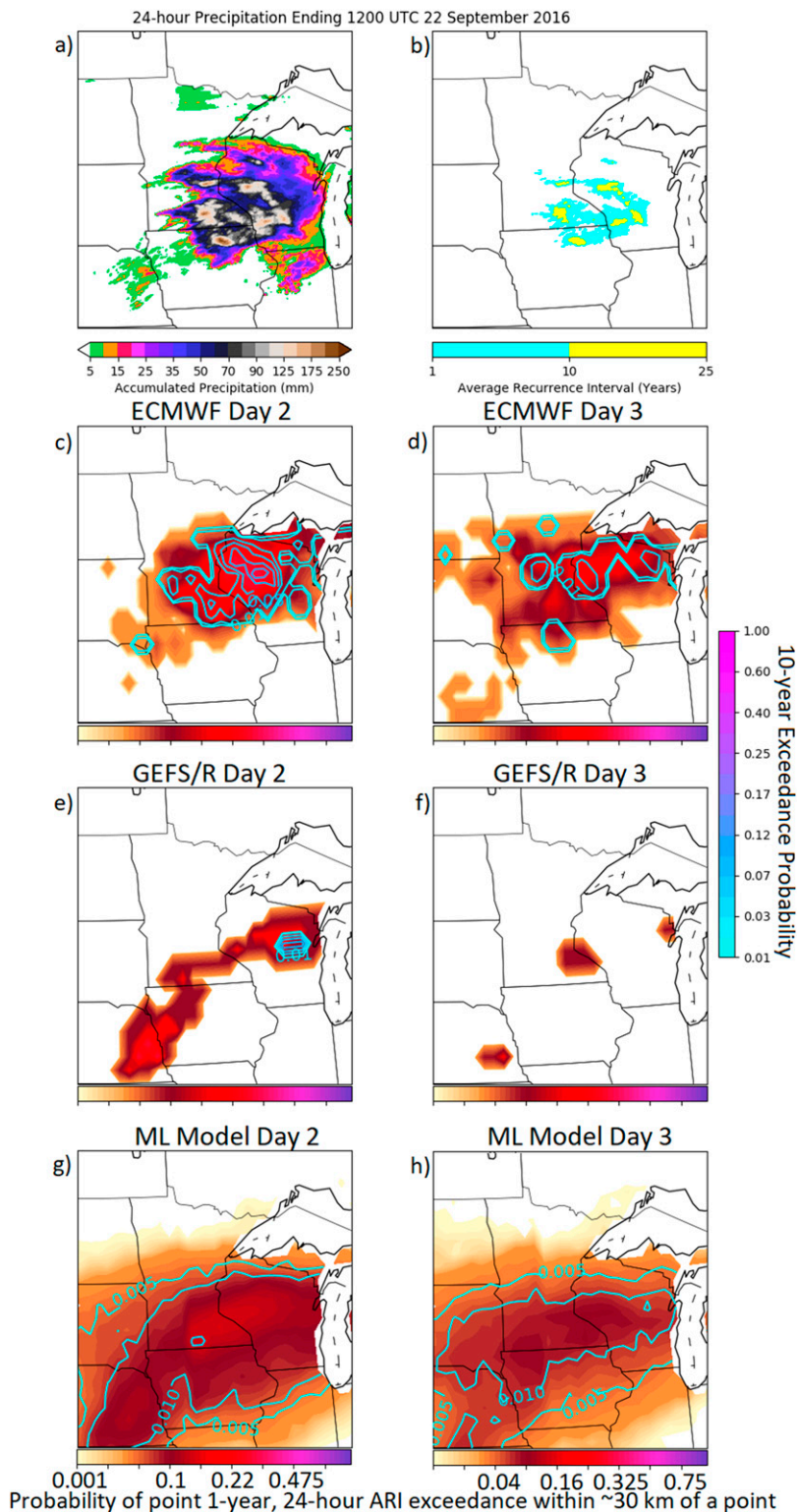


FIG. 13. As in Fig. 12, but for the 24-h period ending 1200 UTC 22 Sep 2016.

development of their excessive rainfall outlooks is currently underway.

This forecast technique presents some advantages over purely dynamical approaches, as dynamical models are inherently limited by two factors by which these statistical techniques are not. First, dynamical models require ever-increasing computational resources for increasing model resolution; constraints on computing power prevent sufficient resolution to directly resolve many small-scale processes, many of which are observed in the highest-impact weather phenomena. Second, dynamical models are limited by our physical understanding of the processes we are attempting to simulate or forecast. Machine learning algorithms, in contrast, can detect predictive patterns in the available information even in places where we do not know or understand the physical connection between the information and the phenomenon that we wish to predict. While they are also limited in complexity by computational and data resources, the strict limits on resolvability are not there: physical resolution can often be gained through postprocessing of larger-scale information. There is thus ample reason to believe that further investigation of these techniques for NWP is a worthwhile venture, and eventual implementation into forecast operations could help forecasters with their tasks by skillfully synthesizing many different sources of forecast information to help alleviate their often time-pressed schedules. This in turn can aid end-user preparedness and, in the case of high-impact events, hopefully help to protect lives and property.

One of the main advantages of the methods explored in this study compared with other popular machine learning methods, in addition to their computational tractability, is the ability to visualize their output and gain insights into detecting and quantifying specific biases in the underlying GEFS/R model and physical insights into the most valuable forecast information for predicting locally extreme precipitation. For focus and brevity, the diagnostics that shed light on these insights have been omitted from this manuscript and are presented instead in a companion paper focused on the diagnostics (Herman and Schumacher 2018) rather than the forecasts and forecast process explored in depth here.

Some limitations of this work are worthy of note. Stage IV precipitation is used as truth for this study; though there is not a clearly better verification source available, it does have its drawbacks. It does have some spurious quality control issues and often struggles in areas of complex terrain due to radar beam blockage, interference, and limited gauge coverage (Herman and

Schumacher 2016a; Nelson et al. 2016). Since the model is trained to forecast stage IV QPE exceedances, this can lead to some idiosyncrasies and other anomalies associated with the biases observed in the stage IV product. One such anomaly is the persistent presence of very small areas of exceedances in some regions of complex terrain during times of favorable convective conditions. This can be removed by quality control procedures to some extent, but some artifacts do remain. This happens most prominently in the terrain of western New Mexico; a small region there has many more instances of ARI exceedances over both the training and test periods than any other part of the CONUS. The ML-based models recognize this and, for the SW region, consistently issue much higher probabilities in this region. In one sense, this is correct—it is correctly predicting what it was trained to predict—but is still undesirable behavior due to a disparity between “truth” in the study and the true extreme rainfall risk. Solutions to this issue and related issues in other parts of the country must be explored in order to maximize operational utility. Additionally, while the choice of using the ARI framework was an intentional decision and provides numerous benefits, it is not an end-all for predicting heavy precipitation impacts. While ARIs often have better correspondence with impacts than a fixed threshold, there are still regional discrepancies in which ARIs have optimal association with impacts, and the framework employed here does not account for antecedent conditions, which can be critical for assessing flash flood risk. More investigation into the relationship between QPE exceedances and rainfall impacts should be performed to maximize the practical significance of the model predictand.

Additionally, the predictors for this study come from a very coarse and otherwise rather antiquated global model. The GEFS/R was used for this study because unlike almost any other dynamical model, it has been nearly static for a very long period of record and has nearly stationary bias characteristics—an essential property for performing this kind of analysis. However, the models trained herein are not working off of the “state of the art” of flash flood predictors. The longer-range day 2 and day 3 lead times were chosen for this study in part because the discrepancy between GEFS/R forecast quality and state of the art is smaller at these longer lead times due to less convection-allowing guidance being available and higher-resolution models degrading in utility with increasing forecast lead time (e.g., Zhang et al. 2003, 2007).

There are also some complications that must be considered for real-time implementation. As one example, the regional models are trained completely

independently of one another, with different training data and different solutions. Consequently, they can occasionally give rather different predictions on nearly identical inputs, resulting in undesirable probability discontinuities across region boundaries. Appropriate methods for removing probability discontinuities in space must be further explored.

Future work will seek to alleviate these limitations in a variety of ways. Exploration of using different predictands, likely combining hydrometeorological information from a variety of sources, will be made for more explicit flash flood prediction. This may involve a regionally varying predictand definition, with some ARI thresholds better corresponding to flash flood impacts in some regions compared with others. Additionally, although a large number of predictors were explored in this study, there are many additional choices for predictors that could ostensibly further improve forecast skill. While atmospheric fields are represented here in absolute terms, it may be beneficial to instead represent some fields relative to the local climatology of the forecast point in terms of standardized anomalies. This is particularly true for fields like PWAT, where standardized anomalies have often shown better correspondence with precipitation impacts across varied regions than absolute values (e.g., Junker et al. 2009; Graham and Grumm 2010; Nielsen et al. 2015). More exploration of derived fields of physical relevance to extreme precipitation processes should also be explored. Some possible examples include upslope flow to gauge forcing for ascent by the horizontal wind, column mean wind to ascertain potential for slow-moving storms, and deep-layer shear as a metric for supercell potential.

This study also focused on a rather specific time interval and took all dynamical predictors from a single, somewhat antiquated ensemble system. Future expansion both to the 12–36-h day 1 period and beyond the day 3 period will be explored, including predictors from more contemporary CAM guidance and potentially including observations as well for the shorter lead time forecasts. Operational models also tend to undergo periodic upgrades and thus do not remain static like the ensemble system used here. The sensitivity of ML model performance to changes in dynamical model bias characteristics that result from these upgrades is a question of considerable operational relevance and an additional factor worthy of future investigation. It was also seen that the ML models suffered to varying degrees from underconfidence and, in some instances, negative bias. Methods of probability calibration of the ML model probabilities as a final post-processing step (e.g., Hagedorn et al. 2008; Hamill et al. 2008; Bentzien and Friederichs 2012; Herman and

Schumacher 2016b) should be explored in future work and parameter choices reconsidered in light of this additional calibration. Finally, this study only explored a subset of available machine learning algorithms. Other choices, including adaptive learning algorithms, may be able to better exploit predictor–predictand relationships, appropriately update to reflect changes in an underlying dynamical model, and produce superior forecasts for the locally extreme precipitation and flash flood forecast problem (e.g., Liu et al. 2001; Roebber 2015; Pelosi et al. 2017).

Acknowledgments. The authors wish to thank Josh Hacker, David John Gagne, and two anonymous reviewers for insights and feedback that greatly improved the quality of the study. The authors would also like to thank Tom Hamill and Gary Bates for generating and providing the GEFS/R data that made this research possible. Erik Nielsen provided helpful assistance in the creation of Fig. 3. We also wish to thank Diana Stovern, Sarah Perfater, Benjamin Albright, Mark Klein, Michael Erickson, and James Nelson at the Weather Prediction Center, who helped improve the operational utility of this research. Funding for this research was supported by NOAA Award NA16OAR4590238 and NSF Grant ACI-1450089.

APPENDIX A

Algorithm Descriptions

a. Random forests

As noted in the main text, RFs are simply an ensemble of decision trees. Decision trees consist of a network of two types of nodes: decision nodes and leaf nodes. Decision nodes each have exactly two children, which may be either decision nodes or leaf nodes, with a binary split based on the numeric value of a single input predictor determining whether to traverse to the left or right child. A leaf node has no children and instead makes a categorical prediction of the outcome of the input example based on the leaf's relationship to its ancestor nodes. For a given forecast, one begins at a decision tree's root, traversing through its children based on the relative value of the forecast's predictors to each decision node's threshold critical value for the predictor associated with the node. This process is repeated until a leaf node is reached; its value corresponding to the leaf becomes the tree's deterministic prediction.

Decision trees can be a powerful approach for a wide array of applications, but they also have several significant drawbacks. In particular, they are very prone to

overfitting (e.g., [Brodley and Utgoff 1995](#)), fitting to the noise of the training data rather than just the underlying relationships. They also do not convey any information about forecast uncertainty, as would be the case in a probabilistic framework. RFs are used instead to alleviate these concerns by producing a probabilistic forecast in a way that can significantly decrease error from overfitting the supplied training error with only a slight increase to error from oversimplistic model assumptions, provided the trees are sufficiently uncorrelated. The difficulty then revolves around generating a large set (forest) of skillful decision trees that are not strongly correlated. The decision tree generating procedure described above is deterministic: a given set of training data will always produce the same decision tree. A forest of identical decision trees, of course, adds no value over using a single decision tree. Two additional processes—tree bagging and feature bagging—are employed to produce unique trees. Tree bagging produces unique trees through a straightforward bootstrapping procedure. Specifically, a forest of size B is formed from the n training examples by creating B samples of size n , with replacement, from the original training data and running the decision tree algorithm on each sample. Overfitting due to correlated trees can still occur under this approach, particularly if a small subset of the original features contains much more robust predictors of the verifying category than the rest ([Breiman 2001](#); [Murphy 2012](#)). To overcome this problem, feature bagging is also employed, whereby only a random subset of the m original input predictors is considered at each decision node. The size of the random subset is denoted here as S : $1 \leq S \leq m$. This combination can result in a set of B largely uncorrelated trees, each of which is individually fairly skillful.

With any machine learning algorithm, there are numerous considerations in the actual model construction that manifest themselves in tunable parameters. Compared with other machine learning algorithms, such as gradient boosting or support vector machines, RFs are often praised for their relative insensitivity to their parameters with respect to model performance, but it is nevertheless important to explore the parameter space in order to realize the full utility of the algorithm. The forest size B is perhaps the most obvious parameter. The general relationship between model performance and B is well known and consistent across all prediction problems; it starts quite low at very low B , initially increases rapidly with increasing B , and then slowly asymptotes to some threshold performance limit as the relationships between input features have been fully explored by the forest and the inclusion of new trees becomes redundant. Larger forest sizes require more

computational expense, so the goal is to select B such that it is small enough to be computationally tractable but large enough to be near the performance limit. Another parameter noted above is S , the number of features to consider at each node split. If this number is too small, model performance may suffer from only considering irrelevant or otherwise uninformative features in the context of the node; if S is too large, performance will also suffer because of underdispersive trees producing an overfit forest solution. Another frequently explored parameter is the splitting criterion evaluation function. Most commonly used is either the Gini impurity or the information gain; past studies have shown that this choice is not important for many forecast problems. Information gain is used in this study; it can be expressed for a training set T , candidate splitting feature x_a , and candidate split value v_a as

$$\text{IG}(T, x_a, v_a) = H(T) - H(T|x_a < v_a), \quad (\text{A1})$$

where $H(T)$ is the so-called entropy of a tree, defined for each of the K verifying categories, with each category i having forecast probability p_i , as

$$H(T) = - \sum_{i=1}^K p_i \log_2 p_i. \quad (\text{A2})$$

The chosen splitting feature and split value are selected among those considered that maximize Eq. (A1) (e.g., [Quinlan 1986](#); [Murphy 2012](#)). However, there are two other parameters that have the most substantial influence on model performance. The first, denoted Z , is the minimum number of training examples required to split a node. Traditionally, RFs create a leaf only once a node is “pure”; that is, all the remaining training examples associated with that node have the same labels (event outcomes). In this way, each tree makes a categorical prediction of the predictand outcome, and probabilities are generated only in counting the proportion of trees in the forest making a particular forecast. However, this can make predictions from an individual tree very susceptible to the outcome of a particular historical case and, in some cases, result in substantial overfitting. Instead, by increasing Z , an RF can be allowed to make “impure” leaves; at these nodes, an individual tree makes a probabilistic prediction based on the proportion of remaining training examples exhibiting each event class rather than continuing to split based on the remaining training data. Making S too large, however, can result in underfitting—lumping data as indistinguishable when there are, in fact, underlying discernible distinctions among remaining training examples with different labels. The last parameter,

denoted P , is not actually an RF algorithm parameter at all. When PCA is performed, there is always a question about the number of components to retain. Though there are some heuristics (e.g., North et al. 1982), there is no definitive method to know a priori how many retained components P will produce the most skillful forecasts (Wilks 2011). If P is too small, valuable forecast data are discarded, and predictive performance consequently suffers. However, if it is too large, the retained PCs eventually become essentially just noise, and the RF, by fitting to these predictor values in the training data, will yield an overfit model that does not generalize to unseen data. Experiments that will not be discussed herein revealed that using information gain to determine splits and letting $B = 1000$ produced skill near that of an infinitely large forest, and skill was insensitive to modifications of these settings, including modest increases in the forest size beyond this point. However, the Z - S - P parameter space is explored for the models trained, and those results are presented in appendix B.

One final consideration concerns the handling of rare event scenarios. For rare event problems, one necessarily has many more examples of the common event class in comparison to the rare class, leaving the rare class somewhat underrepresented in the learning problem, and model fitting that is done with respect to the rare class is often too dependent on a small number of examples. An approach that has been applied with some success in past studies (e.g., Ahijevych et al. 2016) is to sample training data disproportionately from the rarer classes so that the number of training examples associated with each event class is approximately equal. A comparison between this so-called “balanced” sampling and unmodified “unbalanced” sampling is also made and the results presented in appendix B.

b. Logistic regression

One sensitivity experiment compares model performance as a function of the model algorithm by comparing skill of forecasts produced by RFs with those produced with logistic regression. LR is in many senses a simpler model than an RF, since the structural form of the relationship between the predictors and the predictand is predefined before training. RFs, in contrast, make few assumptions about the relationships between the predictors and the predictand, allowing more diverse diagnoses of underlying relationships. However, this lack of assumptions can result in overfitting. As an application of the generalized linear model, LR assumes a linear predictors–predictand relationship via the logit function. In LR, a single regression equation, or K equations for a multicategory problem with K categories, is computed to represent the probability of the

outcome being category k , given the set of input predictors \mathbf{x} . In particular, verifying probabilities are computed using the softmax function:

$$P(y = k|\mathbf{x}) = \frac{e^{\mathbf{x}^T \mathbf{w}_k}}{\sum_{j=1}^K e^{\mathbf{x}^T \mathbf{w}_j}}. \quad (\text{A3})$$

In training an LR model, the goal is to determine the optimal weights \mathbf{w}_k associated with each predictor in order to yield the most accurate predictions for each event class. As with RF models, LR can be prone to overfitting if unconstrained. For RFs, one aforementioned approach to alleviate this problem is to increase the above-termed Z parameter, which stops node splitting earlier on and makes the model less tailored to the specific training data supplied to it. Complexity in LR can be thought of as being analogously represented by large weights, or regression coefficients. To ensure better generalizability of the trained regression equations, it is often good practice to penalize large weights through a process known as regularization. When this is done, the computation of optimal weights can be represented as a minimization problem with two terms. For 1) a matrix \mathbf{Y} with binary elements that are nonzero if and only if training example i has associated verifying category k and 2) a model outputting a probability matrix \mathbf{P} for each training example and category, the multinomial loss J to be minimized can be computed as

$$J[\mathbf{Y}, \mathbf{P}(\mathbf{w})] = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \frac{1}{CN} \sum_{i=1}^N \sum_{k=1}^K \mathbf{Y}_{i,k} \log(\mathbf{P}_{i,k}), \quad (\text{A4})$$

where C represents the extent of regularization, with smaller values indicating that large weights are penalized more than with larger values of C . Alternative approaches to regularization exist (e.g., Pedregosa et al. 2011; Murphy 2012) and are explored to some degree in sensitivity experiments of appendix B.

c. Computational considerations

Other machine learning algorithms do not scale well to the high dimensionality of the forecast problem explored here. While time to train a model is not of primary concern for operational forecasting since it is performed only once (or periodically) offline, there are nevertheless some practical considerations; for example, models that take months or longer to train would be unlikely to be realistic choices. The “online” forecasting component—that is, the time required to take a new forecast, input it into a trained model, and receive a forecast—is of operational concern, but all of the

TABLE B1. Optimal RF parameters obtained in cross validation for the Z – S – P parameter space. SQRT indicates the square root of the total number of predictors; symbols are otherwise as described in the text. Evaluated values were 1, 2, 4, 8, 16, 30, 60, 120, 240, and 480 for Z and 20, 25, 30, 40, 50, 60, 70, 80, 90, and 100 for P .

Region	S parameter	Z parameter	P parameter
ROCK	SQRT	30	30
NGP	SQRT	120	40
MDWST	SQRT	240	30
NE	SQRT	120	30
PCST	SQRT	4	60
SW	SQRT	30	30
SGP	SQRT	120	30
SE	SQRT	240	30

forecast techniques considered here can produce forecasts in a matter of minutes, and the small differences are not considered to be of practical concern. Using the random forest classification heuristic of considering the square root of the total number of features at each node split (Geurts et al. 2006), the computational complexity of training an RF of size B from N training examples with F features ($N > F$) may be expressed as $O[B\sqrt{F}N \log(N)]$ and may be readily parallelized across trees or within trees. Some algorithms are quadratic or even cubic (e.g., Cortes and Vapnik 1995) in the number of training examples and do not parallelize as readily. LR is linear in the number of training examples but requires matrix multiplication, a process that yields a computational complexity of $O(NF^2)$. PCA preprocessing, and dimensionality reduction more generally, acts both to make learning algorithms more computationally tractable and also to reduce overfitting by alleviating the so-called “curse of dimensionality.”

APPENDIX B

Results: Parameter Tuning

RF model parameters were tuned for each region and lead time separately through the fourfold cross-validation

procedure employed throughout the study. Overall, the optimal parameters were found not to vary with the two different lead times, but did vary for two of the parameters as a function of forecast region, at least to an extent; the full results appear in Table B1. For the S parameter—the number of predictors considered for each node split—the default heuristic of the square root of the total number of features was found to maximize RPSS for all regions and lead times. In all instances where both were tested, unbalanced sampling from the event classes in proportion to their true observed frequencies outperformed balanced equal sampling from each event class, in contrast to Ahijevych et al. (2016) and others; the finding appeared to be attributable to biased probabilities produced from the balanced sampling technique. For the Z parameter, the minimum number of remaining training examples in an impure parameter subspace required to perform a further node split was generally found to be around 120. Lesser values maximized skill in the western regions, with values of 30 maximizing skill in the SW and ROCK regions and $Z = 4$ producing the best skill over PCST. A couple of the larger regions of the east, SE and MDWST, maximized RPSS with a value of 240, although the sensitivity between $Z = 120$ and 240 was small for all regions. For P in the CTL_PCA models, skill was generally maximized with $P = 30$, that is, retaining the 30 PCs that explain the most variance of the entire GEFS/R predictor set. For most regions, there was very limited sensitivity in the $P = 30$ – 40 interval—although there was larger sensitivity outside this interval—and $P = 40$ was found to produce slightly better skill in the NGP region. The PCST region was again the main exception, where $P = 60$ was found to maximize cross-validation RPSS.

LR model parameters were tuned using an identical framework to ascertain the type of regularization, either based on a L1 norm, which penalizes nonzero weights, or L2 norm—described in appendix A—which penalizes large magnitude weights. L2 regularization was consistently found to produce superior results (Table B2), perhaps because the number of

TABLE B2. Optimal LR parameters obtained in cross validation for the C parameter and regularization type for all lead times and regions. Evaluated for C were 0.0001, 0.0008, 0.0060, 0.0464, 0.359, 2.78, 21.54, 167.8, 1291, and 10 000.

Region	Regularization	C parameter, day 2	C parameter, day 3
ROCK	L2	0.0001	0.0001
NGP	L2	10 000	0.0008
MDWST	L2	0.0001	0.0464
NE	L2	2.78	10 000
PCST	L2	0.0001	0.0001
SW	L2	0.359	0.0001
SGP	L2	0.0008	0.0464
SE	L2	21.54	0.0008

retained PCs was already taken from the P parameter in the RF experiments, acting to nullify many potential nonzero weights of higher-numbered PCs. Unlike the RF experiments, there were occasionally some large differences in the obtained optimal regularization parameter value C between lead times within the same region. Generally, models performed better with more regularized solutions, but there were some notable exceptions, with the day 2 NGP model and day 3 NE model obtaining optimal C parameter values on the other end of the spectrum.

REFERENCES

- Ahijevych, D., J. O. Pinto, J. K. Williams, and M. Steiner, 2016: Probabilistic forecasts of mesoscale convective system initiation using the random forest data mining technique. *Wea. Forecasting*, **31**, 581–599, <https://doi.org/10.1175/WAF-D-15-0113.1>.
- Alvarez, F. M., 2014: Statistical calibration of extended-range probabilistic tornado forecasts with a reforecast dataset. Ph.D. thesis, Saint Louis University, 210 pp.
- Antolik, M. S., 2000: An overview of the National Weather Service's centralized statistical quantitative precipitation forecasts. *J. Hydrol.*, **239**, 306–337, [https://doi.org/10.1016/S0022-1694\(00\)00361-9](https://doi.org/10.1016/S0022-1694(00)00361-9).
- Applequist, S., G. E. Gahrs, R. L. Pfeffer, and X.-F. Niu, 2002: Comparison of methodologies for probabilistic quantitative precipitation forecasting. *Wea. Forecasting*, **17**, 783–799, [https://doi.org/10.1175/1520-0434\(2002\)017<0783:COMFPQ>2.0.CO;2](https://doi.org/10.1175/1520-0434(2002)017<0783:COMFPQ>2.0.CO;2).
- Baars, J. A., and C. F. Mass, 2005: Performance of National Weather Service forecasts compared to operational, consensus, and weighted model output statistics. *Wea. Forecasting*, **20**, 1034–1047, <https://doi.org/10.1175/WAF896.1>.
- Barthold, F. E., T. E. Workoff, B. A. Cosgrove, J. J. Gourley, D. R. Novak, and K. M. Mahoney, 2015: Improving flash flood forecasts: The HMT-WPC Flash Flood and Intense Rainfall Experiment. *Bull. Amer. Meteor. Soc.*, **96**, 1859–1866, <https://doi.org/10.1175/BAMS-D-14-00201.1>.
- Bentzien, S., and P. Friederichs, 2012: Generating and calibrating probabilistic quantitative precipitation forecasts from the high-resolution NWP model COSMO-DE. *Wea. Forecasting*, **27**, 988–1002, <https://doi.org/10.1175/WAF-D-11-00101.1>.
- Bermowitz, R. J., 1975: An application of model output statistics to forecasting quantitative precipitation. *Mon. Wea. Rev.*, **103**, 149–153, [https://doi.org/10.1175/1520-0493\(1975\)103<0149:AAOMOS>2.0.CO;2](https://doi.org/10.1175/1520-0493(1975)103<0149:AAOMOS>2.0.CO;2).
- Bonnin, G. M., D. Todd, B. Lin, T. Parzybok, M. Yekta, and D. Riley, 2004: *Precipitation-Frequency Atlas of the United States*. NOAA Atlas 14, Vol. 2, U.S. Department of Commerce, NOAA/NWS, 71 pp.
- , D. Martin, B. Lin, T. Parzybok, M. Yekta, and D. Riley, 2006: *Precipitation-Frequency Atlas of the United States*. NOAA Atlas 14, Vol. 3, U.S. Department of Commerce, NOAA/NWS, 143 pp.
- Bougeault, P., and Coauthors, 2010: The THORPEX Interactive Grand Global Ensemble. *Bull. Amer. Meteor. Soc.*, **91**, 1059–1072, <https://doi.org/10.1175/2010BAMS2853.1>.
- Breiman, L., 2001: Random forests. *Mach. Learn.*, **45**, 5–32, <https://doi.org/10.1023/A:1010933404324>.
- Brimelow, J. C., G. W. Reuter, R. Goodson, and T. W. Krauss, 2006: Spatial forecasts of maximum hail size using prognostic model soundings and HAILCAST. *Wea. Forecasting*, **21**, 206–219, <https://doi.org/10.1175/WAF915.1>.
- Brodley, C. E., and P. E. Utgoff, 1995: Multivariate decision trees. *Mach. Learn.*, **19**, 45–77, <https://doi.org/10.1023/A:1022607123649>.
- Brooks, H. E., and D. J. Stensrud, 2000: Climatology of heavy rain events in the United States from hourly precipitation observations. *Mon. Wea. Rev.*, **128**, 1194–1201, [https://doi.org/10.1175/1520-0493\(2000\)128<1194:COHREI>2.0.CO;2](https://doi.org/10.1175/1520-0493(2000)128<1194:COHREI>2.0.CO;2).
- Buizza, R., A. Hollingsworth, F. Lalauette, and A. Ghelli, 1999: Probabilistic predictions of precipitation using the ECMWF Ensemble Prediction System. *Wea. Forecasting*, **14**, 168–189, [https://doi.org/10.1175/1520-0434\(1999\)014<0168:PPOPUT>2.0.CO;2](https://doi.org/10.1175/1520-0434(1999)014<0168:PPOPUT>2.0.CO;2).
- Clark, A. J., W. A. Gallus Jr., and T.-C. Chen, 2007: Comparison of the diurnal precipitation cycle in convection-resolving and non-convection-resolving mesoscale models. *Mon. Wea. Rev.*, **135**, 3456–3473, <https://doi.org/10.1175/MWR3467.1>.
- , —, and M. L. Weisman, 2010: Neighborhood-based verification of precipitation forecasts from convection-allowing NCAR WRF Model simulations and the operational NAM. *Wea. Forecasting*, **25**, 1495–1509, <https://doi.org/10.1175/2010WAF2222404.1>.
- Cortes, C., and V. Vapnik, 1995: Support-vector networks. *Mach. Learn.*, **20**, 273–297, <https://doi.org/10.1007/BF00994018>.
- Davis, C. A., K. W. Manning, R. E. Carbone, S. B. Trier, and J. D. Tuttle, 2003: Coherence of warm-season continental rainfall in numerical weather prediction models. *Mon. Wea. Rev.*, **131**, 2667–2679, [https://doi.org/10.1175/1520-0493\(2003\)131<2667:COWCRI>2.0.CO;2](https://doi.org/10.1175/1520-0493(2003)131<2667:COWCRI>2.0.CO;2).
- Delrieu, G., and Coauthors, 2005: The catastrophic flash-flood event of 8–9 September 2002 in the Gard Region, France: A first case study for the Cévennes–Vivarais Mediterranean Hydrometeorological Observatory. *J. Hydrometeorol.*, **6**, 34–52, <https://doi.org/10.1175/JHM-400.1>.
- Duda, J. D., and W. A. Gallus, 2013: The impact of large-scale forcing on skill of simulated convective initiation and upscale evolution with convection-allowing grid spacings in the WRF. *Wea. Forecasting*, **28**, 994–1018, <https://doi.org/10.1175/WAF-D-13-00005.1>.
- Eckel, F. A., 2003: Effective mesoscale, short-range ensemble forecasting. Ph.D. thesis, University of Washington, 64 pp.
- Friedman, J. H., 1997: On bias, variance, 0/1—loss, and the curse-of-dimensionality. *Data Min. Knowl. Discovery*, **1**, 55–77, <https://doi.org/10.1023/A:1009778005914>.
- Fritsch, J. M., and R. Carbone, 2004: Improving quantitative precipitation forecasts in the warm season: A USWRP research and development strategy. *Bull. Amer. Meteor. Soc.*, **85**, 955–966, <https://doi.org/10.1175/BAMS-85-7-955>.
- Gagne, D. J., A. McGovern, and M. Xue, 2014: Machine learning enhancement of storm-scale ensemble probabilistic quantitative precipitation forecasts. *Wea. Forecasting*, **29**, 1024–1043, <https://doi.org/10.1175/WAF-D-13-00108.1>.
- , —, J. Brotzge, M. Coniglio, J. Correia Jr., and M. Xue, 2015: Day-ahead hail prediction integrating machine learning with storm-scale numerical weather models. *Proc. 27th Conf. on Innovative Applications of Artificial Intelligence*, Austin, TX, AAAI, 3954–3960, <https://www.aaai.org/ocs/index.php/IAAI/IAAI15/paper/viewFile/9724/9898>.
- , —, S. E. Haupt, R. A. Sobash, J. K. Williams, and M. Xue, 2017: Storm-based probabilistic hail forecasting with machine learning applied to convection-allowing ensembles. *Wea. Forecasting*, **32**, 1819–1840, <https://doi.org/10.1175/WAF-D-17-0010.1>.

- Geurts, P., D. Ernst, and L. Wehenkel, 2006: Extremely randomized trees. *Mach. Learn.*, **63**, 3–42, <https://doi.org/10.1007/s10994-006-6226-1>.
- Glahn, H. R., and D. A. Lowry, 1972: The use of model output statistics (MOS) in objective weather forecasting. *J. Appl. Meteor.*, **11**, 1203–1211, [https://doi.org/10.1175/1520-0450\(1972\)011<1203:TUOMOS>2.0.CO;2](https://doi.org/10.1175/1520-0450(1972)011<1203:TUOMOS>2.0.CO;2).
- Gochis, D., and Coauthors, 2015: The great Colorado flood of September 2013. *Bull. Amer. Meteor. Soc.*, **96**, 1461–1487, <https://doi.org/10.1175/BAMS-D-13-00241.1>.
- Gourley, J. J., J. M. Erlingis, Y. Hong, and E. B. Wells, 2012: Evaluation of tools used for monitoring and forecasting flash floods in the United States. *Wea. Forecasting*, **27**, 158–173, <https://doi.org/10.1175/WAF-D-10-05043.1>.
- Graham, R. A., and R. H. Grumm, 2010: Utilizing normalized anomalies to assess synoptic-scale weather events in the western United States. *Wea. Forecasting*, **25**, 428–445, <https://doi.org/10.1175/2009WAF2222273.1>.
- Grams, J. S., W. A. Gallus Jr., S. E. Koch, L. S. Wharton, A. Loughe, and E. E. Ebert, 2006: The use of a modified Ebert–McBride technique to evaluate mesoscale model QPF as a function of convective system morphology during IHOP 2002. *Wea. Forecasting*, **21**, 288–306, <https://doi.org/10.1175/WAF918.1>.
- Hagedorn, R., T. M. Hamill, and J. S. Whitaker, 2008: Probabilistic forecast calibration using ECMWF and GFS ensemble reforecasts. Part I: Two-meter temperatures. *Mon. Wea. Rev.*, **136**, 2608–2619, <https://doi.org/10.1175/2007MWR2410.1>.
- Hall, T., H. E. Brooks, and C. A. Doswell III, 1999: Precipitation forecasting using a neural network. *Wea. Forecasting*, **14**, 338–345, [https://doi.org/10.1175/1520-0434\(1999\)014<0338:PFUANN>2.0.CO;2](https://doi.org/10.1175/1520-0434(1999)014<0338:PFUANN>2.0.CO;2).
- Hamill, T. M., 2017: Changes in the systematic errors of global reforecasts due to an evolving data assimilation system. *Mon. Wea. Rev.*, **145**, 2479–2485, <https://doi.org/10.1175/MWR-D-17-0067.1>.
- , and J. Juras, 2006: Measuring forecast skill: Is it real skill or is it the varying climatology? *Quart. J. Roy. Meteor. Soc.*, **132**, 2905–2923, <https://doi.org/10.1256/qj.06.25>.
- , and J. S. Whitaker, 2006: Probabilistic quantitative precipitation forecasts based on reforecast analogs: Theory and application. *Mon. Wea. Rev.*, **134**, 3209–3229, <https://doi.org/10.1175/MWR3237.1>.
- , —, and X. Wei, 2004: Ensemble reforecasting: Improving medium-range forecast skill using retrospective forecasts. *Mon. Wea. Rev.*, **132**, 1434–1447, [https://doi.org/10.1175/1520-0493\(2004\)132<1434:ERIMFS>2.0.CO;2](https://doi.org/10.1175/1520-0493(2004)132<1434:ERIMFS>2.0.CO;2).
- , 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, <https://doi.org/10.1175/2007MWR2411.1>.
- , G. T. Bates, J. S. Whitaker, D. R. Murray, M. Fiorino, T. J. Galarneau Jr., Y. Zhu, and W. Lapenta, 2013: NOAA’s second-generation global medium-range ensemble reforecast dataset. *Bull. Amer. Meteor. Soc.*, **94**, 1553–1565, <https://doi.org/10.1175/BAMS-D-12-00014.1>.
- , M. Scheuerer, and G. T. Bates, 2015: Analog probabilistic precipitation forecasts using GEFS reforecasts and climatology-calibrated precipitation analyses. *Mon. Wea. Rev.*, **143**, 3300–3309, <https://doi.org/10.1175/MWR-D-15-0004.1>.
- Hapuarachchi, H., Q. Wang, and T. Pagano, 2011: A review of advances in flash flood forecasting. *Hydrol. Processes*, **25**, 2771–2784, <https://doi.org/10.1002/hyp.8040>.
- Herman, G. R., and R. S. Schumacher, 2016a: Extreme precipitation in models: An evaluation. *Wea. Forecasting*, **31**, 1853–1879, <https://doi.org/10.1175/WAF-D-16-0093.1>.
- , and —, 2016b: Using reforecasts to improve forecasting of fog and visibility for aviation. *Wea. Forecasting*, **31**, 467–482, <https://doi.org/10.1175/WAF-D-15-0108.1>.
- , and —, 2018: “Dendrology” in numerical weather prediction: What random forests and logistic regression tell us about forecasting extreme precipitation. *Mon. Wea. Rev.*, <https://doi.org/10.1175/MWR-D-17-0307.1>, in press.
- Hershfield, D. M., 1961: Rainfall frequency atlas of the United States. Weather Bureau, Department of Commerce Tech. Paper 40, 65 pp., http://www.nws.noaa.gov/oh/hdsc/PF_documents/TechnicalPaper_No40.pdf.
- Hitchens, N. M., H. E. Brooks, and R. S. Schumacher, 2013: Spatial and temporal characteristics of heavy hourly rainfall in the United States. *Mon. Wea. Rev.*, **141**, 4564–4575, <https://doi.org/10.1175/MWR-D-12-00297.1>.
- Hong, T., P. Pinson, S. Fan, H. Zareipour, A. Troccoli, and R. J. Hyndman, 2016: Probabilistic energy forecasting: Global Energy Forecasting Competition 2014 and beyond. *Int. J. Forecast.*, **32**, 896–913, <https://doi.org/10.1016/j.ijforecast.2016.02.001>.
- Hou, D., and Coauthors, 2014: Climatology-calibrated precipitation analysis at fine scales: Statistical adjustment of stage IV toward CPC gauge-based analysis. *J. Hydrometeor.*, **15**, 2542–2557, <https://doi.org/10.1175/JHM-D-11-0140.1>.
- Jacks, E., J. B. Bower, V. J. Dagostaro, J. P. Dallavalle, M. C. Erickson, and J. C. Su, 1990: New NGM-based MOS guidance for maximum/minimum temperature, probability of precipitation, cloud amount, and surface wind. *Wea. Forecasting*, **5**, 128–138, [https://doi.org/10.1175/1520-0434\(1990\)005<0128:NNBMGF>2.0.CO;2](https://doi.org/10.1175/1520-0434(1990)005<0128:NNBMGF>2.0.CO;2).
- Junker, N. W., M. J. Brennan, F. Pereira, M. J. Bodner, and R. H. Grumm, 2009: Assessing the potential for rare precipitation events with standardized anomalies and ensemble guidance at the Hydrometeorological Prediction Center. *Bull. Amer. Meteor. Soc.*, **90**, 445–454, <https://doi.org/10.1175/2008BAMS2636.1>.
- Kain, J. S., S. J. Weiss, J. J. Levit, M. E. Baldwin, and D. R. Bright, 2006: Examination of convection-allowing configurations of the WRF Model for the prediction of severe convective weather: The SPC/NSSL Spring Program 2004. *Wea. Forecasting*, **21**, 167–181, <https://doi.org/10.1175/WAF906.1>.
- , and Coauthors, 2008: Some practical considerations regarding horizontal resolution in the first generation of operational convection-allowing NWP. *Wea. Forecasting*, **23**, 931–952, <https://doi.org/10.1175/WAF2007106.1>.
- Klein, W. H., B. M. Lewis, and I. Enger, 1959: Objective prediction of five-day mean temperatures during winter. *J. Meteor.*, **16**, 672–682, [https://doi.org/10.1175/1520-0469\(1959\)016<0672:OPOFDM>2.0.CO;2](https://doi.org/10.1175/1520-0469(1959)016<0672:OPOFDM>2.0.CO;2).
- Lackmann, G. M., 2013: The south-central U.S. flood of May 2010: Present and future. *J. Climate*, **26**, 4688–4709, <https://doi.org/10.1175/JCLI-D-12-00392.1>.
- Lean, H. W., P. A. Clark, M. Dixon, N. M. Roberts, A. Fitch, R. Forbes, and C. Halliwell, 2008: Characteristics of high-resolution versions of the Met Office Unified Model for forecasting convection over the United Kingdom. *Mon. Wea. Rev.*, **136**, 3408–3424, <https://doi.org/10.1175/2008MWR2332.1>.
- Lin, Y., and K. E. Mitchell, 2005: The NCEP Stage II/IV hourly precipitation analyses: Development and applications. *19th Conf. on Hydrology*, San Diego, CA, Amer. Meteor. Soc., 1.2, https://ams.confex.com/ams/Annual2005/techprogram/paper_83847.htm.

- Liu, H., V. Chandrasekar, and G. Xu, 2001: An adaptive neural network scheme for radar rainfall estimation from WSR-88D observations. *J. Appl. Meteor.*, **40**, 2038–2050, [https://doi.org/10.1175/1520-0450\(2001\)040<2038:AANNSF>2.0.CO;2](https://doi.org/10.1175/1520-0450(2001)040<2038:AANNSF>2.0.CO;2).
- McGovern, A., K. L. Elmore, D. J. Gagne, S. E. Haupt, C. D. Karstens, R. Lagerquist, T. Smith, and J. K. Williams, 2017: Using artificial intelligence to improve real-time decision-making for high-impact weather. *Bull. Amer. Meteor. Soc.*, **98**, 2073–2090, <https://doi.org/10.1175/BAMS-D-16-0123.1>.
- Miller, J., R. Frederick, and R. Tracey, 1973: *Precipitation-Frequency Atlas of the Western United States*. NOAA Atlas 2, Vol. 9, U.S. Department of Commerce, NOAA/NWS, 35 pp.
- 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, <https://doi.org/10.1002/qj.49712252905>.
- Murphy, A. H., 1973: A new vector partition of the probability score. *J. Appl. Meteor.*, **12**, 595–600, [https://doi.org/10.1175/1520-0450\(1973\)012<0595:ANVPOT>2.0.CO;2](https://doi.org/10.1175/1520-0450(1973)012<0595:ANVPOT>2.0.CO;2).
- Murphy, K. P., 2012: *Machine Learning: A Probabilistic Perspective*. MIT Press, 1102 pp.
- Nelson, B. R., O. P. Prat, D.-J. Seo, and E. Habib, 2016: Assessment and implications of NCEP Stage IV quantitative precipitation estimates for product intercomparisons. *Wea. Forecasting*, **31**, 371–394, <https://doi.org/10.1175/WAF-D-14-00112.1>.
- Nielsen, E. R., and R. S. Schumacher, 2016: Using convection-allowing ensembles to understand the predictability of an extreme rainfall event. *Mon. Wea. Rev.*, **144**, 3651–3676, <https://doi.org/10.1175/MWR-D-16-0083.1>.
- , G. R. Herman, R. C. Tournay, J. M. Peters, and R. S. Schumacher, 2015: Double impact: When both tornadoes and flash floods threaten the same place at the same time. *Wea. Forecasting*, **30**, 1673–1693, <https://doi.org/10.1175/WAF-D-15-0084.1>.
- North, G. R., T. L. Bell, R. F. Cahalan, and F. J. Moeng, 1982: Sampling errors in the estimation of empirical orthogonal functions. *Mon. Wea. Rev.*, **110**, 699–706, [https://doi.org/10.1175/1520-0493\(1982\)110<0699:SEITEO>2.0.CO;2](https://doi.org/10.1175/1520-0493(1982)110<0699:SEITEO>2.0.CO;2).
- Novak, D. R., C. Bailey, K. F. Brill, P. Burke, W. A. Hogsett, R. Rausch, and M. Schichtel, 2014: Precipitation and temperature forecast performance at the Weather Prediction Center. *Wea. Forecasting*, **29**, 489–504, <https://doi.org/10.1175/WAF-D-13-00066.1>.
- NWS, 2017a: Service change notice 17-100. National Centers for Environmental Prediction, Weather Prediction Center Rep., http://www.nws.noaa.gov/os/notification/scn17-100wpc_excessive_rainfall.htm.
- , 2017b: Summary of natural hazard statistics in the United States. Office of Climate, Weather, and Water Services, National Weather Service, <http://www.nws.noaa.gov/om/hazstats.shtml>.
- Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.*, **12**, 2825–2830, <http://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>.
- Pelosi, A., H. Medina, J. Van den Bergh, S. Vannitsem, and G. B. Chirico, 2017: Adaptive Kalman filtering for postprocessing ensemble numerical weather predictions. *Mon. Wea. Rev.*, **145**, 4837–4854, <https://doi.org/10.1175/MWR-D-17-0084.1>.
- Perica, S., and Coauthors, 2011: *Precipitation-Frequency Atlas of the United States*. NOAA Atlas 14, Vol. 6, 241 pp.
- , and Coauthors, 2013: *Precipitation-Frequency Atlas of the United States*. NOAA Atlas 14, Vol. 9, 171 pp.
- Pinto, J. O., J. A. Grim, and M. Steiner, 2015: Assessment of the High-Resolution Rapid Refresh model's ability to predict mesoscale convective systems using object-based evaluation. *Wea. Forecasting*, **30**, 892–913, <https://doi.org/10.1175/WAF-D-14-00118.1>.
- Quinlan, J. R., 1986: Induction of decision trees. *Mach. Learn.*, **1**, 81–106.
- Ramshaw, J. D., 1985: Conservative rezoning algorithm for generalized two-dimensional meshes. *J. Comput. Phys.*, **59**, 193–199, [https://doi.org/10.1016/0021-9991\(85\)90141-X](https://doi.org/10.1016/0021-9991(85)90141-X).
- Reed, S., J. Schaake, and Z. Zhang, 2007: A distributed hydrologic model and threshold frequency-based method for flash flood forecasting at ungauged locations. *J. Hydrol.*, **337**, 402–420, <https://doi.org/10.1016/j.jhydrol.2007.02.015>.
- Roebber, P. J., 2015: Adaptive evolutionary programming. *Mon. Wea. Rev.*, **143**, 1497–1505, <https://doi.org/10.1175/MWR-D-14-00095.1>.
- Ross, D. A., J. Lim, R.-S. Lin, and M.-H. Yang, 2008: Incremental learning for robust visual tracking. *Int. J. Comput. Vis.*, **77**, 125–141, <https://doi.org/10.1007/s11263-007-0075-7>.
- Rutz, J. J., W. J. Steenburgh, and F. M. Ralph, 2014: Climatological characteristics of atmospheric rivers and their inland penetration over the western United States. *Mon. Wea. Rev.*, **142**, 905–921, <https://doi.org/10.1175/MWR-D-13-00168.1>.
- Scheuerer, M., and T. M. Hamill, 2015: Statistical postprocessing of ensemble precipitation forecasts by fitting censored, shifted gamma distributions. *Mon. Wea. Rev.*, **143**, 4578–4596, <https://doi.org/10.1175/MWR-D-15-0061.1>.
- Schumacher, R. S., and R. H. Johnson, 2006: Characteristics of U.S. extreme rain events during 1999–2003. *Wea. Forecasting*, **21**, 69–85, <https://doi.org/10.1175/WAF900.1>.
- , and —, 2008: Mesoscale processes contributing to extreme rainfall in a midlatitude warm-season flash flood. *Mon. Wea. Rev.*, **136**, 3964–3986, <https://doi.org/10.1175/2008MWR2471.1>.
- , A. J. Clark, M. Xue, and F. Kong, 2013: Factors influencing the development and maintenance of nocturnal heavy-rain-producing convective systems in a storm-scale ensemble. *Mon. Wea. Rev.*, **141**, 2778–2801, <https://doi.org/10.1175/MWR-D-12-00239.1>.
- Shlens, J., 2014: A tutorial on principal component analysis. ArXiv, <https://arxiv.org/abs/1404.1100v1>.
- Stevenson, S. N., and R. S. Schumacher, 2014: A 10-year survey of extreme rainfall events in the central and eastern United States using gridded multisensor precipitation analyses. *Mon. Wea. Rev.*, **142**, 3147–3162, <https://doi.org/10.1175/MWR-D-13-00345.1>.
- Vislocky, R. L., and J. M. Fritsch, 1997: Performance of an advanced MOS system in the 1996–97 National Collegiate Weather Forecasting Contest. *Bull. Amer. Meteor. Soc.*, **78**, 2851–2858, [https://doi.org/10.1175/1520-0477\(1997\)078<2851:POAAMS>2.0.CO;2](https://doi.org/10.1175/1520-0477(1997)078<2851:POAAMS>2.0.CO;2).
- Wang, S.-Y., T.-C. Chen, and S. E. Taylor, 2009: Evaluations of NAM forecasts on midtropospheric perturbation-induced convective storms over the U.S. northern plains. *Wea. Forecasting*, **24**, 1309–1333, <https://doi.org/10.1175/2009WAF2222185.1>.
- Wei, M., Z. Toth, R. Wobus, and Y. Zhu, 2008: Initial perturbations based on the ensemble transform (ET) technique in the NCEP global operational forecast system. *Tellus*, **60A**, 62–79, <https://doi.org/10.1111/j.1600-0870.2007.00273.x>.
- Weisman, M. L., C. Davis, W. Wang, K. W. Manning, and J. B. Klemp, 2008: Experiences with 0–36-h explicit convective forecasts with the WRF-ARW Model. *Wea. Forecasting*, **23**, 407–437, <https://doi.org/10.1175/2007WAF2007005.1>.
- Welles, E., S. Sorooshian, G. Carter, and B. Olsen, 2007: Hydrologic verification: A call for action and collaboration. *Bull. Amer. Meteor. Soc.*, **88**, 503–512, <https://doi.org/10.1175/BAMS-88-4-503>.
- Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. 3rd ed. International Geophysics Series, Vol. 100, Academic Press, 704 pp.
- Wilson, J. W., and R. D. Roberts, 2006: Summary of convective storm initiation and evolution during IHOP: Observational

- and modeling perspective. *Mon. Wea. Rev.*, **134**, 23–47, <https://doi.org/10.1175/MWR3069.1>.
- Wolter, K., J. K. Eischeid, L. Cheng, and M. Hoerling, 2016: What history tells us about 2015 U.S. daily rainfall extremes. *Bull. Amer. Meteor. Soc.*, **97**, S9–S13, <https://doi.org/10.1175/BAMS-D-16-0166.1>.
- Zeng, J., and W. Qiao, 2011: Support vector machine-based short-term wind power forecasting. *Power Systems Conf. and Exposition (PSCE)*, Phoenix, AZ, IEEE, 1–8.
- Zhang, F., C. Snyder, and R. Rotunno, 2003: Effects of moist convection on mesoscale predictability. *J. Atmos. Sci.*, **60**, 1173–1185, [https://doi.org/10.1175/1520-0469\(2003\)060<1173:EOMCOM>2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060<1173:EOMCOM>2.0.CO;2).
- , N. Bei, R. Rotunno, C. Snyder, and C. C. Epifanio, 2007: Mesoscale predictability of moist baroclinic waves: Convection-permitting experiments and multistage error growth dynamics. *J. Atmos. Sci.*, **64**, 3579–3594, <https://doi.org/10.1175/JAS4028.1>.